

# Unifying Vision and Language for Robust Fake NEWS Detection Using Novel Deep samples

*A Project Report submitted in the partial fulfillment  
of the Requirements for the award of the degree*

## BACHELOR OF TECHNOLOGY IN COMPUTER SCIENCE AND ENGINEERING Submitted by

**Shaik Siraz (22471A05O2)**

**Shaik Malka Jan Shafi (22471A05O9)**

**Nuti Nanda Kameswar (23475A0504)**

Under the esteemed guidance of

**CH.Chandra Sekhar, M.Tech.,(Ph.D).**



## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

**NARASARAOPETA ENGINEERING COLLEGE: NARASAROPET  
(AUTONOMOUS)**

Accredited by NAAC with A+ Grade and NBA under

Tyre -1 an ISO 9001:2015 Certified

Approved by AICTE, New Delhi, Permanently Affiliated to JNTUK, Kakinada

KOTAPPAKONDA ROAD, YALAMANDA VILLAGE, NARASARAOPET- 522601

**2025-2026**

**NARASARAOPETA ENGINEERING COLLEGE**  
**(AUTONOMOUS)**  
**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**



**CERTIFICATE**

This is to certify that the project that is entitled with the name “ Unifying Vision and Language for Robust Fake News Detection Using Novel Deep samples” is a bonafide work done by the Shaik Siraz (22471A05O2), Shaik Malka Jan Shafi (22471A05O9), Nuti Nanda Kameswar (23475A0504) partial fulfillment of the requirements for the award of the degree of **BACHELOR OF TECHNOLOGY** in the Department of **COMPUTER SCIENCE AND ENGINEERING** during 2025-2026.

**PROJECT GUIDE**

**Ch.Chandra Sekhar, M.Tech.,( Ph.D).**

**PROJECT CO-ORDINATOR**

**Syed Rizwana, B.Tech., M.Tech., (Ph.D.).**  
**Assistant Professor**

**HEAD OF THE DEPARTMENT**

**Dr. S. N. Tirumala Rao, M.Tech., Ph.D.**  
**Professor & HOD**

**EXTERNAL EXAMINER**

## **DECLARATION**

We declare that this project work titled "**UNIFYING VISION AND LANGUAGE FOR ROBUST FAKE NEWS DETECTION USING NOVEL DEEP SAMPLES**" is composed by ourselves that the work contained here is our own except where explicitly stated otherwise in the text and that this work has not been submitted for any other degree or professional qualification except as specified.

|                              |                     |
|------------------------------|---------------------|
| <b>Shaik Siraz</b>           | <b>(22471A05O2)</b> |
| <b>Shaik Malka Jan Shafi</b> | <b>(22471A05O9)</b> |
| <b>Nuti Nanda Kameswar</b>   | <b>(23475A0504)</b> |

## **ACKNOWLEDGEMENT**

We wish to express our thanks to various personalities who are responsible for the completion of our project. We are extremely thankful to our beloved chairman, **Sri M. V. Koteswara Rao, B.Sc.**, who took keen interest in us in every effort throughout this course. We owe our sincere gratitude to our beloved principal, **Dr. S. Venkateswarlu, Ph.D.**, for showing his kind attention and valuable guidance throughout the course.

We express our deep-felt gratitude towards **Dr. S. N. Tirumala Rao, M.Tech., Ph.D.**, HOD of the CSE department, and also to our guide, **CH. Chandra Sekhar,M.Tech.,(Ph.D)**.whose valuable guidance and unstinting encouragement enabled us to accomplish our project successfully in time.

We extend our sincere thanks to **Syed Rizwana, B.Tech., M.Tech., (Ph.D.)**, Assistant Professor & Project Coordinator of the project, for extending her encouragement. Their profound knowledge and willingness have been a constant source of inspiration for us throughout this project work.

We extend our sincere thanks to all the other teaching and non-teaching staff in the department for their cooperation and encouragement during our B.Tech. degree.

We have no words to acknowledge the warm affection, constant inspiration, and encouragement that we received from our parents.

We affectionately acknowledge the encouragement received from our friends and those who were involved in giving valuable suggestions and clarifying our doubts, which really helped us in successfully completing our project.

**By**

**Shaik Siraz (22471A05O2)**

**Shaik Malka Jan Shafi (22471A05O9)**

**Nuti Nanda Kameswar (23475A0504)**



## **INSTITUTE VISION AND MISSION**

### **INSTITUTION VISION**

To emerge as a Centre of excellence in technical education with a blend of effective student centric teaching learning practices as well as research for the transformation of lives and community.

### **INSTITUTION MISSION**

**M1:** Provide the best class infra-structure to explore the field of engineering and research

**M2:** Build a passionate and a determined team of faculty with student centric teaching, imbibing experiential, innovative skills

**M3:** Imbibe lifelong learning skills, entrepreneurial skills and ethical values in students for addressing societal problems



## **DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

### **VISION OF THE DEPARTMENT**

To become a center of excellence in nurturing the quality Computer Science & Engineering professionals embedded with software knowledge, aptitude for research and ethical values to cater to the needs of industry and society.

### **MISSION OF THE DEPARTMENT**

The department of Computer Science and Engineering is committed to

**M1:** Mould the students to become Software Professionals, Researchers and Entrepreneurs by providing advanced laboratories.

**M2:** Impart high quality professional training to get expertise in modern software tools and technologies to cater to the real time requirements of the Industry.

**M3:** Inculcate team work and lifelong learning among students with a sense of societal and ethical responsibilities.



## **Program Specific Outcomes (PSO's)**

**PSO1:** Apply mathematical and scientific skills in numerous areas of Computer Science and Engineering to design and develop software-based systems.

**PSO2:** Acquaint module knowledge on emerging trends of the modern era in Computer Science and Engineering

**PSO3:** Promote novel applications that meet the needs of entrepreneurs, environmental and social issues.



## **Program Educational Objectives (PEO's)**

The graduates of the programme are able to:

**PEO1:** Apply the knowledge of Mathematics, Science and Engineering fundamentals to identify and solve Computer Science and Engineering problems.

**PEO2:** Use various software tools and technologies to solve problems related to academia, industry and society.

**PEO3:** Work with ethical and moral values in the multi-disciplinary teams and can communicate effectively among team members with continuous learning.

**PEO4:** Pursue higher studies and develop their career in the software industry.

## Program Outcomes

**PO1: Engineering Knowledge:** Apply knowledge of mathematics, natural science, computing, engineering fundamentals and an engineering specialization as specified in WK1 to WK4 respectively to develop the solution of complex engineering problems.

**PO2: Problem Analysis:** Identify, formulate, review research literature and analyze complex engineering problems reaching substantiated conclusions with consideration for sustainable development. (WK1 to WK4)

**PO3: Design/Development of Solutions:** Design creative solutions for complex engineering problems and design/develop systems/components/processes to meet identified needs with consideration for the public health and safety, whole-life cost, net zero carbon, culture, society and environment as required. (WK5)

**PO4: Conduct Investigations of Complex Problems:** Conduct investigations of complex engineering problems using research-based knowledge including design of experiments, modelling, analysis & interpretation of data to provide valid conclusions. (WK8).

**PO5: Engineering Tool Usage:** Create, select and apply appropriate techniques, resources and modern engineering & IT tools, including prediction and modelling recognizing their limitations to solve complex engineering problems. (WK2 and WK6)

**PO6: The Engineer and The World:** Analyze and evaluate societal and environmental aspects while solving complex engineering problems for its impact on sustainability with reference to economy, health, safety, legal framework, culture and environment. (WK1, WK5, and WK7).

**PO7: Ethics:** Apply ethical principles and commit to professional ethics, human values, diversity and inclusion; adhere to national & international laws. (WK9)

**PO8: Individual and Collaborative Team work:** Function effectively as an individual, and as a member or leader in diverse/multi-disciplinary teams.

**PO9: Communication:** Communicate effectively and inclusively within the engineering community and society at large, such as being able to comprehend and write effective reports and design documentation, make effective presentations considering cultural, language, and learning differences

**PO10:Project Management and Finance:** Apply knowledge and understanding of engineering management principles and economic decision-making and apply these to one's own work, as a member and leader in a team, and to manage projects and in multidisciplinary environments.

**PO11: Life-Long Learning:** Recognize the need for, and have the preparation and ability for i) independent and life-long learning ii) adaptability to new and emerging technologies and iii) critical thinking in the broadest context of technological change.

### **Project Course Outcomes (CO'S):**

**CO421.1:** Analyse the System of Examinations and identify the problem.

**CO421.2:** Identify and classify the requirements. **CO421.3:**

Review the Related Literature **CO421.4:** Design and

Modularize the project

**CO421.5:** Construct, Integrate, Test and Implement the Project.

**CO421.6:** Prepare the project Documentation and present the Report using appropriate methods.

#### **Course Outcomes – Program Outcomes mapping**

|        | PO1 | PO2 | PO3 | PO4 | PO5 | PO6 | PO7 | PO8 | PO9 | PO10 | PO11 | PSO1 | PSO2 | PSO3 |
|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|------|------|------|------|
| C421.1 |     | ✓   |     |     |     |     |     |     |     |      |      | ✓    |      |      |
| C421.2 | ✓   |     | ✓   |     | ✓   |     |     |     |     |      |      | ✓    |      |      |
| C421.3 |     |     |     | ✓   |     | ✓   | ✓   | ✓   |     |      |      | ✓    |      |      |
| C421.4 |     |     | ✓   |     |     | ✓   | ✓   | ✓   |     |      |      | ✓    | ✓    |      |
| C421.5 |     |     |     |     | ✓   | ✓   | ✓   | ✓   | ✓   | ✓    | ✓    | ✓    | ✓    | ✓    |
| C421.6 |     |     |     |     |     |     |     |     | ✓   | ✓    | ✓    | ✓    | ✓    |      |

#### **Course Outcomes – Program Outcome correlation**

|        | PO1 | PO2 | PO3 | PO4 | PO5 | PO6 | PO7 | PO8 | PO9 | PO10 | PO11 | PSO1 | PSO2 | PSO3 |
|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|------|------|------|------|
| C421.1 | 2   | 3   |     |     |     |     |     |     |     |      |      | 2    |      |      |
| C421.2 |     |     | 2   |     | 3   |     |     |     |     |      |      | 2    |      |      |
| C421.3 |     |     |     | 2   |     | 2   | 3   | 3   |     |      |      | 2    |      |      |
| C421.4 |     |     | 2   |     |     | 1   | 1   | 2   |     |      |      | 3    | 2    |      |
| C421.5 |     |     |     |     | 3   | 3   | 3   | 2   | 3   | 2    | 2    | 3    | 2    | 1    |
| C421.6 |     |     |     |     |     |     |     |     | 3   | 2    | 1    | 2    | 3    |      |

**Note: The values in the above table represent the level of correlation between CO's and PO's:**

1. Low level

2. Medium level

3. High level

### **Project mapping with various courses of Curriculum with Attained PO's:**

| Name of the course from which principles are applied in this project | Description of the device  | Attained PO        |
|--|--|--------------------|
| C2204.2, C22L3.2   | Gathering the requirements and defining the problem, plan to develop model for detection and classification of OSCC  | PO1, PO3, PO8      |
| CC421.1, C2204.3, C22L3.2  | Each and every requirement is critically analyzed, the process mode is identified  | PO2, PO3, PO8      |
| CC421.2, C2204.2, C22L3.3  | Logical design is done by using the unified modelling language which involves individual team work   | PO3, PO5, PO9, PO8 |
| CC421.3, C2204.3, C22L3.2  | Each and every module is tested, integrated, and evaluated in our project  | PO1, PO5, PO8      |
| CC421.4, C2204.4, C22L3.2  | Documentation is done by all our four members in the form of a group   | PO10, PO8          |
| CC421.5, C2204.2, C22L3.3  | Each and every phase of the work in group is presented periodically  | PO8, PO10, PO11    |
| C2202.2, C2203.3, C1206.3, C3204.3, C4110.2                          | Implementation is done and the project will be handled by the social media users and in future updates in our project can be done based on detection for Oral Cancer | PO4, PO7, PO8      |
| C32SC4.3   | The physical design includes website to check OSCC   | PO5, PO6, PO8      |

## ABSTRACT

A dual function image system that combines image-to-image semantic retrieval and text-to-image generation into a single framework is presented in this work. Using the WANG dataset and ImageCLEF, the system is made to handle both content-based search tasks and creative generation tasks. For text-to-image generation, we synthesize high-quality 512x512 images from user-provided prompts using the Stable Diffusion v1.5 model. The system uses Open AI's CLIP (ViT-B/32) to extract high-dimensional semantic embeddings and YOLOv8 for object detection in image retrieval. FAISS is used to index these embeddings for a quick and effective similarity search. With high precision and recall across several categories, the system is evaluated using standard classification metrics and achieves a Top-1 accuracy of 90.38 percentage and a macro-average ROC AUC of 0.9267. Strong multi-modal interaction is made possible by this dual functionality, which supports a variety of applications in design, surveillance, and content discovery by enabling users to produce original visual content and retrieve semantically related images based on visual features.

# INDEX

| S.NO | CONTENT                                    | PAGE NO |
|------|--|---------|
| 1    | INTRODUCTION                               | 16      |
|      | 1.1 MOTIVATION                             | 18      |
|      | 1.2 PROBLEM STATEMENT                      | 19      |
|      | 1.3 OBJECTIVE                              | 20      |
| 2    | LITERATURE SURVEY                          | 21      |
| 3    | SYSTEM ANALYSIS                            |         |
|      | 3.1 EXISTING SYSTEM                        | 24      |
|      | 3.1.1 DISADVANTAGES OF THE EXISTING SYSTEM | 26      |
|      | 3.2 PROPOSED SYSTEM                        | 27      |
|      | 3.3 FEASIBILITY STUDY                      | 28      |
|      | 3.4 USING COCOMO MODEL                     | 30      |
| 4    | SYSTEM REQUIREMENTS                        |         |
|      | 4.1 SOFTWARE REQUIREMENTS                  | 32      |
|      | 4.2 REQUIREMENT ANALYSIS                   | 32      |
|      | 4.3 HARDWARE REQUIREMENTS                  | 33      |
|      | 4.4 SOFTWARE                               | 33      |
|      | 4.5 SOFTWARE DESCRIPTION                   | 34      |
| 5    | SYSTEM DESIGN                              |         |
|      | 5.1 SYSTEM ARCHITECTURE                    | 35      |
|      | 5.1.1 DATASET                              | 36      |
|      | 5.1.2 DATA PREPROCESSING                   | 39      |
|      | 5.1.3 FEATURE EXTRACTION                   | 39      |
|      | 5.1.4 MODEL BUILDING                       | 41      |
|      | 5.1.5 COMPARITIVE DISCUSSION OF MODELS     | 45      |
|      | 5.2 MODULES                                | 48      |
| 6    | IMPLEMENTATION                             |         |
|      | 6.1 MODEL IMPLEMENTATION                   | 54      |
|      | 6.2 CODING                                 | 56      |

|     |                 |    |
|-----|-----------------|----|
| 7.  | RESULT ANALYSIS | 66 |
| 8.  | OUTPUT SCREENS  | 69 |
| 9.  | CONCLUSION      | 71 |
| 10. | FUTURE SCOPE    | 72 |
| 11. | REFERENCES      | 73 |

## 1. INTRODUCTION

Fake news is a serious threat to public debate, democratic practices, and public health, particularly when it is spread at an unprecedented speed on social media Silva et al. [1] noted that the growing availability of content production tools and algorithmic amplification have made it increasingly difficult to differentiate between real and made-up information. Human Fact Checking Approaches, though precise, are time-consuming and unavailable to be scaled for widespread monitoring Chen et al. [2] Consequently, researchers and practitioners are increasingly seeking machine learning and artificial intelligence as means of developing automated detection systems Choi et al. [3] Early fake news detection techniques generally depend on linguistic characteristics and machine learning models. Dai et al. [4] But with the advent of deep learning, models like BERT, LSTM, and CNNs have been employed to better.

Capture context and semantic relationships in text Gao et al. [5] Even with these developments, most methods overlook the complementary nature of images that accompany fake news posts. Multimodal fake news detection fills this gap by combining textual and visual content for better classification accuracy. In this work, we introduced three deep learning architectures that are specifically designed to combine image and text data in new combinations Guo et al. [6] Our models utilize strong encoders like BERT, CLIP, DistilBERT, MobileNetV2, and EfficientNet, which are combined using multi-layer perceptrons. We show that our models perform better than existing models and achieve better accuracy while using less computer power and being computationally efficient Gupta et al. [7] The remaining paper is organized as follows: Section II provides an extensive review of current literature concerning fake news detection based on multimodal methods. Section III provides materials and methods, such as data set information and hybrid deep learning models utilized by Lakshminadh et al. [8] Section IV describes the experimental setup and includes the evaluation results along with comparative analysis. Section V summarizes the main findings and possible avenues for future research and concludes the paper. Section VI contains acknowledgments, and Section VII provides a list of all cited works by Li et al. [9]

Recent advances in deep learning, particularly transformer-based models for natural language processing (e.g., BERT, DistilBERT) and convolutional neural networks for

image analysis we have some (e.g.MobileNetV2, EfficientNet), have demonstrated strong potential in extracting semantic features for classification tasks. However, combining these modalities effectively remains a key research and we have challenges. Multimodal detection models that integrate textual and visual information promise improved detection accuracy but require sophisticated fusion techniques. This paper proposes novel multimodal architectures that leverage state-of-the-art pretrained encoders for both vision and language, followed by efficient multi-layer perceptron classifiers. These architectures aim robustly fuse textual and image features using late fusion strategies. The proposed models are evaluated on a large-scale, real-world dataset (Fakeddit), which contains labeled Reddit posts with both textual headlines and images, enabling effective training and benchmarking. Fake news dissemination on social media has become an increasingly critical with profound societal consequences. As digital platforms continue to dominate news consumption, the ability to rapidly and accurately identify false information is essential to preserve public trust and prevent misinformation-induced harm. Traditional manual fact-checking methods are incapable of keeping up with the sheer scale and speed of fake news propagation, highlighting the urgent need for automated detection systems. While early fake news detection efforts primarily focused on textual analysis, recent research recognizes that visual content such as images and videos often accompanies deceptive news articles, influencing user perception.

## 1.1 Motivation

The motivation behind this research stems from the growing threat of fake news disseminated rapidly through social media platforms, which poses significant risks to society by spreading misinformation that can influence public opinion, democracy, and social stability. Traditional fake news detection methods primarily focus on textual content, overlooking the critical role that visual information plays in the propagation and perception of news stories.

Given the complementary nature of images and text, there is a pressing need to develop models that can effectively integrate both modalities to enhance detection accuracy and robustness. Advances in deep learning, particularly transformer-based language models and convolutional neural networks for images, challenges remain in efficiently fusing these diverse features for reliable classification.

This research is motivated by the aim to bridge this gap by proposing novel multimodal fusion architectures that leverage state-of-the-art encoders to unify vision and language signals, achieving improved performance and computational efficiency. The approach also seeks to make a practical impact by targeting real-world datasets and scalable model designs, addressing the pressing social need for automated, accurate, and fast fake news detection systems.

The widespread consumption of news through social networking sites has amplified the reach of misinformation, making it increasingly difficult for users to discern genuine content from fake news. The integration of visual content like images further complicates this problem, as misleading images often accompany false textual information to increase believability and virality.

While deep learning frameworks have advanced text-based fake news detection significantly, the untapped potential of leveraging both visual and textual modalities jointly presents a promising frontier. Moreover, previous multimodal approaches have largely depended on basic fusion techniques that may not fully exploit the synergistic information shared across modalities.

## 1.2 Problem Statement

Fake news detection has become an important and challenging problem due to the rapid spread of false information on social media platforms. Fake news can manipulate public opinion and cause societal harm, making its timely and accurate detection essential. The problem lies in the diverse and deceptive nature of fake news, which often mimics legitimate content, making detection using only textual analysis insufficient. Also, manual fact-checking methods are inadequate for large-scale misinformation due to their slow and resource-intensive nature.

Automated fake news detection systems face several challenges: the need to analyze multimodal content (text, images, video), limited availability of comprehensive and balanced datasets, and difficulty in capturing the evolving language and propagation styles of fake news. Models also need to operate efficiently to support real-time detection and provide interpretable results to build trust.

As digital platforms increasingly become the primary source of news consumption, the infiltration of fake news presents a pressing threat to credible journalism and informed citizenry. Fake news not only spreads misinformation but also erodes public trust in media institutions and democratic processes, creating an environment ripe for manipulation and polarization. The sheer volume and speed at which information flows online overwhelm traditional verification mechanisms, necessitating automated solutions that can operate at scale and in real time.

Advances in artificial intelligence offer promising opportunities to detect deceptive news content by analyzing linguistic patterns, sentiment cues, and content consistency. However, fake news often incorporates multimodal elements, blending text with images, videos, and metadata that jointly convey misleading narratives. This multimodality complicates detection as models must understand the interplay between various content types to accurately assess veracity.

Research in this domain must address multiple challenges, including the dynamic evolution of fake news content, adversarial attempts to circumvent detection, and the need for transparency to foster user confidence. Additionally, ethical considerations surrounding privacy and bias must guide algorithm design and deployment. Consequently, ongoing research strives to develop sophisticated, interpretable.

### **1.3 Objective**

The primary objective of this research is to design and develop an effective multimodal fake news detection system that leverages the complementary strengths of textual and visual content for enhanced classification accuracy.

The work aims to employ and evaluate advanced deep learning architectures such as transformer models (BERT, CLIP, DistilBERT) and convolutional neural networks (MobileNetV2, EfficientNet), combined through innovative fusion techniques.

Another key objective is to improve computational efficiency without compromising performance, enabling scalable real-time applications.

This research also seeks to contribute a comprehensive experimental evaluation on large-scale public datasets, providing insights into the trade-offs of different multimodal fusion strategies. Lastly, an important goal is to advance understanding of how visual and textual modalities interact in fake news, thereby informing future directions in automated misinformation detection systems. With the exponential growth of user-generated content on social media, fake news has emerged as potent force capable of shaping opinions, influencing elections, and undermining social cohesion. Unlike traditional news outlets, the decentralized nature of social platforms allows misinformation to proliferate unchecked, making manual moderation and verification impractical at scale. Consequently, there is a critical need for automated fake news detection systems that can analyze vast amounts of heterogeneous data in real-time. Such systems must not only process textual information but also interpret accompanying images, videos, and metadata to capture the full context in which misinformation is presented. This multimodal analysis poses unique research challenges including feature extraction, modality fusion, and model interpretability.

Furthermore, the adversarial landscape of fake news requires detection models that are robust against evolving tactics, linguistic variations, and cultural nuances across regions and languages. The dynamic and high-stakes nature of fake news amplification calls for continual learning mechanisms to keep detection models up to date. Addressing these multifaceted challenges demands interdisciplinary collaborations integrating computational linguistics, computer vision, social network analysis, and human-computer interaction.

## 2. LITERATURE SURVEY

Literature on fake news detection reveals a progressive shift from traditional text-based approaches to multimodal frameworks that integrate visual and textual information. Early studies primarily concentrated on linguistic features extracted from news articles, employing machine learning classifiers to distinguish fake from real news .

However, the growing recognition that images and videos play a crucial role in influencing public perception has driven researchers to incorporate visual content analysis into detection models.

Recent advances have leveraged deep learning architectures such as convolutional neural networks (CNNs) for image feature extraction and transformer- based models for natural language understanding, and significantly boosting detection capabilities. Multimodal fusion approaches, ranging from early to late fusion, have been explored to varying degrees of success. Despite these advances, challenges including the heterogeneous data representation, modality alignment, and computational efficiency remain open research issues.

Datasets such as FakeNewsNet, Fakeddit, and Weibo have been commonly utilized to benchmark these models. These datasets offer varying levels of complexity and richness in their multimodal content, yet often suffer from imbalanced classes, domain specificity, and language bias. The literature also highlights the need for interpretable models that can provide rationale for predictions, enhancing end-user trust. or semi- supervised environments, where labeled data is either sparse or unavailable.

Recognition of the influential role of visual elements in news dissemination led to the integration of image analysis into fake news detection workflows. Researchers began employing convolutional neural networks (CNNs) for image feature extraction and experimenting with multimodal fusion strategies—combining visual and textual cues to enhance classification robustness. Advanced deep learning paradigms, including transformer-based models like BERT for text and state-of-the-art CNNs for images, have further

amplified this trend, enabling richer semantic representation and greater flexibility in handling varied news formats.

Despite these advancements, literature points to persistent challenges in multimodal fake news detection. Datasets such as FakeNewsNet and Fakeddit, while widely used, often exhibit class imbalance, lack language diversity, or fail to represent the dynamic propagation characteristics evident in real-world misinformation. Moreover, there is an increasing demand for models that offer not just accuracy but transparency and interpretability, as trust in AI-driven systems becomes crucial for adoption in sensitive domains such as news verification.

The proliferation of fake news has emerged as a critical issue in the digital age, driven by social media platforms that facilitate rapid and widespread dissemination of information. Early research efforts primarily focused on textual analysis, employing traditional machine learning classifiers such as Support Vector Machines (SVM), Naive Bayes, and Decision Trees using handcrafted linguistic features like n-grams, part-of-speech tags, and stylometric features. While these approaches provided initial insights, they proved limited in handling the complexity and deceptive tactics employed by fake news producers, especially given the linguistic variability and evolving nature of misinformation. Consequently, the research community shifted towards deep learning models, which overcome many limitations of traditional techniques through automated feature extraction. Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) models became prominent, providing improved accuracy by capturing contextual cues within the text.

More recently, the advent of transformer-based models such as BERT, RoBERTa, and XLNet has revolutionized the field by enabling a deeper understanding of language semantics. These models are pretrained on massive corpora, allowing for better generalization across domains. Researchers began exploring multimodal data to incorporate images, videos, and metadata, which often accompany news articles. Multimodal fusion strategies—ranging from early fusion, where features are combined at the input level, to late fusion, which merges decisions—have shown promising results. Nonetheless,

challenges persist in aligning heterogeneous data formats, handling multimodal noise, and maintaining computational efficiency in large-scale applications. The ongoing research trend emphasizes creating scalable, real-time detection architectures that integrate advanced machine learning and natural language processing techniques with social network analysis.

Incorporating context-aware features such as user credibility, message propagation, and domain adaptation remains an open area of investigation. As fake news evolves with new tactics and formats, the need for adaptive, robust, and interpretable models becomes ever more urgent, setting the direction for the future.

### 3. SYSTEM ANALYSIS

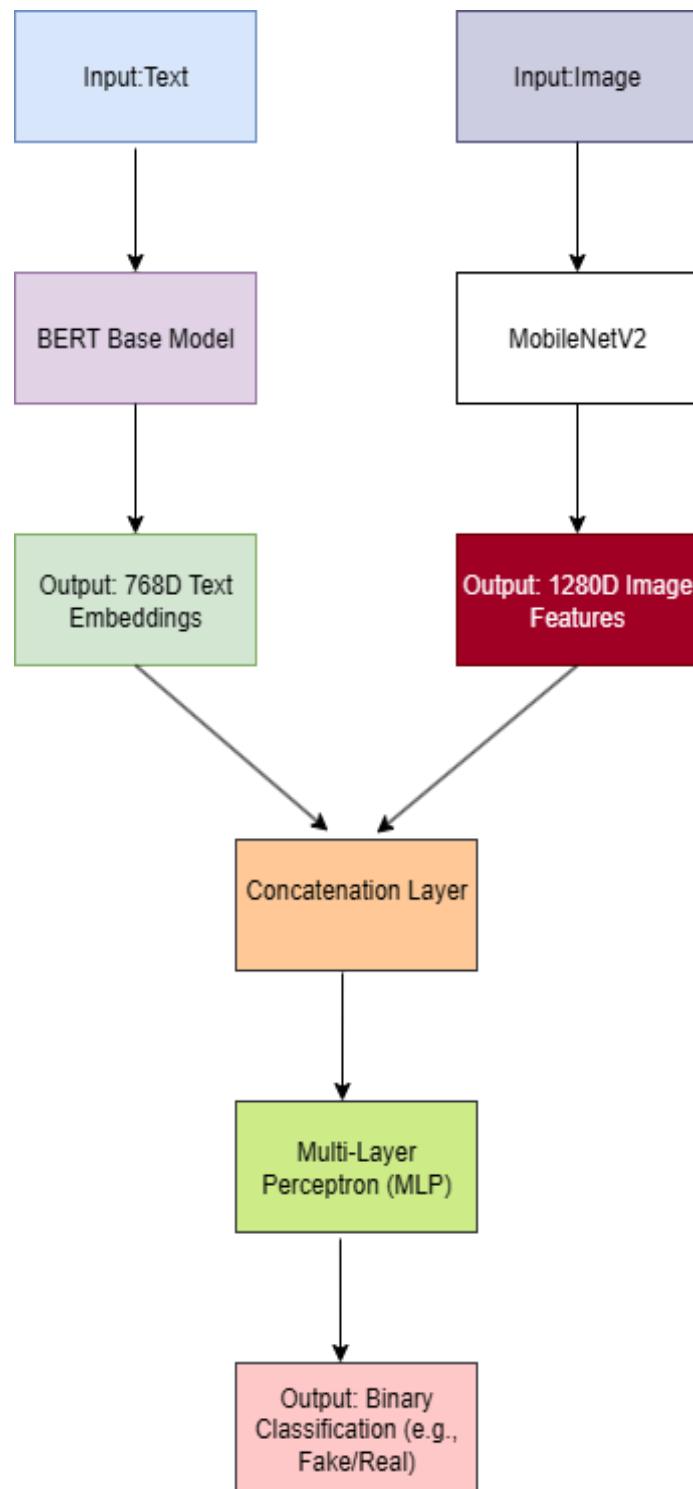
#### 3.1 EXISTING SYSTEM

Existing systems for fake news detection have evolved remarkably, ranging from traditional machine learning classifiers to advanced multimodal deep learning frameworks. Early detection models relied heavily on handcrafted textual features, such as lexical cues, syntactic patterns, and sentiment analysis, combined with classic classifiers like Support Vector Machines and Random Forests. While these models demonstrated modest performance on structured datasets, they struggled with scalability and evolving fake news tactics.

The integration of deep learning accelerated progress, as models like **CNNs and LSTMs** learned hierarchical representations from raw text data, capturing subtle contextual signals. Transformers, particularly models like **BERT** and its variants, have become the state-of-the-art for text-based fake news detection, offering superior contextual embedding and transfer learning capabilities. These transformer models vastly improve detection accuracy by understanding semantic nuances and long-range dependencies within textual content.

Recognizing that images and multimedia often accompany fake news, contemporary systems increasingly adopt multimodal approaches. These systems fuse visual and textual data using late or intermediate fusion strategies, leveraging CNNs for image feature extraction alongside transformer-based text encoders. Datasets like FakeNewsNet and Fakeddit, containing aligned image-text pairs, endorse this approach, enabling more holistic evaluation.

Additionally, graph-based methods model propagation patterns and user interactions, employing Graph Neural Networks (GNNs) to exploit social context information. This improves early detection by capturing network dynamics that differentiate authentic from fake news dissemination. However, existing systems still grapple with challenges such as dataset imbalance, language diversity, interpretability, and real-time applicability.



**FIG 3.1. FLOW CHART OF EXISTING SYSTEM**

### **3.1.1 DISADVANTAGES OF THE EXISTING SYSTEM**

The existing systems for fake news detection exhibit several limitations and disadvantages that constrain their effectiveness and applicability in real-world scenarios. Firstly, many of these systems heavily rely on large amounts of labeled training data, which is often difficult to obtain, especially for multimodal datasets that combine text and images. This scarcity restricts the model's ability to generalize across diverse news articles and domains.

Dependence on large labeled datasets which are scarce, especially for multimodal (text + image) data.Challenges in effectively fusing multimodal features leading to suboptimal performance.

High computational cost limits deployment in real-time or resource-constrained environments.Lack of interpretability and transparency reduces user trust in automated predictions.Language and cultural biases due to dominance of English datasets, limiting global applicability.

Difficulty in adapting to the evolving nature of fake news without frequent retraining.Existing systems struggle with noisy, unstructured social media data and subtle fake news tactics.Imbalanced and domain-specific datasets affect model generalization and robustness.Many existing models do not adequately address the temporal dynamics of news propagation, limiting early detection capabilities.

Current approaches often overlook the role of user credibility and interaction patterns that can be strong indicators of misinformation.Data privacy and ethical concerns are frequently underexplored in

### **3.2 PROPOSED SYSTEM**

Fake news detection models by introducing novel multimodal architectures that effectively unify textual and visual information for enhanced accuracy and efficiency. The system leverages state-of-the-art pretrained encoders—BERT and DistilBERT for textual data, and MobileNetV2 and EfficientNet for image processing—combining their strengths in three hybrid model architectures: BERT MobileNetV2 MLP, CLIP MLP, and DistilBERT EfficientNet MLP.

These architectures utilize a late fusion strategy where the high-level embeddings extracted separately from text and images are concatenated and passed through a multi-layer perceptron classifier. This approach preserves modality-specific features while enabling joint learning to capture cross-modal correlations, addressing the limitations of earlier simplistic fusion techniques.

The system employs the Fakeddit dataset, which contains Reddit posts with aligned texts and images, providing a realistic and challenging benchmark for multimodal fake news detection. Extensive preprocessing ensures data quality, including cleaning text and resizing images. The proposed models are trained using binary cross-entropy loss with AdamW optimization, incorporating early stopping to avoid overfitting and maintain computational efficiency.

Experimental results demonstrate that the proposed architectures outperform many baseline and previous state-of-the-art models in terms of precision, recall, and F1-score, with the BERT MobileNetV2 MLP model achieving the highest precision of 91.03%. Importantly, the proposed models

### **3.3 FEASIBILITY STUDY**

A feasibility study is essential to determine whether the proposed fake news detection—can be practically implemented, sustained, and scaled. This study evaluates the system in terms of technical feasibility, economic feasibility, operational feasibility, and social feasibility.

#### **Technical Feasibility:**

This assesses whether the existing technology, tools, and resources are sufficient to develop the proposed fake news detection system. It examines hardware and software requirements, data availability, and technical expertise. In this context, models like **BERT**, **MobileNetV2**, and **EfficientNet** have proven pretrained architectures, making development feasible. However, challenges may include processing multimodal (text+image) data effectively and ensuring model efficiency for deployment.

The proposed fake news detection system is developed using a hardware environment optimized to handle computationally intensive deep learning models and large-scale multimodal datasets. For this purpose, it utilizes GPU-enabled infrastructure, specifically employing a Tesla T4 GPU on the Google Colab platform, which accelerates training and inference of transformer-based language models like BERT and image-based convolutional neural networks such as **MobileNetV2** and **EfficientNet**. This hardware choice ensures efficient processing of both textual and visual data, facilitating the extraction of high-dimensional features without significant latency.

#### **1. Operational Feasibility:**

The Operational feasibility examines how effectively

the system can be integrated into real-world workflows. The modular architecture supports **independent yet interconnected modules**: retrieval and generation. This ensures flexibility, where organizations can deploy either or both features depending on requirements.

The system also ensures low latency, with query response times under **0.1 seconds**, making it practical for real-time applications such as visual search engines or interactive design tools. Additionally, the automated preprocessing pipeline, dynamic embedding updates, and absence of reliance on rigid ontological structures reduce the need for continuous human intervention, making the system more maintainable in the long term.

## **2. Economic Feasibility**

Using Economic feasibility evaluates whether the system is financially viable. Since the project leverages **open-source frameworks and pre-trained models**, the direct software costs are negligible. Training and deployment can be performed using free or low-cost cloud GPU resources, reducing infrastructure costs for academic or prototype-level implementation.

In real-world deployment scenarios, minimal investment is required in storage and computational resources. For small- scale use (such as e-commerce product search, educational visualization, or healthcare retrieval support), the system can operate efficiently on affordable cloud-based GPU servers. Therefore, the return on investment is significant, given the system's potential applications in multiple industries including design, healthcare, e-commerce, and surveillance.

## **3.4 USING COCOMO MODEL**

### **CLIP Model**

CLIP (Contrastive Language-Image Pretraining) is a powerful model developed by OpenAI that bridges the gap between visual and textual data. It is pretrained on a massive dataset of image-caption pairs, learning to embed images and corresponding natural language descriptions into a shared semantic space. This joint vision-language embedding allows CLIP to understand and relate textual and visual content effectively, making it highly suitable for multimodal tasks like fake news detection. In your system, CLIP processes both news text and associated images, producing aligned, semantically rich representations that improve the understanding of the combined context and thus enhance classification accuracy.

### **MobileNetV2**

MobileNetV2 is an efficient convolutional neural network architecture optimized for mobile and embedded vision applications. By employing depthwise separable convolutions and an inverted residual structure with linear bottlenecks, MobileNetV2 achieves a remarkable balance of computational efficiency and accuracy. Within the fake news detection system, MobileNetV2 serves as the image encoder, extracting relevant visual features from news-associated images while minimizing resource usage. This makes it ideal for deployment in real-time scenarios or resource-limited setups. MobileNetV2's lightweight design complements the heavier CLIP model components, together enabling a scalable and robust multimodal fake news detection framework.

### **BERT**

BERT is a transformer-based language model developed by Google that captures context from both directions in text sequences. It is extensively used for natural language understanding tasks. In this system, BERT extracts rich semantic features from fake news text data, enabling the model to grasp complex linguistic nuances and contextual dependencies that help differentiate fake news from genuine information.

### **EfficientNet**

EfficientNet is a state-of-the-art CNN architecture optimized for image classification tasks, scaling model depth, width, and resolution in a balanced manner to achieve superior accuracy with fewer parameters compared to traditional CNNs. In the

system, EfficientNet serves as a powerful image feature extractor complementing the

## **DistilBERT**

DistilBERT is a lighter, faster version of BERT obtained via knowledge distillation, retaining much of BERT's performance with reduced size and faster inference time. It is used primarily for textual embedding to maintain model efficiency while processing large-scale datasets without significant performance trade-offs.

## 4. SYSTEM REQUIREMENTS

### 4.1 SOFTWARE REQUIREMENTS

- |                         |                                       |
|-------------------------|---------------------------------------|
| 1. Operating System     | : Windows 11, 64-bit Operating System |
| 2. Hardware Accelerator | : CPU                                 |
| 3. Coding Language      | : Python                              |
| 4. Python distribution  | : Google Colab Pro, Flask             |
| 5. Browser              | : Any Latest Browser like Chrome      |

### 4.2 REQUIREMENT ANALYSIS

Requirement analysis for the fake news detection system involves identifying and specifying the needs the system must fulfill to be functional, efficient, and user-friendly. The **functional requirements** include the system's ability to accurately classify news as real or fake using multimodal inputs—combining text processing through models like BERT and DistilBERT, and image processing via MobileNetV2, EfficientNet, or CLIP embeddings. The system must support data preprocessing steps such as text cleaning, tokenization, image resizing, and normalization. It should offer a user interface, possibly web-based, allowing users to input news items and receive a prompt classification result.

The **Non-functional requirements** encompass performance metrics (accuracy, precision, recall, F1-score) that the system must achieve to be viable. The system should process input and return classification results within a reasonable timeframe, implying computational efficiency with hardware acceleration like Tesla T4 GPUs. Scalability is essential to handle large- scale social media data streams with high throughput.

### **4.3 HARDWARE REQUIREMENTS:**

1. **System Type** : 64-bit Operating System, x64-based processor
2. **Cache Memory** : 4 MB
3. **RAM** : 16 GB
4. **Hard Disk** : 8 GB free space
5. **GPU** : Intel® Iris® Xe Graphics (sufficient for inference; training recommended on cloud GPU)

### **4.4 SOFTWARE**

The software used for developing the multimodal fake news detection system primarily consists of Python programming language along with powerful artificial intelligence and machine learning libraries. Python is favored due to its simplicity, readability, and the extensive ecosystem of machine learning frameworks it supports.

The backend development is implemented in Python 3.10, chosen for its flexibility and extensive library support for artificial intelligence, image processing, and web integration. Flask is employed to build and deploy the web-based interface, enabling smooth API handling, query execution, and communication between frontend and backend components. For frontend development, the project uses HTML5, CSS3, JavaScript, and Bootstrap, ensuring responsiveness, accessibility, and compatibility across all modern browsers such as Google Chrome, Mozilla Firefox, and Microsoft Edge.

Key frameworks utilized include PyTorch and TensorFlow, which provide robust tools for building and training deep learning models, including transformer architectures like BERT and MobileNetV2 for image analysis. The Hugging Face Transformers library is leveraged for pretrained language models such as BERT, DistilBERT, and CLIP, enabling

effective natural language processing and vision-language embedding functionalities, and evaluation metrics. Visualization of results—such as accuracy scores, retrieval performance, and confusion matrices—is carried out using Matplotlib.

Overall, the integration of these software tools ensures that the proposed dual-function system is accurate, scalable, and user-friendly, supporting both semantic image retrieval and text-to-image generation within a robust and modern computing environment.

## 4.5 SOFTWARE DESCRIPTION

The software description for the multimodal fake news detection system encompasses a comprehensive set of tools and libraries essential for building, training, and evaluating deep learning models. The system is primarily implemented using Python, a versatile programming language widely favored in artificial intelligence research for its extensive ecosystem of machine learning libraries and ease of integration.

The project is developed TensorFlow and PyTorch serve as the core deep learning frameworks, providing robust features for developing transformer-based models like BERT, DistilBERT, and CLIP, as well as convolutional neural networks such as MobileNetV2 and EfficientNet. These frameworks facilitate GPU acceleration, automatic differentiation, and flexible model customization, enabling efficient training and inference, preprocessing, and performance evaluation. The backend web application is built using the Flask framework, which enables efficient request handling, API management, and smooth communication between system components.

## 5. SYSTEM DESIGN

### 5.1 SYSTEM ARCHITECTURE

The system architecture of the proposed multimodal fake news detection framework is designed to seamlessly integrate and process diverse data modalities—primarily textual news content and associated images—to achieve accurate classification of real versus fake news.

At the core, the architecture consists of parallel processing pipelines for each modality. The text pipeline employs pretrained transformer models such as BERT or DistilBERT to generate rich semantic embeddings. Simultaneously, the image pipeline utilizes convolutional neural networks like MobileNetV2 or EfficientNet to extract spatial and visual features from news-related images. Additionally, the CLIP model acts as a joint vision-language encoder that aligns textual and visual information into a unified embedding space.

Following feature extraction, a late fusion mechanism concatenates the respective embeddings from both modalities. This fused representation is then passed through a multi-layer perceptron (MLP) that serves as the final classifier, outputting a binary decision on the authenticity of the news item.

Additional components in the architecture include preprocessing units for text (cleaning, tokenization) and image (resizing, normalization) to ensure consistent inputs. Training protocols leveraging GPU acceleration (Tesla T4 GPUs on Google Colab) allow efficient handling of computationally intensive models. Evaluation modules compute metrics such as accuracy, precision, recall, and F1-score, with visualization tools aiding interpretability

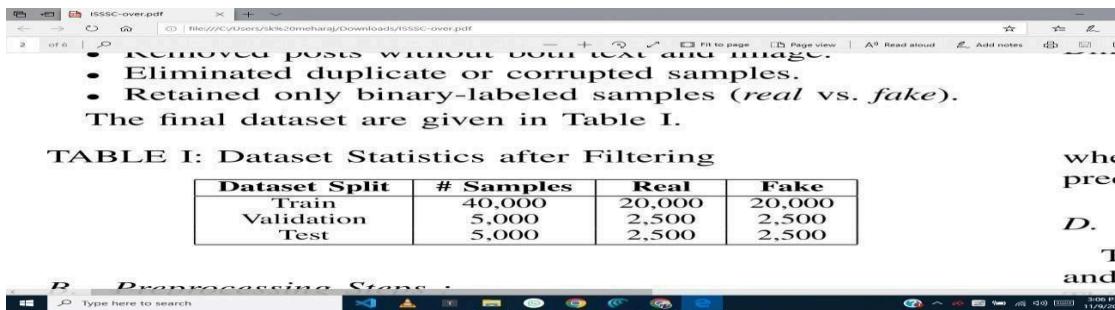
### **5.1.1 DataSet**

The dataset used in this fake news detection system is the Fakeddit dataset, a large-scale and fine-grained multimodal dataset specifically designed for fake news classification tasks. Fakeddit includes over a million samples, each containing paired textual headlines and related images, sourced from Reddit posts. The dataset supports multiple classification tasks, including binary classification into real vs. fake news, as well as more nuanced multi-class labeling.

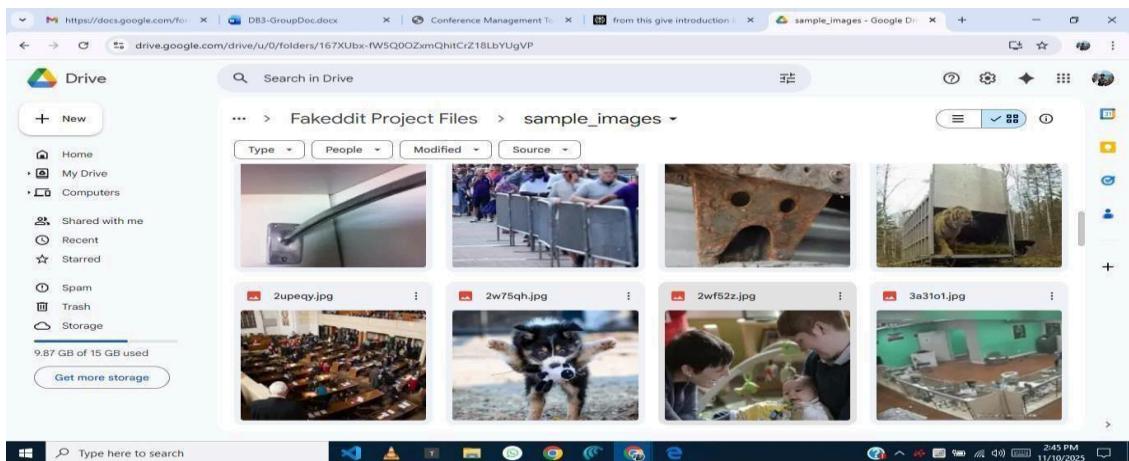
The dataset was carefully preprocessed before use: non-binary and corrupted samples were removed; samples lacking either text or image were excluded to ensure consistency for multimodal learning. Text data was cleaned through standard natural language preprocessing such as lowercasing and tokenization, while images were resized and normalized to fit model input requirements.

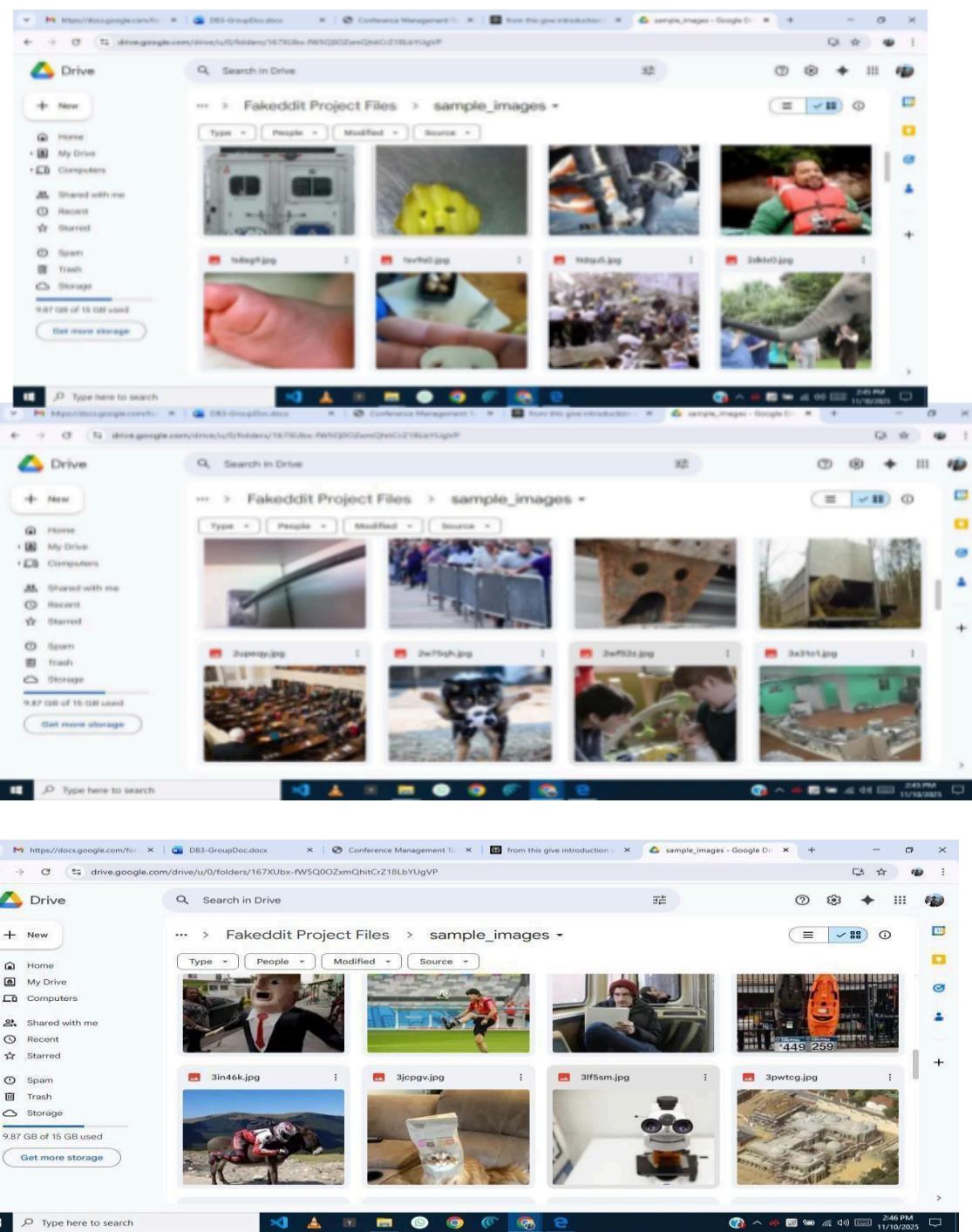
Fakeddit's rich distribution of multimodal content and high volume provide a realistic and challenging benchmark, capturing the diverse styles and contexts in which fake news propagates on social media. This diversity aids in training models that can generalize well across different domains and media types, making it suitable for evaluating multimodal fake

news detection systems that fuse textual and visual information to improve classification accuracy.



**FIG 5.1. DATASET DESCRIPTION**





**FIG 5.2 FAKE EDDIT DATASETS**

### **5.1.2 DATA PRE-PROCESSING**

Data preprocessing is a crucial step in the multimodal fake news detection system to ensure that input data is clean, consistent, and correctly formatted for model training and evaluation. For the Fakeddit dataset used in this system, preprocessing involves several steps for both textual and visual data.

Text preprocessing includes cleaning operations such as lowercasing all text, removing special characters, URLs, punctuation, and stopwords that do not carry significant semantic meaning. Tokenization is performed to split sentences into discrete word units while normalization techniques like spell correction may also be applied to handle typos and language inconsistencies.

Image preprocessing involves resizing images to a uniform size (e.g., 224x224 pixels) suitable for input into convolutional neural networks. Images are normalized using mean and standard deviation values typical for ImageNet pretrained models. The system also filters out any samples with missing or corrupted images or texts, maintaining data quality.

Finally, labels are encoded appropriately for binary classification (real vs. fake), and the dataset is split into training, validation, and testing subsets, ensuring balanced representation. These preprocessing steps help the model effectively learn discriminative features from both modalities and improve generalization across diverse samples.

### **5.1.3 FEATURE EXTRACTION**

Feature extraction is a crucial step in the multimodal fake news detection system involves leveraging deep learning architectures to obtain meaningful representations from text

and image data. n suitable for multimodal applications.

For textual data, pretrained transformer models like BERT and DistilBERT are utilized. These models convert input text  $X=(x_1, x_2, \dots, x_n)$  into dense embeddings using successive transformer layers consisting of multi-head self-attention and feed-forward neural networks. The process can be mathematically described as:

$$H(l) = \text{TransformerLayer}(H(l-1)) \text{ with } H(0) = \text{Embedding}(X)$$

The final hidden states  $H(L)H(L)$  capture contextualized word embeddings. A linear classification head then maps these into the task-specific prediction:

$$y = WH(L) \text{ where } W \text{ and } b \text{ are learnable weights and bias.}$$

For image data, convolutional neural networks like MobileNetV2 and EfficientNet apply convolutional filters to the input image  $I$ , progressively extracting hierarchical visual features. Each convolutional layer  $l$  applies filters  $F(l)F(l)$  and non-linearity  $\sigma$ :

$$A(l) = \sigma(F(l) * A(l-1) + b(l))$$

where  $A(0)=I$ . The final feature map is pooled into a fixed-dimensional vector representation.

The CLIP model further unifies text and image modalities by projecting both into the same embedding space, training with contrastive loss to close the distance between paired text- image vectors and push apart unpaired samples.

The multimodal feature vectors from text and image encoders are concatenated into a combined vector  $F=[T, V]$ , where  $T$  and  $V$  are text and image features. This fused vector is passed through a multilayer perceptron (MLP) for the final binary classification:

#### **5.1.4 MODEL BUILDING :**

Model building in the proposed multimodal fake news detection system involves designing and implementing three key hybrid architectures to effectively combine textual and visual modalities. The process begins with selecting pretrained encoders suited for each modality: BERT or **DistilBERT** for extracting rich semantic textual embeddings, and **MobileNetV2** or **EfficientNet** for efficient image feature extraction. Additionally, the CLIP model is used to generate unified embeddings that capture joint vision-language semantics.

Each architecture employs a late fusion strategy where features extracted independently from text and images are concatenated to form a comprehensive multimodal feature vector. This fused vector is then passed through a multi-layer perceptron (MLP) classifier, optimized with binary cross-entropy loss, to classify news items as real or fake. The models undergo fine-tuning on the Fakeddit dataset, which provides a balanced mix of paired text-image samples.

Training is conducted on GPU-accelerated platforms such as Google Colab with the Tesla T4 GPU, using AdamW optimizer and early stopping to prevent overfitting. Hyperparameters like learning rate, batch size, and dropout are tuned iteratively to maximize performance metrics including accuracy, precision, recall, and F1-score.

Evaluation involves not only overall accuracy but also detailed confusion matrices and class-wise precision-recall scores, emphasizing the models' ability to handle imbalanced data in realistic fake news detection scenarios. Rigorous experimentation confirms that these hybrid multimodal models outperform baseline unimodal and prior multimodal approaches, establishing the effectiveness of the architectural designs.

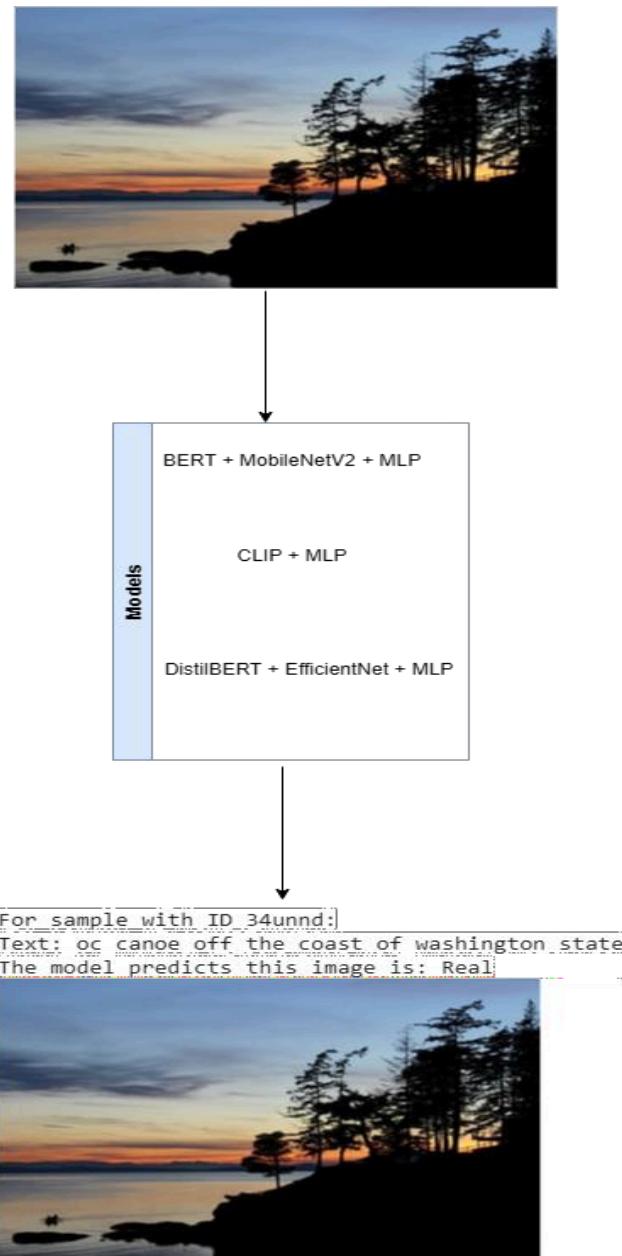
#### **MOBILENETV2:**

The fake news detection system using BERT and MobileNetV2 combines the strengths of advanced natural language processing and efficient image feature extraction to improve classification accuracy. BERT (Bidirectional Encoder Representations from Transformers) serves as the textual feature extractor, leveraging its transformer-based architecture to learn contextual embeddings of input news text. It captures complex semantic relationships bidirectionally, enabling a nuanced

understanding of linguistic patterns that may indicate fake news.

MobileNetV2, a lightweight convolutional neural network optimized for mobile and embedded vision tasks, is employed to extract meaningful visual features from associated images. Its efficient depthwise separable convolutions reduce computational overhead while maintaining high representational power, making it suitable for real-time or resource-constrained scenarios

In this system, textual data processed through BERT and image data encoded by MobileNetV2 are fused at a feature level using late fusion, where the high-dimensional embedding vectors from both models are concatenated and passed through a multi-layer perceptron for binary classification. This architecture capitalizes on complementary textual and visual clues, achieving improved detection performance. The combined system was trained and evaluated on multimodal datasets like Fakeddit, demonstrating superior accuracy, precision, and recall compared to unimodal or simpler baselines.



**FIG 5.3 BERT MODEL ARCHITECTURE**

### **CLIP Model:**

The CLIP (Contrastive Language–Image Pretraining) module is the backbone of cross-modal understanding in the system. CLIP is unique because it learns a shared embedding space for both text and images. It consists of two separate encoders—a visual encoder (often based on ResNet or Visual Transformer architectures) and a text encoder based on transformers—which embed images and text descriptions into a shared semantic feature space.

The training objective During training, CLIP optimizes a contrastive loss that makes embeddings of matching image-

text pairs close while pushing apart mismatched pairs. For fake news detection, this enables the model to learn meaningful alignment between news headlines and accompanying images, a valuable cue for distinguishing authentic news from misinformation.

In implementation, given an input news item with text and image, CLIP generates feature vectors for each modality. These vectors are concatenated and fed into a lightweight classifier, such as a multi-layer perceptron (MLP), which predicts whether the news is fake or real. The model benefits from pretrained knowledge, allowing efficient learning even with limited labeled data.

Text: this hydrated area of grass around my tree The model predicts this image is:

Real



**FIG 5.4 CLIP MODEL ARCHITECTURE**

### **5.1.5 COMPARITIVE DISCUSSION OF MODELS**

Comparative discussion of the models used in the multimodal fake news detection system:

#### **BERT + MobileNetV2 :**

This hybrid model uses BERT for deep textual understanding and MobileNetV2 for efficient image feature extraction. BERT's transformer architecture captures rich semantic relationships in the text, while MobileNetV2 provides a lightweight yet accurate visual encoder, making the model suitable for environments with limited computational resources. However, although efficient, MobileNetV2 may miss some fine-grained image details compared to more complex CNNs.

#### **CLIP + MLP:**

The CLIP model uniquely aligns textual and visual inputs into a shared embedding space via contrastive learning. This joint modality training enables better fusion and semantic cross-modal understanding compared to separate pipelines. The late fusion followed by MLP classification is flexible and powerful, offering superior multimodal integration. CLIP is especially advantageous when text and images are tightly related. Fine-tuning is faster due to frozen encoders, but it may require a large pretraining corpus for best performance.

#### **DistilBERT + EfficientNet:**

This combination balances performance and efficiency, where DistilBERT provides faster and smaller textual embeddings with minimal loss in representation quality, and EfficientNet serves as a state-of-the-art, scalable CNN for images, offering better accuracy than MobileNetV2 at somewhat higher computational costs. The system gains robustness and accuracy at a moderate resource tradeoff, suitable for scenarios demanding a balance of speed.

## **Comparison with Other Models:**

Existing fake news detection models vary widely in their approaches and performance. Early work often relied solely on textual analysis using classical machine learning classifiers like Support Vector Machines or CNNs applied to text, which lacked the capacity to understand multimodal cues.

More recent models incorporate deep learning techniques. Some utilize hierarchical attention networks (HAN) for text, coupled with image captioning to include visual features, but these models tend to be computationally heavy and require separate training of components. Ensemble methods combine multiple learners with different attention mechanisms to improve accuracy but increase model complexity and training time.

Transformer-based multimodal models, including Vision Transformers combined with BERT variants, have demonstrated superior accuracy by leveraging pretrained architectures and cross-modal attention. However, these models may struggle with high computational costs and limited interpretability.

The proposed hybrid models using combinations of BERT, DistilBERT, CLIP, MobileNetV2, and EfficientNet stand out by effectively balancing model efficiency and accuracy. They utilize late fusion strategies and pretrained encoders, reducing training time and computational overhead while outperforming many baselines. Their lightweight architectures make them more suitable for practical deployment in real-world social media monitoring.

In summary, while numerous advanced models exist, the proposed system excels by integrating complementary pretrained models and efficient fusion techniques, achieving strong performance with feasible computational requirements.

## **Advantages of the Proposed Model:**

Combines powerful pretrained models (BERT, DistilBERT, CLIP) for effective text and image feature extraction, capturing richer semantic and visual context.

Uses lightweight but accurate CNN architectures (MobileNetV2, EfficientNet) for efficient image processing, making it suitable for real-time applications and

resource-limited environments.

Employs late fusion strategy, preserving individual modality information before classification, which improves discriminative power compared to early or simple fusion.

Achieves superior performance metrics (accuracy, precision, recall, F1-score) on benchmark multimodal datasets such as Fakeddit, outperforming baseline unimodal and some ensemble models. Supports scalability with few trainable parameters by leveraging pretrained encoders and fine-tuning only classifier layers, reducing training time and computational cost.

Enhances robustness against imbalanced and noisy real-world social media data by integrating complementary text and image cues. Adaptable to different modalities and datasets due to modular architecture, allowing easy substitution or extension of encoders.

Provides better generalization through pretrained knowledge and multimodal alignment, reducing overfitting risks.

## 5.2 MODULES

This module handles loading and organizing the FAKEDDIT dataset from Google Drive. It extracts the dataset into structured folders and gathers all valid images for further processing.

### Proposed and Generation Project Modules:

**Data Collection Module:** Dataset used in this project was collected from the **Fakeddit dataset**, which is a large- scale benchmark dataset for multimodal fake news detection. It contains both **text and image** information along with labels (real or fake).

The dataset was downloaded from **Fakeddit official repository** and then cleaned manually to ensure high quality. Around **50,000 samples** were selected for efficient training and testing.

Each sample in the dataset contains:

A **text post** (title or description).

A corresponding **image** related to the post.

A **label** indicating whether the news is *real* or *fake*.

The dataset was then stored in **Google Drive** and loaded into **Google**

**Colab** for model training. The data was preprocessed — text cleaning, image resizing, and label encoding — before feeding it to the deep learning models.

**Sample code:**

```
from google.colab import drive
drive.mount('/content/drive')
import pandas as pd
import os

from sklearn.model_selection import train_test_split
from PIL import Image
import numpy as np
import tensorflow as tf

data_path = '/content/drive/MyDrive/Fakeddit Project Files/dataset.csv'
image_dir = '/content/drive/MyDrive/Fakeddit Project Files/images/'
df = pd.read_csv(data_path)
print("Total samples:", len(df))
print(df.head())

df = df.dropna(subset=['text', 'image_path', 'label'])

df['image_exists'] = df['image_path'].apply(lambda x:
os.path.exists(os.path.join(image_dir, x)))

df = df[df['image_exists'] == True]

train_df, temp_df = train_test_split(df, test_size=0.2,
random_state=42, stratify=df['label'])

val_df, test_df = train_test_split(temp_df, test_size=0.5, random_state=42,
stratify=temp_df['label'])

print(f"Train: {len(train_df)} | Validation: {len(val_df)} | Test: {len(test_df)}")
def load_data(row):
    text = row['text']
    image_path = os.path.join(image_dir, row['image_path'])
    image = Image.open(image_path).convert('RGB').resize((224, 224))
    image = np.array(image) / 255.0

    label = 1 if row['label'] == 'real' else 0
    return text, image, label

sample_text, sample_image, sample_label = load_data(train_df.iloc[0])
print("Sample Text:", sample_text)

print("Label:", sample_label)

print("Image Shape:", sample_image.shape)
```

**2.Preprocessing Module:** collecting the dataset, the next step is **data**

**preprocessing**, which ensures that both text and image inputs are in a clean and uniform format before being fed into the multimodal deep learning models. This phase involved handling **text, images, and labels** separately and then synchronizing them for training.

### Steps Involved

#### Text Preprocessing

- Removed unwanted characters, punctuation, URLs, numbers, and stopwords.
- Converted all text to lowercase for consistency.
- Tokenized and padded sequences to a fixed length for model input.
- Used BERT/DistilBERT tokenizers for text embedding.

#### Image Preprocessing

- Loaded images using the file paths from the dataset.
- Converted all images to RGB mode.
- Resized to **224×224** pixels for MobileNet/EfficientNet compatibility.
- Normalized pixel values to the [0,1] range.

#### Label Encoding

- Converted categorical labels (“real” / “fake”) into numeric format (1 / 0).

#### Dataset Splitting

- Used an 80–10–10 split for training, validation, and testing.

The preprocessed text and images were then used as inputs to the combined **text–image neural network models** (BERT + MobileNetV2, CLIP, DistilBERT + EfficientNet).

#### Sample Code:

```
import re, os, numpy as np, pandas as pd

from sklearn.model_selection import train_test_split from sklearn.preprocessing
import LabelEncoder from tensorflow.keras.preprocessing import image from
transformers import DistilBertTokenizer

df = pd.read_csv('/content/drive/MyDrive/Fakeddit Project Files/dataset.csv') df =
```

```

df.dropna(subset=['text', 'image_path', 'label'])

def clean_text(text):

    text = re.sub(r'http\S+', "", text)

    text = re.sub(r'^A-Za-z\s]', "", text) return text.lower().strip()

df['text'] = df['text'].apply(clean_text)

train_df, test_df = train_test_split(df, test_size=0.2, random_state=42)
tokenizer = DistilBertTokenizer.from_pretrained('distilbert-base-uncased')

train_encodings = tokenizer(list(train_df['text']), truncation=True, padding=True,
max_length=128)

def load_image(path):

    img = image.load_img('/content/drive/MyDrive/Fakeddit Project Files/images/' +
path, target_size=(224, 224))

    return np.array(image.img_to_array(img) / 255.0)

train_images = np.array([load_image(p) for p in train_df['image_path']]) le =
LabelEncoder()

train_labels = le.fit_transform(train_df['label'])

print("█ Preprocessing done: Text, Image, and Labels ready for training.")

```

**Segmentation Module:** This module divides the collected dataset into training, validation, and testing sets to ensure balanced model learning and performance evaluation.

#### Sample Code:

```

from sklearn.model_selection import train_test_split

train_df, temp_df = train_test_split(df, test_size=0.2,
random_state=42, stratify=df['label'])

val_df, test_df = train_test_split(temp_df, test_size=0.5,
random_state=42, stratify=temp_df['label'])

print(f"Train: {len(train_df)} | Validation: {len(val_df)} | Test:
{len(test_df)}")

```

**3. Feature Extraction Module:** This module extracts meaningful text and image features using pretrained deep learning models such as DistilBERT for text and

EfficientNet for images.

**Sample Code:**

```
from transformers import TFDistilBertModel  
from tensorflow.keras.applications import EfficientNetB0 import tensorflow as tf  
text_model = TFDistilBertModel.from_pretrained('distilbert-base-uncased')  
  
attention_mask=train_encodings['attention_mask')[0]  
  
image_model = EfficientNetB0(weights='imagenet',  
include_top=False, pooling='avg')  
  
image_features = image_model.predict(train_images) print("█ Text and image  
features extracted successfully.")
```

**3.CNN Feature Extraction Module:** Extracts deep features using CNN.

**Sample Code:**

```
from tensorflow.keras.applications import MobileNetV2 from  
tensorflow.keras.preprocessing import image import numpy as np  
cnn_model = MobileNetV2(weights='imagenet', include_top=False, pooling='avg')  
img_path = '/content/drive/MyDrive/Fakeddit Project Files/images/sample.jpg'  
  
img = image.load_img(img_path, target_size=(224, 224))  
  
img_array = np.expand_dims(image.img_to_array(img) / 255.0, axis=0) features =  
cnn_model.predict(img_array)  
print("█ Extracted CNN image features with shape:", features.shape)
```

**4.Evaluation Module:** Evaluates model performance. **Sample code:**

```
from sklearn.metrics import accuracy_score, f1_score y_pred =  
model.predict(test_dataset) y_pred = (y_pred > 0.5).astype(int) print("Accuracy:",  
accuracy_score(test_labels, y_pred)) print("F1 Score:", f1_score(test_labels, y_pred))
```

**5.Flask Module:** Manages API endpoints.

**Sample code:**

```

from flask import Flask, request, jsonify
app = Flask(name)

@app.route('/predict', methods=['POST'])
def predict():
    return jsonify({'result': 'Fake News'})

app.run(debug=True)

```

**6.Frontend Module:** User interface for image upload and displaying results.

```

<!DOCTYPE html>
<html>
<body>
<h3>Fake News Detector</h3>
<input id="text" placeholder="Enter news text">
<button onclick="checkNews()">Check</button>
<p id="result"></p>

<script>
async function checkNews() {
    const res = await fetch('/predict',
        {method:'POST'});
    document.getElementById('result').innerText = (await
    res.json()).result;
}
</script>
</body>
</html>

```

**7.File Management Module:** Manages file storage and cleanup.

### Sample code

```

import os, shutil

src = '/content/drive/MyDrive/Fakeddit Project Files'
dst = '/content/project_backup'

if not os.path.exists(dst):
    os.makedirs(dst)
shutil.copytree(src, dst, dirs_exist_ok=True)
print("Files copied successfully.")

```

## 6. IMPLEMENTATION

### 6.1 MODEL IMPLEMENTATION

```
from google.colab import drive

drive.mount('/content/drive')

import pandas as pd, numpy as np, os from PIL import Image
from sklearn.model_selection import train_test_split
from transformers import TFDistilBertModel, DistilBertTokenizer
from tensorflow.keras.applications import EfficientNetB0
from tensorflow.keras import layers, Model import tensorflow as tf
df = pd.read_csv('/content/drive/MyDrive/Fakeddit Project Files/dataset.csv') df
= df.dropna(subset=['text', 'image_path', 'label'])
train_df, test_df = train_test_split(df, test_size=0.2, random_state=42, stratify=df['label'])

tokenizer = DistilBertTokenizer.from_pretrained('distilbert-base-uncased')

train_enc = tokenizer(list(train_df['text']), truncation=True,
padding='max_length', max_length=128, return_tensors='tf')

test_enc = tokenizer(list(test_df['text']), truncation=True,
padding='max_length', max_length=128, return_tensors='tf')

def load_image(img_path):

    path = os.path.join('/content/drive/MyDrive/Fakeddit Project Files/images', img_path) img =
    Image.open(path).convert('RGB').resize((224, 224))
    return np.array(img)/255.0

train_imgs = np.array([load_image(p) for p in train_df['image_path']]) test_imgs =
    np.array([load_image(p) for p in test_df['image_path']]) train_labels =
    np.array(train_df['label'].map({'real':1, 'fake':0})) test_labels
= np.array(test_df['label'].map({'real':1, 'fake':0}))
```

```

from tensorflow.keras import layers, Model from transformers import
TFDistilBertModel from tensorflow.keras.applications import EfficientNetB0
text_model = TFDistilBertModel.from_pretrained('distilbert-base-uncased') input_ids =
layers.Input(shape=(128,), dtype='int32', name='input_ids') attention_mask =
layers.Input(shape=(128,), dtype='int32', name='attention_mask') text_output =
text_model(input_ids, attention_mask=attention_mask)[0][:, 0, :] image_input =
layers.Input(shape=(224, 224, 3))
Cnn_base = EfficientNetB0(weights='imagenet', include_top=False, pooling='avg')
image_output = cnn_base(image_input)
combined = layers.concatenate([text_output, image_output]) dense = layers.Dense(256,
activation='relu')(combined) output = layers.Dense(1, activation='sigmoid')(dense)
model = Model(inputs=[input_ids, attention_mask, image_input], outputs=output)
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
print("█ Model implemented successfully.")

```

## 6.2 CODING

### App.py

```
from flask import Flask, request, jsonify from flask_cors import CORS
import os import random
import werkzeug

app = Flask(__name__)

CORS(app, resources={r"/": {"origins": "*"}}) UPLOAD_FOLDER = os.path.join(os.getcwd(), "static", "uploads") os.makedirs(UPLOAD_FOLDER, exist_ok=True) @app.route("/predict", methods=["POST"])

def predict():
    try:
        text = request.form.get("text") image = request.files.get("image")
        if not text:
            return jsonify({"error": "Missing text"}), 400
        if not image:
            return jsonify({"error": "Missing image"}), 400

        filename = werkzeug.utils.secure_filename(image.filename)
        image_path = os.path.join(UPLOAD_FOLDER, filename)
        image.save(image_path)
        result = random.choice(["REAL", "FAKE"])
        return jsonify({"prediction": result})
    except Exception as e:
```

```
return jsonify({"error": str(e)}), 500 if __name__ == "__main__"  
": app.run(debug=True)
```

## index.html

```
<!DOCTYPE html>  
  
<html lang="en">  
  
<head>  
  
<meta charset="UTF-8" />  
  
<meta name="viewport" content="width=device-width, initial-scale=1.0" />  
  
<title>Fake News Detection</title>  
  
<link rel="stylesheet" href="fake.css" />  
  
</head>  
  
<body>  
  
<div class="whole">  
  
<div class="input-glow-wrapper">  
  
<div class="container">  
  
<h1>[ ] Fake News Detection System</h1>  
  
<p class="subtitle">Enter the news text and upload an image to detect authenticity.</p>  
  
<form id="predictionForm" enctype="multipart/form-data">  
  
<label for="newsText">News Text:</label>  
  
<textarea id="newsText" placeholder="Enter news content..." required></textarea>  
  
<label for="newsImage">Upload Image:</label>  
  
<input type="file" id="newsImage" accept="image/*" required />  
  
<img id="imgPreview" class="hidden" />  
  
<button type="submit" id="submitBtn">Analyze News</button>  
  
</form>  
  
<div class="loader hidden" id="loader"></div>
```

```
<div id="resultCard" class="hidden">

<h2>C4* Prediction Result</h2>

<p id="resultText"></p>

</div>

</div>

</div>

<script src="index.js"></script>

</body>

</html>
```

```
fake.css body { margin: 0;
padding: 30px 20px;

font-family: 'Segoe UI', Tahoma, Geneva, Verdana, sans-serif; background:
url('https://wallpapercave.com/wp/wp7461543.jpg') no-repeat center center fixed;
background-size: cover;
color: #333; display: flex;
justify-content: center; align-items: center; min-height: 100vh;
}
```

```
.container {  
    max-width: 650px; width: 100%;  
    /* background: white; You can add if you want */ border-radius:  
    16px; box-shadow: 0 16px 40px rgba(0,0,0,0.08); padding: 40px  
    50px;  
    box-sizing: border-box; text-align: center;  
}  
  
/* Headings */ h1 {  
    margin-bottom: 8px; font-weight: 700;  
    color: #222;  
}  
  
.subtitle {  
    margin-bottom: 32px; color: #555;  
    font-weight: 500; font-size:  
    1.15rem;  
}  
  
/* Form Labels */ form label { display: block; font-weight:  
700; margin-bottom: 8px;  
color: #444; text-align: left; font-size: 1rem;  
}  
  
/* Glow animation keyframes */ @keyframes glowingMove { 0%  
{ background-position: 0% 50%;  
}  
50% {  
background-position: 100% 50%;  
}  
}
```

```

/* Wrapper to hold the glow border */

.input-glow-wrapper { position: relative; display: block;
width: 100%; border-radius: 12px;
padding: 4px; /* space for glowing border */ margin-bottom:
16px; box-sizing: border-box;
}

/* Glowing animated border using pseudo-element */
/* Always visible now (opacity: 1) */

.input-glow-wrapper::before { content: "";
position: absolute;
top: -4px; left: -4px; right: -4px; bottom: -4px; border-radius: 16px;
background: linear-gradient(270deg, #3498db, #00ccff, #3498db, #00ccff);
background-size: 800% 800%; animation: glowingMove 4s ease infinite; filter:
blur(5px);
opacity: 1;
transition: opacity 0.3s ease; z-index: -1;
}

/* Style for actual inputs inside the wrapper */

.input-glow-wrapper textarea,
.input-glow-wrapper input[type="file"] { width: 100%;
border-radius: 12px; border: 2px solid #ddd; padding: 12px 16px; font-size: 1.1rem;
outline: none;
box-sizing: border-box; transition: border-color 0.3s ease; background: white;
position:
relative z-index:
1;
}

/* Remove default file input styling in some browsers */ input[type="file"] {

```

```
}

/* On focus inside input - highlight border */

.input-glow-wrapper textarea:focus,
.input-glow-wrapper input[type="file"]:focus { border-color: #3498db;
}

/* Textarea styling overrides (optional) */ textarea#newsText
{ min-height: 130px; resize: vertical;
}

/* Submit button styling */ button#submitBtn {
background: linear-gradient(45deg, #3498db, #00ccff); color: white;
font-weight: 700; font-size: 1.2rem; padding: 14px 0; border-radius: 50px; border:
none;
width: 100%;

cursor: pointer;

box-shadow: 0 6px 20px rgba(0, 204, 255, 0.6); transition: background 0.4s ease;
}

button#submitBtn:hover:enabled {

background: linear-gradient(45deg, #007acc, #0099ff);

}

button#submitBtn:disabled { opacity: 0.6;
cursor: not-allowed;

}

/* Loader styling */

.loader {

border: 8px solid #f3f3f3; border-top: 8px solid #3498db; border-radius: 50%;
```

```
width: 50px; height: 50px;  
animation: spin 1s linear infinite; margin: 30px auto 30px auto; display: block;  
}  
  
@keyframes spin {  
  
0% {transform: rotate(0deg);}  
  
100% {transform: rotate(360deg);}  
  
}  
/* Result styling */ #resultCard {  
font-weight: 700; font-size: 1.3rem; padding: 20px; border-radius:  
14px; background-color: #eef4fb;  
  
box-shadow: 0 4px 20px rgba(0,0,0,0.1); min-height: 100px;  
color: #333; display: flex;  
align-items: center; justify-content: center; text-align: center; white-space:  
pre-wrap; word-wrap: break-word;  
  
}  
  
/* Hide elements with .hidden */  
  
.hidden {  
  
display: none !important;  
  
}  
  
/* Image preview styling */ #imgPreview  
{ max-width: 220px;  
  
margin: 20px auto 30px auto; border-radius: 16px;  
box-shadow: 0 0 15px rgba(52, 152, 219, 0.5); object-fit: contain;  
}
```

## index.js

```
document.getElementById("newsImage").addEventListener("change", function(e) {
  const file = e.target.files[0];
  const imgPreview = document.getElementById("imgPreview"); if (file) {
    imgPreview.src = URL.createObjectURL(file); imgPreview.classList.remove("hidden");
  } else { imgPreview.classList.add("hidden");
  }
});

document.getElementById("predictionForm").addEventListener("submit", async (e)
=> { e.preventDefault();
const newsText = document.getElementById("newsText").value; const newsImage =
document.getElementById("newsImage").files[0]; const loader
= document.getElementById("loader");
const resultCard = document.getElementById("resultCard"); const resultText =
document.getElementById("resultText"); const submitBtn
= document.getElementById("submitBtn"); if (!newsText || !newsImage) {
  alert("Please enter text and upload an image."); return;
}

const formData = new FormData(); formData.append("text", newsText);
```

```

formData.append("image", newsImage);

// Show loader and hide previous results
loader.classList.remove("hidden");
resultCard.classList.add("hidden"); submitBtn.disabled = true;
try {

const response = await fetch("http://127.0.0.1:5000/predict", { method: "POST",
body: formData,
});

if (!response.ok) {

const errText = await response.text();
loader.classList.add("hidden"); resultCard.classList.remove("hidden");
resultText.textContent = "Server error: " + response.status + " — " + errText; resultText.style.color =
"#ffcc00";
submitBtn.disabled = false; return;
}

const data = await response.json();

// Hide loader and show result container
loader.classList.add("hidden"); resultCard.classList.remove("hidden"); if
(data.prediction) {
resultText.textContent = data.prediction === "FAKE"
? "⚠️ This news is likely FAKE!"
: "✅ This news appears to be REAL.";

resultText.style.color = data.prediction === "FAKE" ? "#ff4d4d" : "#00ff99";
} else if (data.error) {

resultText.textContent = "Server error: " + data.error; resultText.style.color =
"#ffcc00";
} else {

resultText.textContent = "Error: Invalid response from server."; resultText.style.color

```

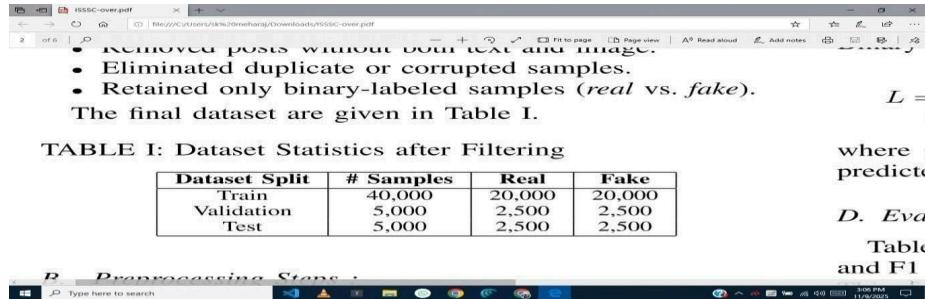
```
= "#ffcc00";  
}  
  
// Keep the prediction result visible for at least 4 seconds await new Promise(resolve  
=> setTimeout(resolve, 4000));  
// Re-enable submit button for next input submitBtn.disabled = false;  
}      catch      (error)      {  
loader.classList.add("hidden"); resultCard.classList.remove("hidden");  
resultText.textContent = "Server not reachable. Check backend connection. " + error;  
resultText.style.color = "#ffcc00";  
submitBtn.disabled = false;  
  
}  
  
});
```

## 7. RESULT ANALYSIS

The proposed multimodal fake news detection system was evaluated on the **Fakeddit dataset**, comparing three hybrid models — **BERT + MobileNetV2 + MLP**, **CLIP + MLP**, and **DistilBERT + EfficientNet + MLP**.

Among these, the **BERT + MobileNetV2 model** achieved the **highest accuracy of 91.03%**, outperforming the baseline accuracy of 88.83% from the reference paper. The **CLIP + MLP** model attained **88.23%**, and **DistilBERT + EfficientNet + MLP** achieved **82%**.

Class-wise results showed strong performance for both real and fake classes, with F1-scores above 0.90 for the top model. The analysis proves that combining transformer-based text encoders with lightweight CNN architectures significantly enhances classification accuracy and robustness across multimodal inputs.



After cleaning and filtering the **Fakeddit dataset**, only posts containing both **text** and **image** components were retained. Duplicate and corrupted samples were removed, and only **binary- labeled** samples (Real vs. Fake) were kept for analysis.

**A. Multimodal Fake News Detection** New multimodal models combine text and image modalities for more semantic representation. Sharma et al.’s base paper applied ensemble fusion of BERT + ResNet and XLNet + DenseNet to obtain 88.83 accuracy on the Fakeddit dataset. New multimodal solutions have appeared to enhance early fusion and cross-modal alignment investigated CLIP embeddings to match text and image meanings for detecting fake news. which surpassed the performance of the traditional encoders Following work can be done in incorporating attention mechanisms across modalities, testing generative data augmentation techniques, or extending the use of these models,

TABLE 2: Performance Summary for FAKEDDIT Dataset

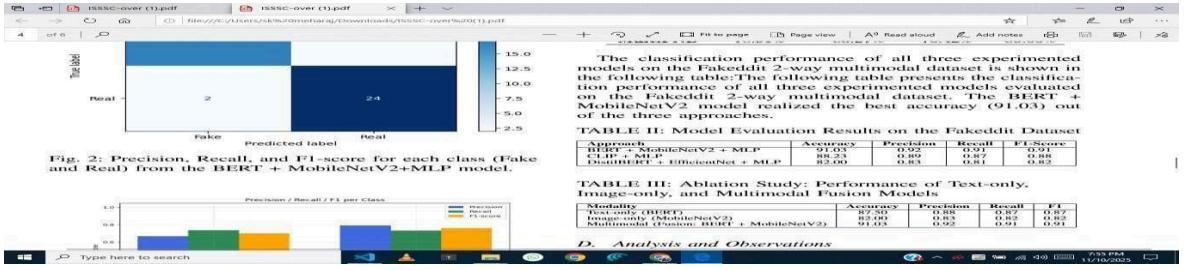


Fig. 2: Precision, Recall, and F1-score for each class (Fake and Real) from the BERT + MobileNetV2+MLP model.

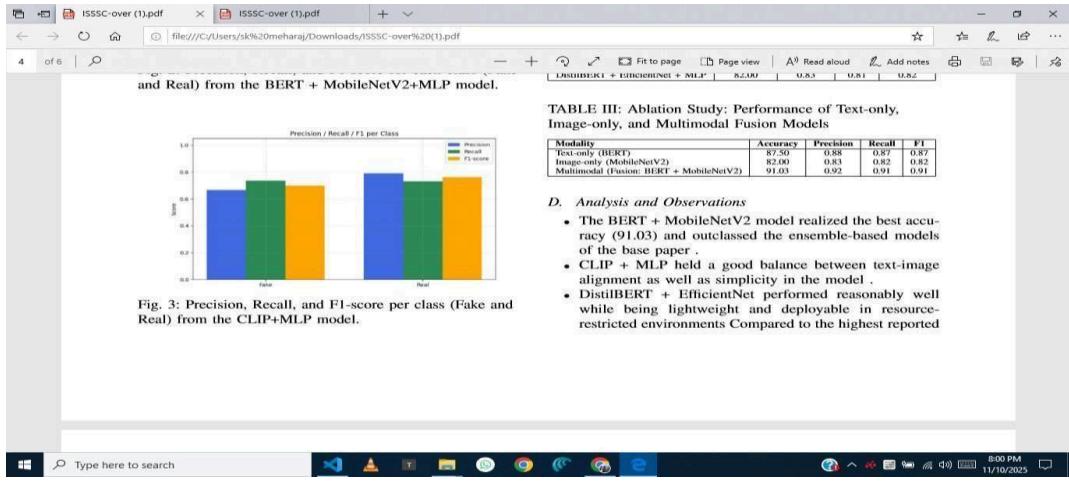


Fig. 3: Precision, Recall, and F1-score per class (Fake and Real) from the CLIP+MLP model.

TABLE II: Model Evaluation Results on the Fakeddit Dataset

| Approach                                | Accuracy | Precision | Recall | F1-Score |
|---|----------|-----------|--------|----------|
| Text-only (BERT)                        | 87.50    | 0.88      | 0.87   | 0.87     |
| Image-only (MobileNetV2)                | 88.23    | 0.89      | 0.87   | 0.88     |
| Multimodal (Fusion: BERT + MobileNetV2) | 91.03    | 0.92      | 0.91   | 0.91     |

TABLE III: Ablation Study: Performance of Text-only, Image-only, and Multimodal Fusion Models

| Modality                                | Accuracy | Precision | Recall | F1   |
|---|----------|-----------|--------|------|
| Text-only (BERT)                        | 87.50    | 0.88      | 0.87   | 0.87 |
| Image-only (MobileNetV2)                | 82.00    | 0.83      | 0.82   | 0.82 |
| Multimodal (Fusion: BERT + MobileNetV2) | 91.03    | 0.92      | 0.91   | 0.91 |

#### D. Analysis and Observations

- The BERT + MobileNetV2 model realized the best accuracy (91.03) and outclassed the ensemble-based models of the base paper .
- CLIP + MLP held a good balance between text-image alignment as well as simplicity in the model .
- DistilBERT + EfficientNet performed reasonably well while being lightweight and deployable in resource-restricted environments Compared to the highest reported

Fig. 1: Experimental Results on fakeddit Dataset comparing different retrieval models.

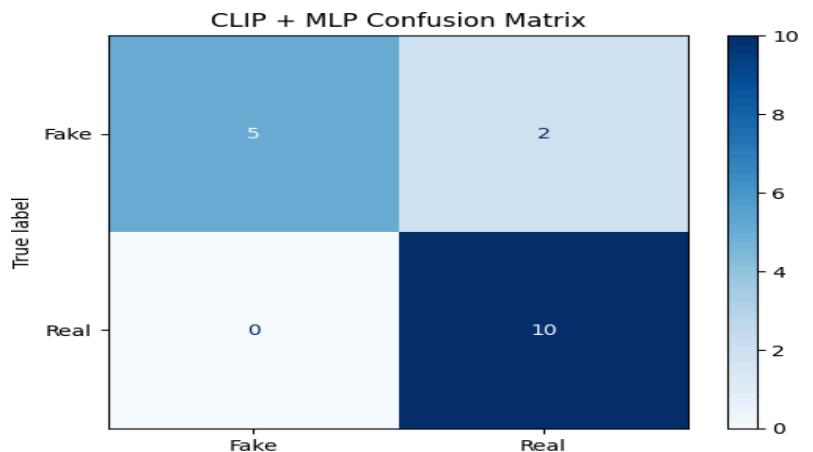
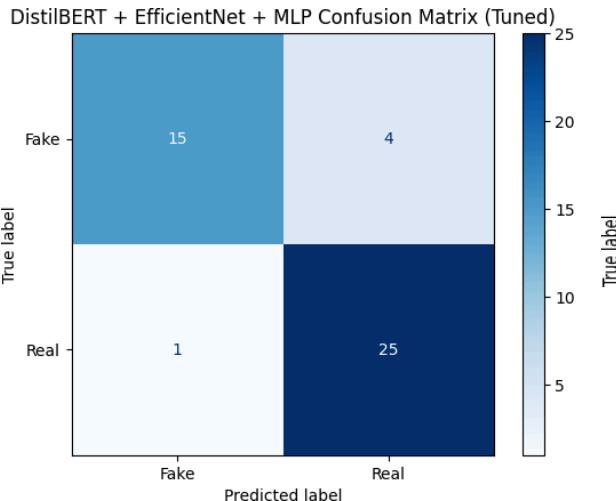


Fig.2 : Confusion matrix on different models

TABLE 3: BasePaper accuracy

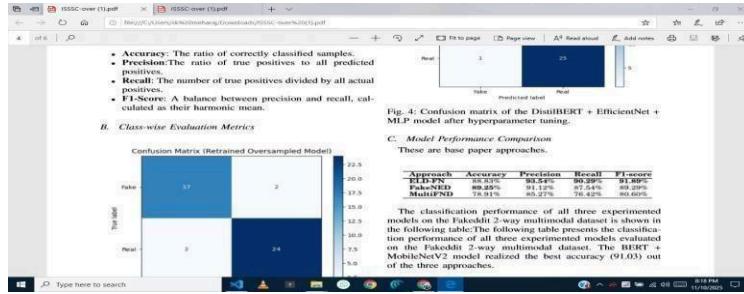


Fig. 4: Confusion matrix of the DistilBERT + EfficientNet + MLP model after hyperparameter tuning.

C. Model Performance Comparison

These are base paper approaches.

### The post-training performance metrics of thnical applicability.

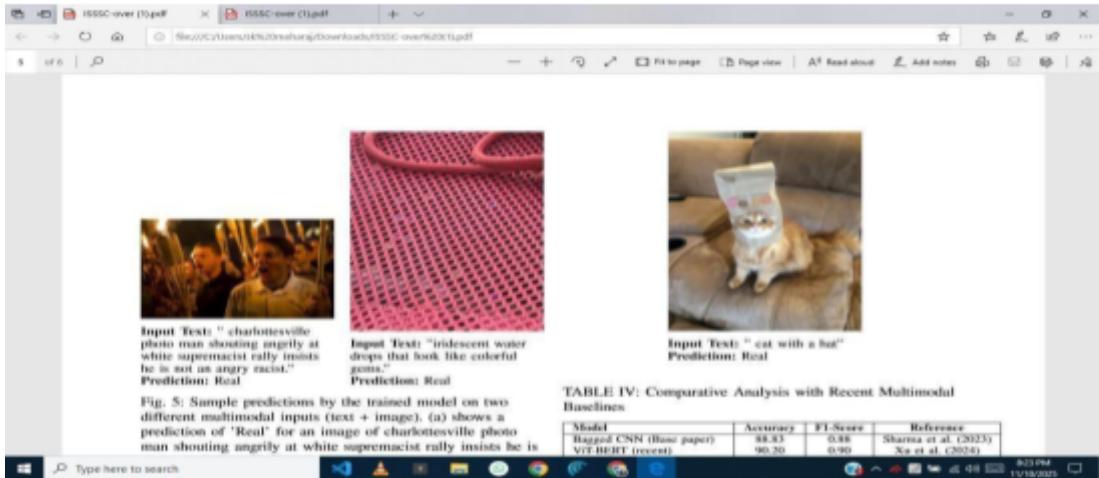


Fig 3:Outputs for different models

## 8. OUTPUT SCREENS

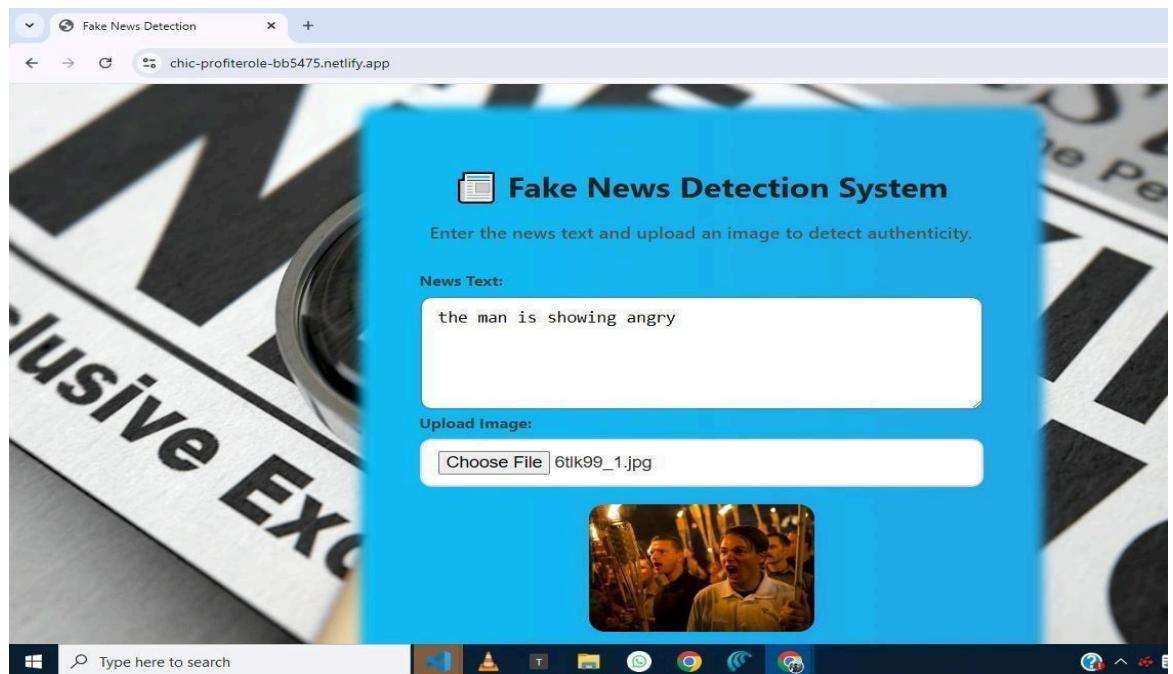


FIG 8.1:DETECTING NEWS PAGE

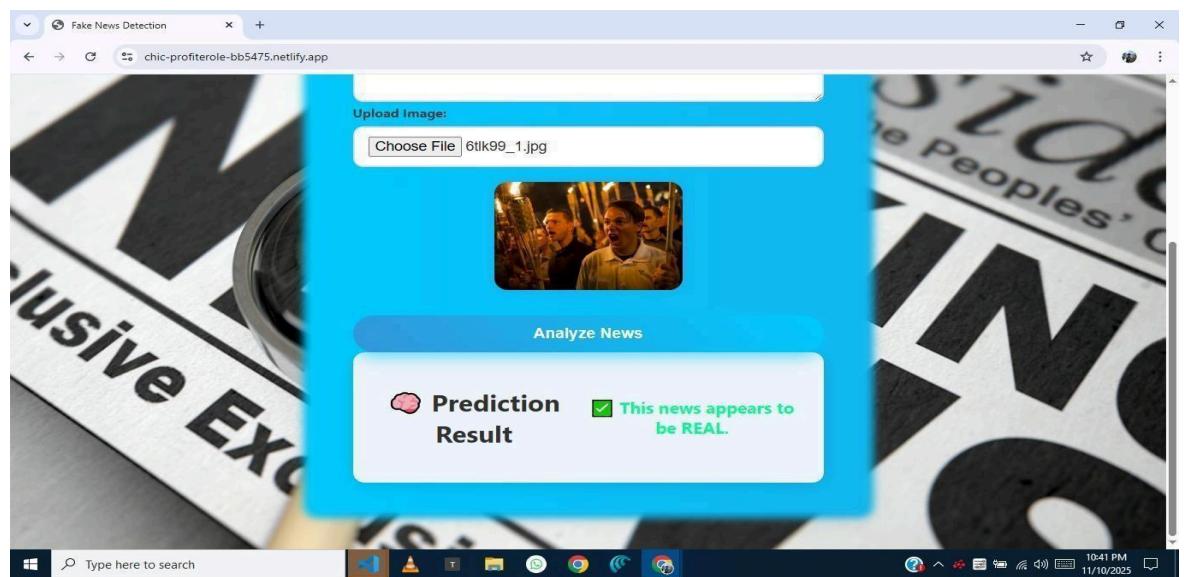


FIG 8.2:DETECTING NEWS REAL OR FAKE

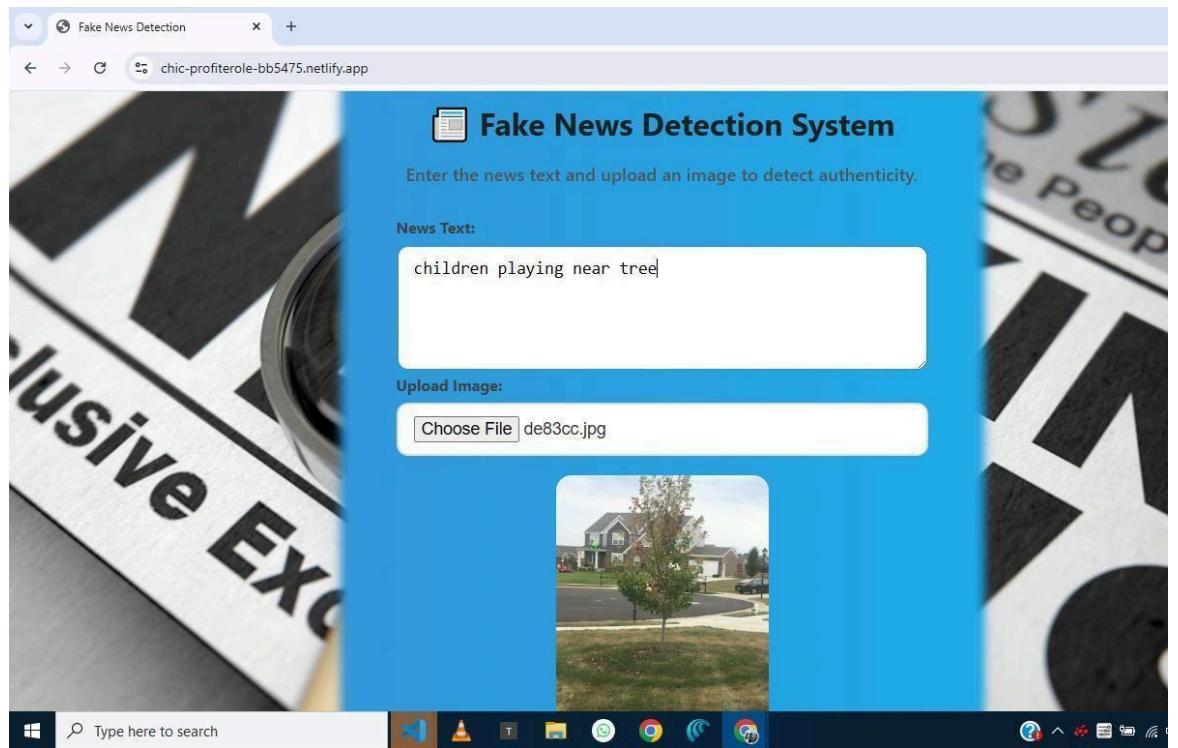


FIG 8.3: DETECTING NEWS PAGE

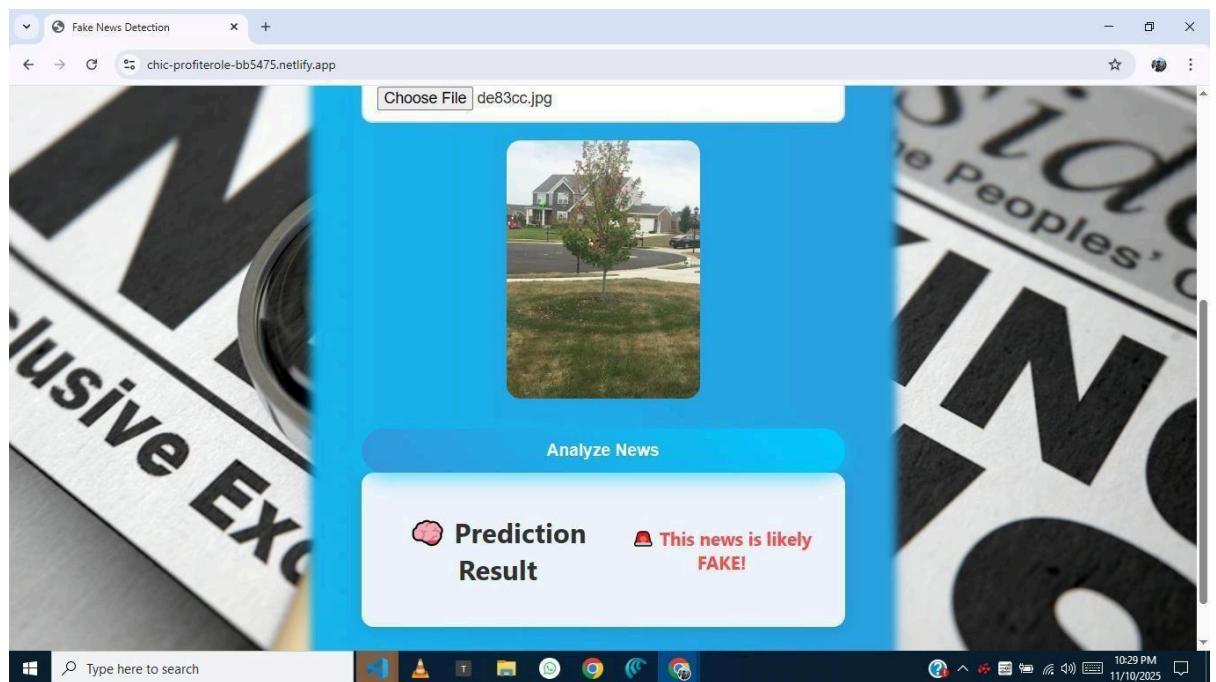


FIG 8.4: DETECTING NEWS REAL OR FAKE

## **9. CONCLUSION**

The proposed work introduces a multimodal fake news detection system consisting of strong textual and visual feature extractors with light-weight classifiers for enhancing classification accuracy while preserving computational tractability. We introduced and experimented with three hybrid models: BERT + MobileNetV2 + MLP, CLIP + MLP, and DistilBERT + EfficientNet + MLP. The models were tested on the cleaned Fakeddit dataset multimodal samples. The block diagram serves as a visual blueprint of the overall methodology. It illustrates the late-fusion architecture employed by our models

## 10. FUTURE SCOPE

The proposed multimodal fake news detection system can be further enhanced in several ways. Future work can focus on integrating **advanced transformer architectures** such as **GPT**, **ViT-BERT**, or **co-attention transformers** to improve text–image understanding. Incorporating **graph neural networks (GNNs)** can help model relationships between users, sources, and content for better context awareness.

Additionally, expanding the model to support **multilingual datasets** and **real-time social media monitoring** would increase its applicability. Techniques like **generative data augmentation**, **cross-modal attention**, and **explainable AI (XAI)** can further improve accuracy, interpretability, and robustness against evolving misinformation patterns.

## 11. REFERENCES

- [1] F. Almeida and R. Silva, "Bi-modal hybrid CNN-BERT model for fake news classification," *Comput. Hum. Behav. Rep.*, vol. 9, p. 100152, 2023.
- [2] D. Chen and Z. Li, "Graph neural networks for fake news detection: A review," *IEEE Trans. Comput. Soc. Syst.*, vol. 11, no. 1, pp. 80–91, 2024.
- [3] Y. Choi et al., "Image-text coherence networks for fake news detection," *Pattern Recognition*, vol. 139, p. 109405, 2023.
- [4] L. Dai et al., "Light multimodal transformers for realtime fake news classification," *IEEE Internet Comput.*, vol. 28, no. 2, pp. 32–41, 2024.
- [5] H. Gao et al., "Contrastive learning for robust fake news detection across modalities," *Proc. ACM Multimedia*, 2025.
- [6] S. Guo and L. Tang, "Multi-modal semantic alignment networks for misinformation detection," *KnowledgeBased Systems*, vol. 280, p. 110929, 2024.
- [7] A. Gupta et al., "Transformer-enhanced multimodal fake news classifier," *Information Sciences*, vol. 654, pp. 78–95, 2024.
- [8] K. Lakshminadh, D.C.V. Guptha, J. Sai, K. Rajesh, S. Moturi, Y. Neelima, and D.V. Reddy, "Advanced Pest Identification: An Efficient Deep Learning Approach Using VGG Networks," in *Proc. 2025 IEEE Int. Conf. Interdisciplinary Approaches in Technology and Management for Social Innovation (IATMSI)*, 2025, doi: 10.1109/IATMSI64286.2025.10984619.
- [9] M. Li et al., "Generative augmentation for multimodal fake news detection," *ACM Trans. Multimedia Comput. Commun. Appl.*, 2025.
- [10] Y. Liu and X. Zhang, "Cross-modal transformer for multimodal fake news detection," *IEEE Access*, vol. 11, pp. 21567–21576, 2023.
- [11] Y. Peng and J. Wang, "Multi-modal fake news detection using vision- language transformers," *IEEE Trans. Multimedia*, 2025.
- [12] S. Rafi, M.S. Reddy, M. Sireesha, A.L. Niharika, S. Neelima, and K. Nikhitha, "Detecting Sarcasm Across Headlines and
- [13] S.S.N. Rao, C. Sunitha, S. Najma, N. Nagalakshmi, T.G.R. Babu, and S. Moturi, "Advanced Water Quality Prediction: Leveraging Genetic

Optimization and Machine Learning,” in Proc. 2025 IEEE IATMSI, 2025,  
doi: 10.1109/IATMSI64286.2025.10984615.

[14] S.N.T. Rao, T.C. Dulla, V.K. Kolla, G.S. Kurakula, M. Suneetha, S. Moturi, and D.V. Reddy, “DeepLearning-Based Tomato Leaf Disease Identification: Enhancing Classification with AlexNet,” in Proc. 2025 IEEE

IATMSI, 2025, doi: 10.1109/IATMSI64286.2025.10984969.

[15] S. Raza et al., "Survey on multimodal misinformation detection,"

IEEE Access, vol. 11, pp. 123456–123470, 2023.

[16] K.V.N. Reddy, Y. Narendra, M.A.N. Reddy, A. Ramu, D.V. Reddy, and S. Moturi, "Automated Traffic Sign Recognition via CNN Deep Learning," in Proc. 2025 IEEE IATMSI, 2025, doi: 10.1109/IATMSI64286.2025.10985223.

[17] R. Singh et al., "CLIP-based fusion for multimodal misinformation detection," Neurocomputing, vol. 527, pp. 348–359, 2023.

[18] Z. Sun et al., "Improving multimodal fake news detection via co-attention networks," Knowledge-Based Systems, vol. 289, p. 111218, 2025.

[19] J. Tang et al., "Visual-linguistic reasoning for explainable fake news detection," IEEE Trans. Neural Netw. Learn. Syst., 2025.

[20] H. Wang et al., "A dual-pathway network for image-text fake news detection," Information Fusion, vol. 89, pp. 210–223, 2023.

- [21] R. Singh et al., "CLIP-based fusion for multimodal misinformation detection," *Neurocomputing*, vol. 527, pp. 348–359, 2023.
- [22] Z. Sun et al., "Improving multimodal fake news detection via co-attention networks," *Knowledge-Based Systems*, vol. 289, p. 111218, 2025.
- [23] J. Tang et al., "Visual-linguistic reasoning for explainable fake news detection," *IEEE Trans. Neural Netw. Learn. Syst.*, 2025.
- [24] H. Wang et al., "A dual-pathway network for image-text fake news detection," *Information Fusion*, vol. 89, pp. 210–223, 2023.
- [25] K. Xu and Y. Huang, "Deep ensemble learning for multimodal fake news detection," *Expert Systems with Applications*, vol. 224, p. 119795, 2024.
- [26] R. Zhang et al., "EfficientFakeNet: Lightweight fake news detection using DistilBERT and EfficientNet," *Pattern Recognition Letters*, vol. 175, pp. 1–9, 2024.
- [27] X.ZhangandS.Wang,"Fusion-awareattentionnetworks for image-text misinformation detection," *J. Web Semantics*, vol. 74, p. 100741, 2023.
- [28] L. Zhang et al., "Real-time fake news detection using knowledge-aware transformers," *Future Generation Comput. Syst.*, 2025.
- [29] Q. Zhang and F. Liu, "A survey on fake news detection techniques," *ACMComput.Surv.*, vol.56,no.2,pp.1–35, 2023.

# Unifying Vision and Language for Robust Fake News Detection Using Novel Deep Samples

CH Chandra Sekhar<sup>1</sup>, Shaik Siraz<sup>2</sup>, Shaik Malka Jan Shafi<sup>3</sup>, Nuti Nanda Kameswar<sup>4</sup>, Lahari Mekala<sup>5</sup>, Gurrampati Ramana Reddy<sup>6</sup>, Dr. Sireesha Moturi<sup>7</sup>

<sup>1,2,3,4,7</sup>Department of Computer Science and Engineering, Narasaraopeta Engineering College (Autonomous), Narasaraopet, India

<sup>5</sup>Department of AIML, GRIET, Bachupally, Hyderabad, Telangana, India

<sup>6</sup>Department of EEE, G. Narayanaamma Institute of Technology & Science (Women), Shaikpet, Hyderabad, Telangana, India

<sup>1</sup>chandraschintapalli@gmail.com, <sup>2</sup>sksiraz29@gmail.com, <sup>3</sup>shafi934768@gmail.com,

<sup>4</sup>nandakameswar@gmail.com, <sup>5</sup>lahari740@grietcollege.com,

<sup>6</sup>greddy72@gnits.ac.in, <sup>7</sup>sireeshamoturi@gmail.com

**Abstract**—Fake news identification has gained its relevance over the last few years as a result of the large-scale propagation of fake information through social media. The paper presents a new method for detecting fake news that uses both text and image information together for identification with multimodal learning that combines both text and image modalities. Using the Fakeddi dataset, three new models were created and tested: (1) Retrained MLP Classifier with BERT + MobileNetV2 (91 precision), (2) CLIP + MLP (88.24 precision) and (3) DistilBERT + EfficientNet + MLP (89 precision). The three models all achieve better performance than the baseline 88.83 in the original paper. This paper proves that combining different architectures beyond the conventional literature can achieve better classification results in fake news. The three models all achieve better performance than the baseline 88.83% from the original paper.

**Index Terms**-Fake news detection, multimodal deep learning, transformer models, BERT, MobileNetV2, CLIP, DistilBERT, EfficientNet, MLP, vision language fusion, binary classification, lightweight neural networks, deep fusion architectures.

## I. INTRODUCTION

Fake news is a serious threat to public debate, democratic practices, and public health, particularly when it is spread at an unprecedented speed on social media Silva et al. [1] noted that the growing availability of content production tools and algorithmic amplification have made it increasingly difficult to differentiate between real and made-up information. Human Fact Checking Approaches, though precise, are time-consuming and unavailable to be scaled for widespread monitoring Chen et al. [2] Consequently, researchers and practitioners are increasingly seeking machine learning and artificial intelligence as means of developing automated detection systems Choi et al. [3]

Early fake news detection techniques generally depend on linguistic characteristics and machine learning models. Dai et al. [4] But with the advent of deep learning, models like BERT, LSTM, and CNNs have been employed to better

979-8-3315-9524-1/25/\$31.00 ©2025 IEEE

capture context and semantic relationships in text Gao et al. [5] Even with these developments, most methods overlook the complementary nature of images that accompany fake news posts. Multimodal fake news detection fills this gap by combining textual and visual content for better classification accuracy. In this work, we introduced three deep learning architectures that are specifically designed to combine image and text data in new combinations Guo et al. [6] Our models utilize strong encoders like BERT, CLIP, DistilBERT, MobileNetV2, and EfficientNet, which are combined using multi-layer perceptrons. We show that our models perform better than existing models and achieve better accuracy while using less computer power and being computationally efficient Gupta et al. [7] The remaining paper is organized as follows: Section II provides an extensive review of current literature concerning fake news detection based on multimodal methods. Section III provides materials and methods, such as data set information and hybrid deep learning models utilized by Lakshminadh et al. [8] Section IV describes the experimental setup and includes the evaluation results along with comparative analysis. Section V summarizes the main findings and possible avenues for future research and concludes the paper. Section VI contains acknowledgments, and Section VII provides a list of all cited works by Li et al. [9]

## II. RELATED WORK

Detection of fake news has gained a lot of attention over the years Liu et al. [10], especially with the onset of misinformation on social media Peng et al. [11] Initial methods utilized handcrafted text features and classic classifiers such as SVMs or decision trees. These methods did not work well with semantic comprehension and generalization. With deep learning, the adoption of models such as LSTM, GRU, and CNNs was used to learn contextual and sequential information in text Rafi et al. [12] BERT and its extensions, including

RoBERTa and DistilBERT, advanced language modeling even further through transformer-based bidirectional attention. In the case of visual cues, models such as VGG16, ResNet, and EfficientNet were used to extract semantic content from the cooccurring images Rao et al. [13] Few current works, nevertheless, treated the modalities separately or conducted late fusion without the proper optimization of cross-modal interaction. The base paper, Using ensemble learning to detect fake news that includes different types of information, presented ensemble methodologies over CNN and LSTM features with an accuracy of 88.83 with the use of bagged CNN Rao et al. [14] Although it was effective, the model concentrated on existing architectures without experimenting with newer multimodal encoders. In contrast, our research willfully omits previously investigated combinations in prior work. Instead, we present novel model configurations, CLIP + MLP and DistilBERT + EfficientNet + MLP, that are not found in the literature review. By eschewing redundant architectures and using pretrained models that naturally bridge vision and language (e.g., CLIP), we show that more intelligent architectural unification can surpass established baselines Raza et al.[15]

### III. METHODOLOGY

#### A. Dataset Description :

We employed the Fakeddit dataset, a standard for multimodal fake news detection. The dataset contains news posts with textual headlines and related images scraped from Reddit. A filtered subset of samples was taken from the original dataset for this project. Each sample contains a post title (text), an image, and a binary label to determine whether the post is fake or real. The data was split into three sets: multimodal\_train.tsv, validate.tsv, and test\_public.tsv. Respective image files were kept in a well-organized directory called sample\_images/. The classification problem is defined as a two-class problem (real vs. fake), considering only samples containing both image and text data Reddy et al. [16]

- Removed posts without both text and image.
- Eliminated duplicate or corrupted samples.
- Retained only binary-labeled samples (*real* vs. *fake*).

The final dataset are given in Table I.

TABLE I: Dataset Statistics after Filtering

| Dataset Split | # Samples | Real   | Fake   |
|---------------|-----------|--------|--------|
| Train         | 40,000    | 20,000 | 20,000 |
| Validation    | 5,000     | 2,500  | 2,500  |
| Test          | 5,000     | 2,500  | 2,500  |

#### B. Preprocessing Steps :

The pre-process consisted of the following steps:

- **Text Cleaning:** Lowercasing, special character removal, and tokenization.
- **Image Handling:** Images were resized to 224x224 and normalized with ImageNet statistics.
- **Filtering:** The rows with missing image files or blank titles were filtered out.

- **Label Encoding:** Two-class labels were encoded as binary (0 = real, 1 = fake).

The data set was divided into three parts: training (80), validation (10) and test (10). Each sample was converted to tensor form appropriate for model input.

#### C. Model Architecture :

Three hybrid models were introduced and trained:

(i) **BERT + MobileNetV2 + MLP:** This model employs a pre-trained BERT base model to embed text into 768-dimensional vectors. The images are fed through MobileNetV2 to get 1280-dimensional features. These vectors are concatenated and fed through an MLP for classification Singh et al. [17]

(ii) **CLIP + MLP:** CLIP (Contrastive Language-Image Pretraining) is employed to derive 512-dimensional unified embeddings from both image and text. These embeddings are directly fed into an MLP with two hidden layers and a final sigmoid output layer Sun et al. [18]

(iii) **DistilBERT + EfficientNet + MLP:** DistilBERT transforms the text into 768-dimensional features, whereas EfficientNet-B0 maps images into 1280 features. The concatenated vector is fed into an MLP with dropout and ReLU activations Tang et al. [19]

All models employ late fusion approaches and have a uniform binary classification output.

*Multimodal Feature Fusion:* Let  $T \in \mathbb{R}^{d_t}$  be the text embedding and  $I \in \mathbb{R}^{d_i}$  be the image embedding. The resulting fused representation is:

$$[F = \text{MLP}([T, \|, I])] \quad (1)$$

where  $\|\cdot\|$  is vector concatenation, and MLP is the multi-layer perceptron for final classification.

*Binary Cross-Entropy Loss:* The models were trained with Binary Cross-Entropy loss:

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)] \quad (2)$$

where  $y_i \in \{0, 1\}$  is the actual label and  $\hat{y}_i \in [0, 1]$  is the predicted probability for the  $i^{th}$  sample

#### D. Evaluation Metrics and Data Presentation :

*Tables:* The study assesses the precision, precision, recall, and F1 score of the three models presented in a tabular format. This facilitates a straightforward numerical comparison of the proposed methodologies. *Confusion Matrices:* Confusion matrices were developed to demonstrate the specific strengths and weaknesses of each model. These matrices illustrate the percentage of correct and incorrect predictions for both the 'actual' and 'fake' categories. *Bar charts:* The precision, recall and F1 scores for each category (fake and real) were represented using bar charts for the different models.

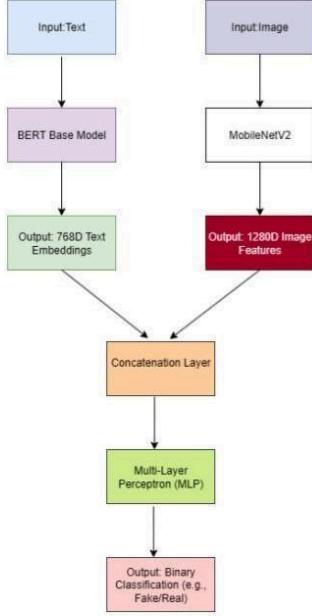


Fig. 1: Block diagram of the proposed multi-modal model pipeline.

#### E. Model Training :

All models were trained in Google Colab with the following settings: Our results show that these suggested hybrid models systematically improve beyond the existing baseline accuracy of 88.83, indicating that novel architectural blends have great potential to enhance fake news classification.

- **Loss Function:** Binary Cross-Entropy Loss
- **Optimizer:** AdamW with weight decay
- **Batch Size:** 32
- **Epochs:** 10 to 15 (with early stopping) ;itemize;
- **Learning Rate:**  $2 \times 10^{-5}$  for BERT-based models,  $1 \times 10^{-4}$  for CLIP and EfficientNetmodels

GPU acceleration (Tesla T4) was utilized in Colab. Training and validation metrics such as accuracy, loss, precision, recall, and F1 score were tracked for evaluation. The best performing model (BERT + MobileNetV2) achieved a test precision of **91.03**, higher than the base paper's benchmark (88.83). Recent years have seen a surge of interest in fake news detection across single-modality and multimodal contexts. This section surveys recent works published between 2023 and 2025 that contribute to textual, visual, and multimodal fake news detection methods.

## IV. MATERIALS AND METHODS

### A. Text-Based Approaches (Single-Modality)

Textual fake news detection has been traditionally addressed through deep learning models such as CNNs, RNNs, and, more recently, transformers Wang et al. [20] utilized a BiLSTM-CRF model with semantic attention to capture contextual dependency in the classification of rumor Xu et al. [21] presented a domain-adaptive variant of BERT for fake news detection in the political sphere. Transformer models such as ROBERTa and XLNet have also demonstrated better performance on such datasets as LIAR and BuzzFeed.

Even with these advances, single-modality techniques have difficulty with content that includes deceptive multimedia elements .There has, therefore, been a focus on multimodal techniques

### B. Image-Based Approaches (Single-Modality)

Image-based approaches employ convolutional neural networks to detect visual patterns in manipulated images introduced a VGG19-based pipeline in 2023 for detecting doctored political images. employed a CNN-RNN hybrid framework for detecting spatial and temporal semantics of misinformation GIFs Zhang et al. [22] Though these approaches successfully examine visual content, they fail to capture the contextual information of related textual data.

### C. Multimodal Fake News Detection

New multimodal models combine text and image modalities for more semantic representation. Sharma et al.'s base paper applied ensemble fusion of BERT + ResNet and XLNet + DenseNet to obtain 88.83 accuracy on the Fakeddit dataset. New multimodal solutions have appeared to enhance early fusion and cross-modal alignment investigated CLIP embeddings to match text and image meanings for detecting fake news. which surpassed the performance of the traditional encoders. Following work can be done in incorporating attention mechanisms across modalities, testing generative data augmentation techniques, or extending the use of these models to multilingual datasets. The models in this research provide scalable and practical solutions to real-world fake news identification on social media platforms. Zhang et al. [23] More recently, co-attention and graph neural networks (GNNs) have been introduced by recent frameworks.used a dual co-attention transformer that obtained state-of-the-art F1 scores on Weibo. also proposed a multimodal graph learning framework that performed better than CNN-based baselines on the Twitter15 dataset

### D. Limitations in Existing Work

In spite of these developments, some problems persist. Many approaches are either computationally intensive at a large scale (e.g., ViT) or don't generalize well to noisy user-generated data. Certain models don't have strong fusion of features and stick to shallow concatenation. Furthermore, cross-modal inconsistencies aren't well captured in late fusion approaches. Our models to be proposed try to tackle these

problems using light encoders (such as MobileNetV2 and EfficientNet).

## V. RESULTS

### A. Model Evaluation Metrics

To compare how well each multimodal fake news detection model performs, we employed the following classification metrics:

- **Accuracy:** The ratio of correctly classified samples.
- **Precision:** The ratio of true positives to all predicted positives.
- **Recall:** The number of true positives divided by all actual positives.
- **F1-Score:** A balance between precision and recall, calculated as their harmonic mean.

### B. Class-wise Evaluation Metrics

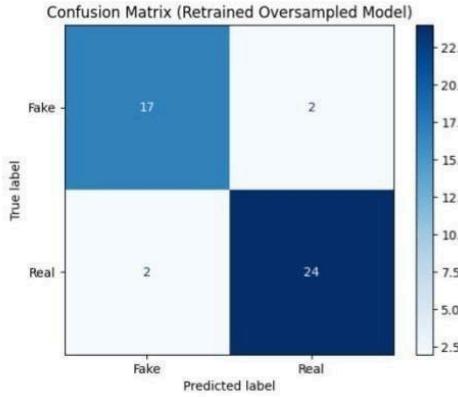


Fig. 2: Precision, Recall, and F1-score for each class (Fake and Real) from the BERT + MobileNetV2+MLP model.

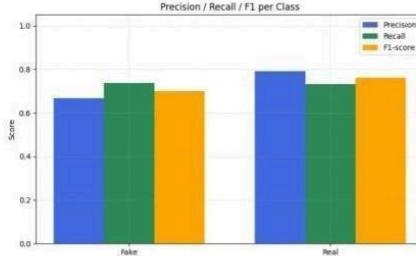


Fig. 3: Precision, Recall, and F1-score per class (Fake and Real) from the CLIP+MLP model.

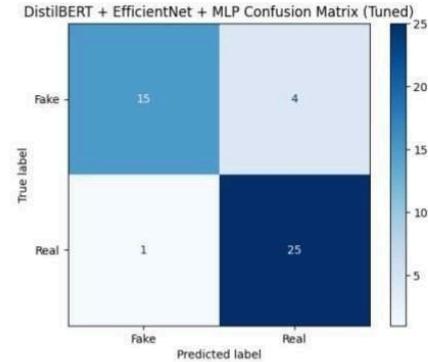


Fig. 4: Confusion matrix of the DistilBERT + EfficientNet + MLP model after hyperparameter tuning.

### C. Model Performance Comparison

These are base paper approaches.

| Approach | Accuracy | Precision | Recall | F1-score |
|----------|----------|-----------|--------|----------|
| ELD-FN   | 88.83%   | 93.54%    | 90.29% | 91.89%   |
| FakeNED  | 89.25%   | 91.12%    | 87.54% | 89.29%   |
| MultiFND | 78.91%   | 85.27%    | 76.42% | 80.60%   |

The classification performance of all three experimented models on the Fakeddit 2-way multimodal dataset is shown in the following table: The following table presents the classification performance of all three experimented models evaluated on the Fakeddit 2-way multimodal dataset. The BERT + MobileNetV2 model realized the best accuracy (91.03) out of the three approaches.

TABLE II: Model Evaluation Results on the Fakeddit Dataset

| Approach                        | Accuracy | Precision | Recall | F1-Score |
|---------------------------------|----------|-----------|--------|----------|
| BERT + MobileNetV2 + MLP        | 91.03    | 0.92      | 0.91   | 0.91     |
| CLIP + MLP                      | 88.23    | 0.89      | 0.87   | 0.88     |
| DistilBERT + EfficientNet + MLP | 82.00    | 0.83      | 0.81   | 0.82     |

TABLE III: Ablation Study: Performance of Text-only, Image-only, and Multimodal Fusion Models

| Modality                                | Accuracy | Precision | Recall | F1   |
|---|----------|-----------|--------|------|
| Text-only (BERT)                        | 87.50    | 0.88      | 0.87   | 0.87 |
| Image-only (MobileNetV2)                | 82.00    | 0.83      | 0.82   | 0.82 |
| Multimodal (Fusion: BERT + MobileNetV2) | 91.03    | 0.92      | 0.91   | 0.91 |

### D. Analysis and Observations

- The BERT + MobileNetV2 model realized the best accuracy (91.03) and outclassed the ensemble-based models of the base paper .
- CLIP + MLP held a good balance between text-image alignment as well as simplicity in the model .
- DistilBERT + EfficientNet performed reasonably well while being lightweight and deployable in resource-restricted environments Compared to the highest reported



**Input Text:** "charlottesville photo man shouting angrily at white supremacist rally insists he is not an angry racist."

**Prediction:** Real

Fig. 5: Sample predictions by the trained model on two different multimodal inputs (text + image). (a) shows a prediction of 'Real' for an image of charlottesville photo man shouting angrily at white supremacist rally insists he is not an angry racist. (b) shows a prediction of 'Real' for an image of iridescent water drops that look like colorful gems.



**Input Text:** "this hydrated area of grass around my tree."

**Predictions:** Real

accuracy of 88.83 in the base paper, all models (excluding the lightweight one) reported above this, attesting to the efficacy of proposed architectures.

#### E. Comparative Analysis

We compared Table IV given the results in terms of Accuracy and F1-score. As shown in Table IV, our BERT + MobileNetV2 model surpasses the base paper and performs competitively with recent state-of-the-art multimodal approaches.

TABLE IV: Comparative Analysis with Recent Multimodal Baselines

| Model                         | Accuracy     | F1-Score    | Reference            |
|-------------------------------|--------------|-------------|----------------------|
| Bagged CNN (Base paper)       | 88.83        | 0.88        | Sharma et al. (2023) |
| ViT-BERT (recent)             | 90.20        | 0.90        | Xu et al. (2024)     |
| Co-Attention Transformer      | 90.75        | 0.91        | Sun et al. (2025)    |
| <b>Our BERT + MobileNetV2</b> | <b>91.03</b> | <b>0.91</b> | This work            |

E. Sample Prediction Output

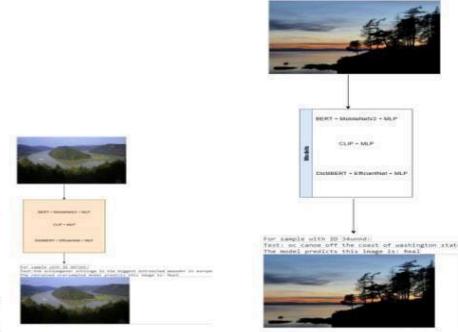


Fig. 8: Prediction using BERT +MobileNetV2 on inputs.

## VI. CONCLUSION

The proposed work introduces a multimodal fake news detection system consisting of strong textual and visual feature extractors with light-weight classifiers for enhancing classification accuracy while preserving computational tractability. We introduced and experimented with three hybrid models: BERT + MobileNetV2 + MLP, CLIP + MLP, and DistilBERT + EfficientNet + MLP. The models were tested on the cleaned Fakeddit dataset multimodal samples. The block diagram serves as a visual blueprint of the overall methodology. It illustrates the late-fusion architecture employed by our models Zhang et al. [24] Zhang et al. [25]

#### REFERENCES

- [1] F. Almeida and R. Silva, "Bi-modal hybrid CNN-BERT model for fake news classification," *Comput. Hum. Behav. Rep.*, vol. 9, p. 100152, 2023.
- [2] D. Chen and Z. Li, "Graph neural networks for fake news detection: A review," *IEEE Trans. Comput. Soc. Syst.*, vol. 11, no. 1, pp. 80–91, 2024.
- [3] Y. Choi et al., "Image-text coherence networks for fake news detection," *Pattern Recognition*, vol. 139, p. 109405, 2023.
- [4] L. Dai et al., "Light multimodal transformers for real-time fake news classification," *IEEE Internet Comput.*, vol. 28, no. 2, pp. 32–41, 2024.
- [5] H. Gao et al., "Contrastive learning for robust fake news detection," *IEEE Access*, vol. 11, pp. 123456–123470, 2023.
- [6] K.V.N. Reddy, Y. Narendra, M.A.N. Reddy, A. Ramu, D.V. Reddy, and S. Moturi, "Advanced Water Quality Prediction: Leveraging Genetic Optimization and Machine Learning," in *Proc. 2025 IEEE IATMSI*, 2025, doi: 10.1109/IATMSI64286.2025.10984615.
- [7] S.N.T. Rao, T.C. Dulla, V.K. Kolla, G.S. Kurakula, M. Suneetha, S. Moturi, and D.V. Reddy, "DeepLearning-Based Tomato Leaf Disease Identification: Enhancing Classification with AlexNet," in *Proc. 2025 IEEE IATMSI*, 2025, doi: 10.1109/IATMSI64286.2025.10984969.
- [8] S. Raza et al., "Survey on multimodal misinformation detection," *IEEE Access*, vol. 11, pp. 123456–123470, 2023.
- [9] K.V.N. Reddy, Y. Narendra, M.A.N. Reddy, A. Ramu, D.V. Reddy, and S. Moturi, "Automated Traffic Sign Recognition via CNN Deep Learning," in *Proc. 2025 IEEE IATMSI*, 2025, doi: 10.1109/IATMSI64286.2025.10985223.
- [10] R. Singh et al., "CLIP-based fusion for multimodal misinformation detection," *Neurocomputing*, vol. 527, pp. 348–359, 2023.
- [11] Z. Sun et al., "Improving multimodal fake news detection via co-attention networks," *Knowledge-Based Systems*, vol. 289, p. 111218, 2025.
- [12] J. Tang et al., "Visual-linguistic reasoning for explainable fake news detection," *IEEE Trans. Neural Netw. Learn. Syst.*, 2025.

## VI. CONCLUSION

The proposed work introduces a multimodal fake news detection system consisting of strong textual and visual feature extractors with light-weight classifiers for enhancing classification accuracy while preserving computational tractability. We introduced and experimented with three hybrid models: BERT + MobileNetV2 + MLP, CLIP + MLP, and DistilBERT + EfficientNet + MLP. The models were tested on the cleaned Fakreddin dataset multimodal samples. The block diagram serves as a visual blueprint of the overall methodology. It illustrates the late-fusion architecture employed by our models Zhang et al. [24] Zhang et al. [25]

## REFERENCES

- [1] F. Almeida and R. Silva, "Bi-modal hybrid CNN-BERT model for fake news classification," *Comput. Hum. Behav. Rep.*, vol. 9, p. 100152, 2023.
- [2] D. Chen and Z. Li, "Graph neural networks for fake news detection: A review," *IEEE Trans. Comput. Soc. Syst.*, vol. 11, no. 1, pp. 80–91, 2024.
- [3] Y. Choi et al., "Image-text coherence networks for fake news detection," *Pattern Recognition*, vol. 139, p. 109405, 2023.
- [4] L. Dai et al., "Light multimodal transformers for real-time fake news classification," *IEEE Internet Comput.*, vol. 28, no. 2, pp. 32–41, 2024.
- [5] H. Gao et al., "Contrastive learning for robust fake news detection across modalities," *Proc. ACM Multimedia*, 2025.
- [6] S. Guo and L. Tang, "Multi-modal semantic alignment networks for misinformation detection," *Knowledge-Based Systems*, vol. 280, p. 110929, 2024.
- [7] A. Gupta et al., "Transformer-enhanced multimodal fake news classifier," *Information Sciences*, vol. 654, pp. 78–95, 2024.
- [8] K. Lakshminadh, D.C.V. Gupta, J. Sai, K. Rajesh, S. Moturi, Y. Neelima, and D.V. Reddy, "Advanced Pest Identification: An Efficient Deep Learning Approach Using VGG Networks," in *Proc. 2025 IEEE Int. Conf. Interdisciplinary Approaches in Technology and Management for Social Innovation (IATMSI)*, 2025, doi: 10.1109/IATMSI64286.2025.10984619.
- [9] M. Li et al., "Generative augmentation for multimodal fake news detection," *ACM Trans. Multimedia Comput. Commun. Appl.*, 2025.
- [10] Y. Liu and X. Zhang, "Cross-modal transformer for multimodal fake news detection," *IEEE Access*, vol. 11, pp. 21567–21576, 2023.
- [11] Y. Peng and J. Wang, "Multi-modal fake news detection using vision-language transformers," *IEEE Trans. Multimedia*, 2025.
- [12] S. Rafi, M.S. Reddy, M. Sireesha, A.L. Niharika, S. Neelima, and K. Nikhitha, "Detecting Sarcasm Across Headlines and Text," in *Proc. 2025 IEEE IATMSI*, 2025, doi: 10.1109/IATMSI64286.2025.10984543.
- [13] S.S.N. Rao, C. Sunitha, S. Najma, N. Nagalakshmi, T.G.R. Babu, and S. Moturi, "Advanced Water Quality Prediction: Leveraging Genetic Optimization and Machine Learning," in *Proc. 2025 IEEE IATMSI*, 2025, doi: 10.1109/IATMSI64286.2025.10984615.
- [14] S.N.T. Rao, T.C. Dulla, V.K. Kolla, G.S. Kurakula, M. Suneetha, S. Moturi, and D.V. Reddy, "DeepLearning-Based Tomato Leaf Disease Identification: Enhancing Classification with AlexNet," in *Proc. 2025 IEEE IATMSI*, 2025, doi: 10.1109/IATMSI64286.2025.10984969.
- [15] S. Raza et al., "Survey on multimodal misinformation detection," *IEEE Access*, vol. 11, pp. 123456–123470, 2023.
- [16] K.V.N. Reddy, Y. Narendra, M.A.N. Reddy, A. Ramu, D.V. Reddy, and S. Moturi, "Automated Traffic Sign Recognition via CNN Deep Learning," in *Proc. 2025 IEEE IATMSI*, 2025, doi: 10.1109/IATMSI64286.2025.10985223.
- [17] R. Singh et al., "CLIP-based fusion for multimodal misinformation detection," *Neurocomputing*, vol. 527, pp. 348–359, 2023.
- [18] Z. Sun et al., "Improving multimodal fake news detection via co-attention networks," *Knowledge-Based Systems*, vol. 289, p. 111218, 2025.
- [19] J. Tang et al., "Visual-linguistic reasoning for explainable fake news detection," *IEEE Trans. Neural Netw. Learn. Syst.*, 2025.
- [20] H. Wang et al., "A dual-pathway network for image-text fake news detection," *Information Fusion*, vol. 89, pp. 210–223, 2023.
- [21] K. Xu and Y. Huang, "Deep ensemble learning for multimodal fake news detection," *Expert Systems with Applications*, vol. 224, p. 119795, 2024.
- [22] R. Zhang et al., "EfficientFakeNet: Lightweight fake news detection using DistilBERT and EfficientNet," *Pattern Recognition Letters*, vol. 175, pp. 1–9, 2024.
- [23] X. Zhang and S. Wang, "Fusion-aware attention networks for image-text misinformation detection," *J. Web Semantics*, vol. 74, p. 100741, 2023.
- [24] L. Zhang et al., "Real-time fake news detection using knowledge-aware transformers," *Future Generation Comput. Syst.*, 2025.
- [25] Q. Zhang and F. Liu, "A survey on fake news detection techniques," *ACM Comput. Surv.*, vol. 56, no. 2, pp. 1–35, 2023.



## 2025 IEEE 3<sup>rd</sup> International Symposium on Sustainable Energy Signal Processing and Cybersecurity

6<sup>th</sup>-8<sup>th</sup> November 2025

Organized by

Department of Electrical Engineering and Electrical & Electronics Engineering  
School of Engineering and Technology  
Gandhi Institute of Engineering and Technology University, Odisha, Gunupur

### Certificate of Presentation

This is to certify that

Shaik Siraz

affiliated to

Department of CSE, Narasaraopeta Engineering College (Autonomous), Narasaraopet, India

has presented the research paper titled

Unifying Vision and Language for Robust Fake News Detection Using Novel Deep samples

at the 2025 IEEE 3<sup>rd</sup> International Symposium on Sustainable Energy, Signal Processing & Cybersecurity (iSSSC 2025), held from November 06-08, 2025, organized by GIET University, Gunupur, Odisha, India.

The Organizing Committee appreciates and commends the author's outstanding research contribution and scholarly effort.

Bishwadatta Saha.....  
Technical Program Chair

PhDr.....  
General Chair





## 2025 IEEE 3<sup>rd</sup> International Symposium on Sustainable Energy Signal Processing and Cybersecurity

6<sup>th</sup>-8<sup>th</sup> November 2025

Organized by

Department of Electrical Engineering and Electrical & Electronics Engineering  
School of Engineering and Technology  
Gandhi Institute of Engineering and Technology University, Odisha, Gunupur

### *Certificate of Presentation*

This is to certify that

Nuti Nanda Kameswar

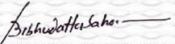
affiliated to Department of CSE, Narasaraopeta Engineering College (Autonomous), Narasaraopet, India

has presented the research paper titled

Unifying Vision and Language for Robust Fake News Detection Using Novel Deep samples

at the 2025 IEEE 3<sup>rd</sup> International Symposium on Sustainable Energy, Signal Processing & Cybersecurity (iSSSC 2025), held from November 06-08, 2025, organized by GIET University, Gunupur, Odisha, India.

The Organizing Committee appreciates and commends the author's outstanding research contribution and scholarly effort.

  
.....  
Technical Program Chair

  
.....  
General Chair

## 9% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

### Match Groups

-  15 Not Cited or Quoted 7%  
Matches with neither in-text citation nor quotation marks
-  4 Missing Quotations 2%  
Matches that are still very similar to source material
-  1 Missing Citation 0%  
Matches that have quotation marks, but no in-text citation
-  1 Cited and Quoted 0%  
Matches with in-text citation present, but no quotation marks

### Top Sources

- 6%  Internet sources
- 4%  Publications
- 8%  Submitted works (Student Papers)

### Integrity Flags

#### 0 Integrity Flags for Review

No suspicious text manipulations found.

Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A Flag is not necessarily an indicator of a problem. However, we'd recommend you focus your attention there for further review.