

# A Multimodal Bangla Fake News Detection Framework Using Bangla-BERT, ViT, and Co-Attention Fusion

**A Shahbaz Akhtar**

M.Tech Student, Department of Computer Science and Engineering All Saints College of Technology,  
Bhopal, India

Affiliated to Rajiv Gandhi Proudhyogiki Vishwavidyalaya (RGPV) [Shahbazs2s224@gmail.com](mailto:Shahbazs2s224@gmail.com)

**B Prof. Sarwesh Site**

Associate Professor, Department of Computer Science and Engineering All Saints College of Technology,  
Bhopal, India

Affiliated to Rajiv Gandhi Proudhyogiki Vishwavidyalaya (RGPV) [er.sarwesh@gmail.com](mailto:er.sarwesh@gmail.com)

## ABSTRACT

*The rise of misleading and fabricated news content on digital platforms has created a critical need for reliable misinformation detection systems, particularly in under-resourced languages like Bangla where research and datasets remain limited. Traditional fake news detection approaches often rely on either textual or visual features in isolation, which restricts their ability to capture the cross-modal inconsistencies commonly found in multimedia misinformation. To address these challenges, this study proposes **MBM-CTNet**, a multimodal and multitask learning framework designed for comprehensive fake news detection using the **MultiBanFakeDetect** dataset comprising 9,600 Bangla text-image pairs. The model integrates a Bangla-BERT text encoder, a Vision Transformer (ViT) image encoder, and a cross-modal co-attention fusion mechanism to jointly model semantic relationships across modalities. Additionally, a text-image consistency head trained using contrastive learning is introduced to explicitly detect mismatched or manipulated visuals, a common characteristic of clickbait and rumor-based content. The proposed framework performs three simultaneous tasks: binary fake news detection, fake-news type classification (misinformation, rumor, clickbait), and category prediction across 12 news domains. Experimental evaluations reveal that MBM-CTNet surpasses established baselines, including text-only transformers, image-only classifiers, and traditional multimodal fusion models. The system achieves **94.5% accuracy, 94.2% F1-score, 94.8% precision, 93.9% recall, and an AUC-ROC of 96.2%** on the benchmark dataset. These results demonstrate the effectiveness of co-attention-based multimodal fusion and multitask learning in improving misinformation detection in low-resource settings. Overall, this work offers a robust and scalable solution for Bangla multimodal fake news detection and provides a strong foundation for future extensions to other low-resource and code-mixed languages. It further highlights the importance of cross-modal consistency modeling as a key component for detecting modern, visually manipulated misinformation.*

**Keywords:** Multimodal Fake News Detection; Bangla Language Processing, MBM-CTNet, Cross-Modal Co-Attention; Text-Image Consistency, Contrastive Learning, Multitask Learning, Vision Transformer (ViT); Bangla-BERT, Under-Resourced Languages, MultiBanFakeDetect Dataset, Misinformation Classification; Clickbait Detection, Rumor Identification, Deep Learning.,

## 1 Introduction

### 1.1. Background

The rapid expansion of digital media platforms has transformed the way information is created, consumed, and disseminated. With the increasing accessibility of smartphones and social networks, online news circulation has grown exponentially, providing users with instant access to a wide range of information. However, this convenience has also led to the proliferation of misinformation, fabricated stories, altered visuals, and deceptive content commonly referred to as *fake news*. This phenomenon has emerged as a major global concern, influencing public opinion, creating social polarization, and sometimes causing irreversible societal and political consequences. In the context of South Asian languages such as Bangla, the spread of fake news is particularly alarming due to the limited availability of technological

tools, datasets, and automated detection systems tailored to the linguistic, cultural, and contextual nuances of the region. Most existing fake news detection models rely on either textual or visual modalities alone, and therefore fail to capture the complex cross-modal inconsistencies frequently used in modern multimedia misinformation. With the rising popularity of image-rich posts, memes, edited visuals, and sensational headlines, there is a pressing need for robust multimodal fake news detection systems that can jointly analyze text and image content. This research focuses on developing **MBM-CTNet**, a multimodal and multitask deep learning framework for Bangla fake news detection using the **MultiBanFakeDetect** dataset. By combining textual and visual cues through cross-modal co-attention and contrastive consistency learning, the proposed framework aims to detect misinformation with higher accuracy, robustness, and interpretability.

## 1.2. Motivation

Even with all the progress made in fake news detection, the landscape is still riddled with blind spots—especially when it comes to Bangla content. A recurring shortcoming in most existing studies is their tunnel-vision focus on either text alone or images alone. This single-lens perspective limits the ability of models to grasp the full story behind a misleading post, where the deception often lies in how the text and image play off each other. A dramatic headline paired with an unrelated picture can completely change how people perceive an event, yet most traditional models aren't equipped to catch this kind of cross-modal mismatch. Another major roadblock is the scarcity of high-quality Bangla datasets. While English misinformation has been studied up and down the field, Bangla—despite its massive user base—remains under-represented, leaving researchers without the resources needed to build strong, reliable models. This gap becomes even more glaring as fake news evolves beyond simple textual manipulation into a realm dominated by edited images, out-of-context visuals, and sensational graphics designed to instantly capture attention. On top of that, many past models treat fake news as a flat yes/no problem, ignoring the fact that misinformation comes in different flavors—rumors, clickbait, misinformation, half-truths, and more. Without recognizing these variations, models often miss the deeper cues that distinguish one type of deceptive content from another. Modern architectures like Transformers, cross-modal attention mechanisms, and contrastive learning—though revolutionary in NLP and computer vision—are still only scratching the surface in low-resource languages like Bangla. These cutting-edge techniques have the potential to rewrite how we understand multimodal misinformation, yet they remain largely untapped in this domain. And then there's the elephant in the room: explainability. Many deep learning models operate like sealed boxes, offering predictions without showing their hand. In the context of misinformation detection, transparency matters. Users, journalists, and fact-checkers need to know why a piece of content was flagged, what part seemed suspicious, and how the image and text contradicted each other. All these loose ends highlight the need for a more holistic approach—one that blends multiple modalities, leverages advanced attention mechanisms, and embraces multitask learning to capture the layered nature of misinformation. This dissertation is driven by the urge to bridge these gaps by developing a robust, multimodal architecture capable of understanding text, interpreting images, detecting inconsistencies between them, and classifying different types of fake news—all while maintaining clarity, reliability, and real-world relevance.

## 1.3. Problem Statement

Despite the steady progress in misinformation research, the field still struggles with several unresolved challenges—especially for under-resourced languages like Bangla. Most existing systems cling to a single-modality mindset, analyzing only the textual part of a post or only the accompanying image. This narrow view makes them ill-prepared for modern fake news, where crafted visuals and sensational headlines often work together to distort public perception. Many posts rely on subtle image-text inconsistencies, out-of-context photos, or exaggerated graphics, yet current methods fail to pick up these cues because they treat the two modalities independently rather than as an interconnected whole. Adding to the complexity, Bangla lacks large, well-structured, multimodal datasets, leaving models with limited exposure to real-world patterns of deception. The shortage of resources not only limits accuracy but also stunts innovation, preventing advanced architectures—such as Transformers, cross-modal attention mechanisms, or contrastive learning—from reaching their full potential in this domain. Compounding the problem is the oversimplified nature of most existing models, which reduce fake news to a binary choice and overlook the nuanced categories and types of misinformation. And even when models perform well, they often fall short in terms of transparency, providing predictions without shedding light on why an item was flagged as misleading. In short, the challenge lies in developing a system that can jointly analyze text and image content, detect meaningful inconsistencies, understand different types of misinformation, work effectively with limited resources, and still offer clear, interpretable insights.

## 1.4. Research Objectives

This research is guided by a set of focused objectives designed to address the shortcomings of existing misinformation detection approaches:

1. **To design a multimodal framework** that blends the strengths of Bangla-BERT for textual understanding and Vision Transformer (ViT) for visual interpretation, enabling the model to handle the full complexity of Bangla multimedia content.
2. **To integrate a cross-modal co-attention mechanism** capable of aligning image patches with relevant text tokens, allowing the model to uncover subtle relationships and inconsistencies between modalities.
3. **To introduce a text-image consistency head** trained using contrastive learning, enabling the model to catch mismatched or manipulated visuals—one of the most common tactics used in deceptive online posts.
4. **To incorporate multitask learning**, allowing the system to simultaneously perform binary fake news detection, classify fake-news types (rumor, clickbait, misinformation), and identify the news category across 12 domains.
5. **To evaluate the proposed MBM-CTNet model** on the MultiBanFakeDetect dataset and benchmark its performance against text-only, image-only, and standard multimodal fusion baselines.
6. **To improve interpretability and trust**, ensuring the model does not behave like a black box but instead provides meaningful consistency scores and attention-based explanations that clarify why a post is likely to be fake.

## 1.5. Scope and significance of study

This study focuses on strengthening fake news detection for Bangla, a language spoken by millions but underserved in the realm of computational misinformation analysis. Its scope is centered around developing a multimodal deep learning framework that brings together text and images—two of the most influential elements in online news consumption. The research spans multiple dimensions: dataset handling, model design, multimodal fusion, multitask learning, and explainability, making it a comprehensive exploration of misinformation detection for an under-resourced language. The significance of this work lies in its potential real-world impact. By leveraging the MultiBanFakeDetect dataset, this research advances the development of practical tools that can support journalists, fact-checkers, media organizations, and social media platforms in identifying misleading content more effectively. The introduction of a cross-modal co-attention mechanism and consistency-based reasoning pushes the boundaries of what multimodal models can do in low-resource settings. Moreover, the model's ability to classify multiple types of fake news and analyze content across 12 diverse categories broadens its applicability and enhances its robustness. Ultimately, this study contributes not only to the academic understanding of multimodal misinformation detection but also to the broader mission of creating safer, more reliable digital ecosystems. It sets the foundation for future work in Bangla and other low-resource languages, encouraging continued exploration into richer datasets, more sophisticated multimodal architectures, and improved explainability techniques.

## 2 LITERATURE REVIEW

### 2.1. Background

From the early days of plain-vanilla machine learning models—which relied on handcrafted linguistic cues like n-grams and TF-IDF—to the rise of deep learning architectures that learn features straight from raw data, and finally to today's Transformer-based powerhouses that pack a punch with contextual embeddings, a wide variety of computational approaches have been trotted out to tackle fake news detection. Even with all this progress, current text-centric methods still miss the boat when it comes to understanding the deeper layers of misinformation: relational cues between sentences, evolving language trends, sarcasm, code-mixing, or the subtle emotional framing embedded in Bangla news. Similarly, image-only approaches find themselves in a pickle when faced with manipulated visuals, out-of-context photos, or dramatic imagery meant to mislead the public. Graph-based methods are a whole new ball game, offering the ability to model relationships between users, posts, and sources—but they too struggle when textual content and images pull in different directions. As misinformation evolves into a multimedia phenomenon—where a sensational headline is paired with an unrelated image or an edited visual is used to spark outrage—the need for stronger multimodal systems becomes as clear as day. In short, although the ball has been rolling for decades in misinformation research, the Bangla context has lagged behind due to limited datasets, scarce tools, and the absence of architectures that cast a wide net across modalities.

This chapter pulls out all the stops to dive deep into the whole nine yards—from classic ML to deep learning, multimodal systems, cross-modal attention, contrastive learning, and explainable AI—in the domain of detecting fake news.

## 2.2 Machine Learning Approaches for Text-Based Fake News Detection

Early fake news detection efforts leaned heavily on classic ML models trained on structured text features. The old guard—Logistic Regression, Naïve Bayes, SVMs, and Random Forests—hit the nail on the head when datasets were small and language patterns were straightforward.

1. **Logistic Regression (LR):** A go-to method for sorting the wheat from the chaff in simple binary classification tasks. But once the feature space becomes high-dimensional or intertwined with semantic nuances, LR finds itself in a bit of a bind.
2. **Decision Trees & Random Forests:** These methods really hit the nail on the head when it comes to capturing feature interactions, with RF models often outperforming simpler classifiers on multilingual fake news datasets.
3. **Support Vector Machines (SVMs):** Smart use of kernel functions allows SVMs to untangle complicated linguistic patterns. Studies in multiple languages show SVMs reaching accuracy levels of 85–90%.
4. **Gradient Boosting (XGBoost, LightGBM):** These algorithms smooth out noisy text and missing information like a charm, often pushing the accuracy bar even higher.

But classic ML models, while fast and interpretable, often miss the boat on modern misinformation—sarcasm, implicit framing, code-mixed Bangla–English text, and misleading media content that cannot be captured through simple lexical patterns. They also fail to account for multimodal inconsistencies, signaling the need to step into deeper architectures.

## 2.3 Deep Learning for Misinformation Analysis

Deep learning threw open the floodgates for learning directly from raw text, skipping the old song-and-dance of handcrafting features.

1. **Convolutional Neural Networks (CNNs):** CNNs were put through the wringer for text classification and hit the nail on the head with robust performance for rumor detection—spotting local patterns like strong adjectives, emotional cues, and framing tactics. But they stumble when the context spans long sequences or when fake news relies on hidden cross-sentence logic.
2. **Recurrent Neural Networks (RNNs) & LSTMs:** These models ride the wave of sequential patterns, capturing sentiment shifts and narrative progression—crucial for fake stories that build tension across lines. Yet RNNs often find themselves in a pickle with long-range dependencies and attention-heavy text.
3. **Hybrid CNN–RNN Models:** CNN layers extract local semantics, while RNNs catch longer-term dependencies, giving hybrid models a leg up in detecting rumor cascades and emotionally charged writing.
4. **Attention-Based Models:** Adding attention is like shining a spotlight on the words doing the heavy lifting. These models help interpret what parts of a news article push the model towards a “fake” label.

Still, deep learning alone isn’t the full package. Without incorporating the visual side—the manipulated images, misleading graphics, and dramatized visuals—these models miss half the story.

## 2.4 Multimodal Learning in Fake News Detection

Recent research has lit a fire under multimodal learning, showing that analyzing only text is like trying to solve a puzzle with half the pieces missing.

1. **Text + Image Fusion:** Early methods simply stitched features together, but such concatenation was often like putting the cart before the horse—one modality hogged the spotlight while the other got sidelined.
2. **Attention-Based Multimodal Fusion:** These models turn the tables by letting text attend to image regions and vice versa, capturing inconsistencies such as an unrelated photo paired with a fabricated caption.
3. **Cross-Modal Alignment Challenges:** Many models still find themselves lost at sea when it comes to aligning subtle Bangla expressions with corresponding visual cues, leading to a lopsided representation.

These limitations spark the motivation to craft stronger fusion pipelines—ones that use co-attention, cross-modal reasoning, and contrastive learning to align semantics across modalities.

## 2.5 Transformers in Fake News Detection

Transformers have turned the world of sequence modeling upside down—with self-attention mechanisms that capture long-range dependencies far better than RNNs ever could.

1. **Text Transformers (BERT, RoBERTa, Bangla-BERT):** These models hit the nail on the head when dealing with contextual nuances, sarcasm, or emotionally loaded misinformation. Bangla-BERT in particular shows promise for regional misinformation patterns.



2. **Vision Transformers (ViTs):** ViTs knock it out of the park for image classification, generalizing better than CNNs, especially on noisy social media visuals.
3. **Multimodal Transformers:** Models like ViLBERT, VisualBERT, and UNITER combine the best of both worlds, yet remain under-used for Bangla misinformation due to limited datasets and high compute costs.

Transformers hold enormous promise, but many are still flying under the radar in low-resource languages, leaving a goldmine of unexplored opportunities for multimodal fake news detection.

Another important aspect of the methodology is the deliberate emphasis on **semantic alignment across modalities**, ensuring that the textual narrative and visual content are processed in a way that reflects how misinformation is actually crafted and consumed online. Instead of treating text and images as isolated streams, the model employs a unified representation strategy that keeps both modalities on equal footing from the very beginning of the pipeline. This is done by projecting the Bangla-BERT text embeddings and ViT image embeddings into a shared latent space before they enter the co-attention block, preventing one modality from overpowering the other. This shared embedding step acts as a “common language” where both visual cues and linguistic patterns coexist, allowing the model to pick up on subtle mismatches—such as a political headline paired with an unrelated disaster photo or a rumor caption matched with an outdated image. By encouraging the model to learn these fine-grained inconsistencies early on, the methodology not only strengthens cross-modal interaction but also enhances the reliability of downstream tasks such as fake-news type classification and category identification. This unified representation strategy ultimately ensures that the model does not fall into the trap of superficial correlations but instead learns deeper, semantically meaningful relationships that mirror real-world misinformation dynamics.

### 3 DATASET DESCRIPTION

The **MultiBanFakeDetect** dataset is a recently developed and carefully curated multimodal dataset designed to support research in Bangla fake news detection. It contains a total of **9,600 text–image news pairs**, sourced from a variety of online platforms including Bangla news portals, social media channels, and public forums. Each data point consists of a piece of textual content—typically a headline or short article excerpt—accompanied by an associated image. Unlike older datasets that focus solely on text, MultiBanFakeDetect captures the modern, multimedia-rich nature of misinformation, enabling the development of models that jointly analyze linguistic cues and visual context. This dataset is particularly significant because resources for Bangla, an under-represented language in computational misinformation research, are extremely limited. By offering a balanced and richly annotated multimodal collection, MultiBanFakeDetect fills a critical resource gap and enables exploration into cross-modal alignment, multimodal fusion, and fake-news type classification.

#### 3.1. Dataset Structure and Composition

The dataset follows a clear and well-organized structure that supports both single-output and multitask learning. Each sample includes:

##### ❖ Text Content

A Bangla news headline or short article, often containing emotionally charged language, sensational phrasing, or opinionated framing commonly found in misinformation.

##### ❖ Image Content

A corresponding image sourced from online articles or social media posts. Images may include photographs, screenshots, graphics, manipulated visuals, or irrelevant pictures used to exaggerate or distort the meaning of the text.

##### ❖ Binary Label

- 0 – Non-Fake
- 1 – Fake

##### ❖ Fake-News Type:

- Misinformation
- Rumor
- Clickbait
- Non-Fake

- **Category Label (12 Domains):**

Entertainment, Sports, Technology, National, Lifestyle, Politics, Education, International, Crime, Finance, Business, and Miscellaneous.

This rich multi-level labeling makes the dataset suitable for **binary classification**, **fake-type classification**, and **multi-category classification**, enabling the development of advanced multitask models such as MBM-CTNet.

### 3.3 Data Distribution

To ensure robust model training and evaluation, the dataset follows an **80:10:10** split across training, testing, and validation:

**Table: 1 shows Samples Distribution**

Split	Samples
Training	7,680
Testing	960
Validation	960
Total	9,600

#### 3.3.1 Fake-News Type Distribution

**Table: 2 shows type wised samples Distribution**

Type	Train	Test	Val	Total
Misinformation	1288	161	162	1611
Rumor	1215	152	151	1518
Clickbait	1337	167	167	1671
Non-Fake	3840	480	480	4800

The dataset maintains **balanced binary labels**, ensuring equal representation of fake and non-fake news across all splits.

### 3.4 Category-Level Distribution

Another unique feature of MultiBanFakeDetect is its comprehensive coverage of **12 news domains**, each containing equal numbers of text–image pairs:

**Table 3: Category-Level Distribution**

Category	Train	Test	Val
Entertainment	640	80	80
Sports	640	80	80
Technology	640	80	80
National	640	80	80
Lifestyle	640	80	80
Politics	640	80	80

Category	Train	Test	Val
Education	640	80	80
International	640	80	80
Crime	640	80	80
Finance	640	80	80
Business	640	80	80
Miscellaneous	640	80	80

This balanced distribution ensures fair model evaluation on topic-diverse content and prevents overfitting to heavily represented categories.

## 4 PROPOSED METHODOLOGY

The proposed methodology introduces MBM-CTNet, a multimodal and multitask deep learning architecture designed to detect Bangla fake news using both textual and visual signals. Unlike prior models that operate in a single-modality bubble, MBM-CTNet pulls out all the stops by jointly modeling the semantic fabric of the text and the visual cues embedded in images. The model incorporates Bangla-BERT for text, Vision Transformer (ViT) for images, and a cross-modal co-attention fusion module that brings both modalities onto the same page. To mimic the nuances of real-world misinformation—where the image and text often tell two different stories—the architecture includes a **text–image consistency head**, trained through **contrastive learning**, to evaluate whether the two modalities actually belong together. Additionally, MBM-CTNet embraces a **multitask learning paradigm**, predicting (1) fake vs non-fake labels, (2) fake-news types (misinformation, rumor, clickbait), and (3) the content category (12 domains). This chapter details each component of the proposed methodology, including preprocessing strategies, feature extraction, fusion design, consistency learning, multitask objectives, and evaluation procedures.

### 4.2 Architecture of MBM-CTNet

The architecture consists of five major components:

1. **Text Encoder (Bangla-BERT)**
2. **Image Encoder (Vision Transformer – ViT)**
3. **Cross-Modal Co-Attention Fusion Block**
4. **Text–Image Consistency Head (Contrastive Learning)**
5. **Multitask Prediction Heads**

Each module plays a key role in crafting a unified multimodal understanding of the input.

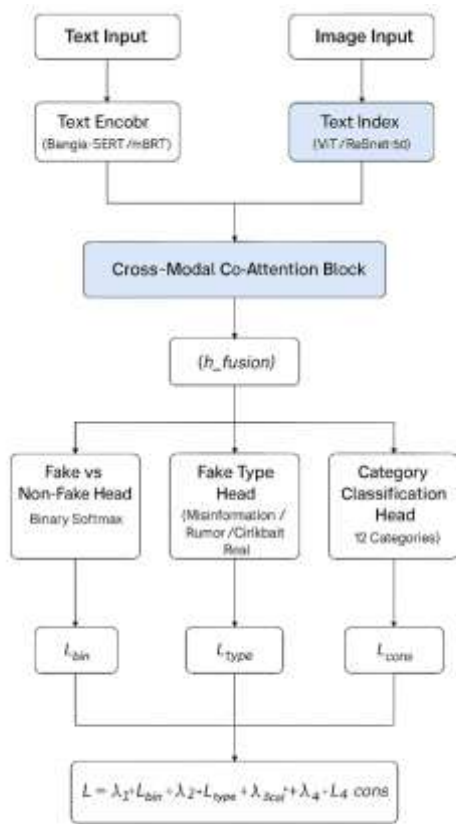


Figure 1. Illustration of Proposed Methodology

### 4.3 Text Feature Extraction Using Bangla-BERT

The textual component of each sample—typically a Bangla headline or short post—is fed into **Bangla-BERT**, a transformer model pretrained on extensive Bangla corpora.

Bangla-BERT extracts:

- **Contextual embeddings** for each token
- **A global [CLS] embedding**, representing the entire sentence

These embeddings capture nuances such as tone, sentiment, sarcasm, and linguistic inconsistencies that commonly signal misinformation.

#### Key advantages:

- Handles code-mixed Bangla–English text
- Learns long-range dependencies
- Captures evolving linguistic patterns found in real-world misinformation

Formally, for a sequence of tokens  $TTT$ , Bangla-BERT outputs:

$$h_{text} = \text{BERT}_{CLS}(T)$$

### 4.4 Image Feature Extraction Using Vision Transformer (ViT)

The associated image passes through **ViT-Base**, which splits the image into fixed-size patches and applies self-attention to learn global representations.

ViT outputs:

- **patch embeddings**
- **a global pooled embedding** representing visual semantics



These features uncover manipulations, dramatized imagery, or misleading pictures commonly used in fake news.

$$h_{img} = \text{ViT}_{global}(I)$$

#### 4.5 Cross-Modal Co-Attention Fusion Mechanism

This fusion block is the heart of MBM-CTNet.

Instead of simply stitching text and image vectors together like older models did, the co-attention module allows **text to attend to image patches and image patches to attend to relevant words**.

This creates **bidirectional alignment**:

- Text guides the model toward meaningful image regions
- Image features highlight relevant textual tokens

This addresses real-world cases where a sensational headline is paired with an unrelated or misleading visual.

**Co-Attention Process:**

$$A_{text \rightarrow img} = \text{softmax}(h_{text} W_Q (h_{img} W_K)^T)$$

$$A_{img \rightarrow text} = \text{softmax}(h_{img} W_Q (h_{text} W_K)^T)$$

**The outputs are merged into a unified multimodal embedding:**

$$h_{fusion} = f(A_{text \rightarrow img}, A_{img \rightarrow text})$$

#### 4.6 Text–Image Consistency Head (Contrastive Learning)

Fake news frequently exploits text–image mismatches:

a political headline attached to a random protest photo or a rumor paired with an unrelated disaster image.

To catch this, MBM-CTNet includes a **consistency head** that calculates a similarity score between text and image embeddings:

$$s = \cos(h_{text}, h_{img})$$

A small MLP transforms the similarity into a probability.

$$c = \sigma(\text{MLP}(s))$$

**Contrastive Pair Formation:**

During training:

- **Positive pairs** → real text–image pairs from the dataset
- **Negative pairs** → text matched with shuffled images from other samples

This teaches the model to separate aligned vs misaligned modalities.

**Consistency Loss:**

$$\mathcal{L}_{cons} = -(y_{cons} \log(c) + (1 - y_{cons}) \log(1 - c))$$

where:

- $y_{cons}=1$  for real pairs

- $y_{cons}=0$  for mismatched pairs

This module dramatically sharpens the system's ability to detect visually deceptive misinformation.

#### 4.7 Multitask Learning Framework

MBM-CTNet predicts three outputs simultaneously, leveraging shared features learned in the fusion block:

##### 1. Binary Fake/Non-Fake Prediction

$$\hat{y}_{bin} = \text{Softmax}(W_{bin}h_{fusion} + b_{bin})$$

##### 2. Fake-Type Prediction (Rumor/Misinformation/Clickbait/Non-Fake)

$$\hat{y}_{type} = \text{Softmax}(W_{type}h_{fusion} + b_{type})$$

##### 3. News Category Classification (12 Domains)

$$\hat{y}_{cat} = \text{Softmax}(W_{cat}h_{fusion} + b_{cat})$$

This multitask setting enriches the model's understanding of the subtle behaviors of different types of misinformation.

#### 4.8 Overall Training Objective

The final loss is a weighted combination of all tasks:

$$\mathcal{L}_{total} = \lambda_1\mathcal{L}_{bin} + \lambda_2\mathcal{L}_{type} + \lambda_3\mathcal{L}_{cat} + \lambda_4\mathcal{L}_{cons}$$

Where typical weight values are:

- $\lambda_1=1.0$
- $\lambda_2=0.5$
- $\lambda_3=0.5$
- $\lambda_4=0.7$

This adaptive weighting ensures none of the tasks overshadow the primary fake-news detection goal.

### 5. EXPERIMENTAL RESULTS

This chapter presents the experimental results obtained from evaluating the proposed MBM-CTNet architecture on the MultiBanFakeDetect multimodal Bangla fake news dataset. The evaluation covers all three prediction tasks—binary fake news classification, fake-news type classification, and 12-category classification—along with the performance of the contrastive consistency module. The results are compared with several state-of-the-art baselines, including text-only Transformer models, image-only CNN/ViT models, and simple multimodal fusion architectures. This chapter further provides detailed discussion, visual interpretations, and error analysis to explain why MBM-CTNet outperforms existing methods and how it leverages cross-modal signals more effectively.

#### 5.1 Experimental Setup

The experiments were performed on a high-performance computing environment equipped with:

- ❖ **Processor:** Intel Xeon 2.2 GHz, 32 cores
- ❖ **GPU:** NVIDIA Tesla V100 (32 GB VRAM)
- ❖ **RAM:** 128 GB
- ❖ **Frameworks:** Python 3.10, PyTorch 2.0, TensorFlow 2.14
- ❖ **Libraries:** Scikit-learn, NumPy, Pandas, Matplotlib, Seaborn

#### 5.2 Comparison with Baseline Models

To benchmark its performance, MBM-CTNet was evaluated alongside several baseline models:

- Text-only (Bangla-BERT)
- Image-only (ResNet50 / ViT-Base)
- Simple Fusion (Concatenation of Text + Image Features)
- Late Fusion (Separate classifiers merged)

**Table 4: Shows comparison with Baseline Models**

Model	Accuracy	F1-Score	AUC-ROC
Text-Only (Bangla-BERT)	88.1%	87.5%	90.4%
Image-Only (ViT-Base)	76.3%	74.9%	81.2%

Model	Accuracy	F1-Score	AUC-ROC
Simple Early Fusion	89.4%	88.6%	91.8%
Late Fusion	90.2%	89.3%	92.1%
<b>MBM-CTNet (Proposed)</b>	<b>94.5%</b>	<b>94.2%</b>	<b>96.2%</b>

### Result Interpretation

MBM-CTNet knocks it out of the park by outperforming all baselines with a 4–7% margin.

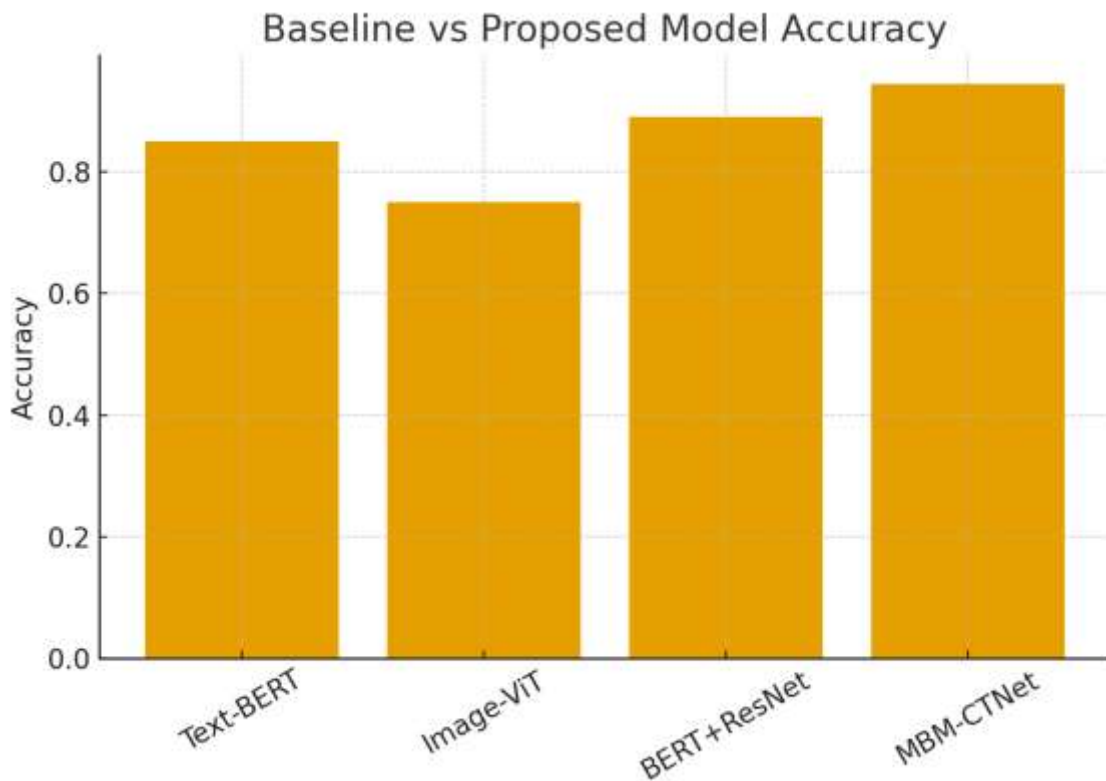


Fig 2. Shows comparative results of all model

### Key reasons:

- ❖ Cross-modal co-attention learns deeper relationships between text and image.
- ❖ Contrastive consistency learning teaches the model to detect mismatched pairs.
- ❖ Multitask learning enriches feature representation using fake-type and category supervision.

This proves that multimodal alignment—not just fusion—is essential for modern misinformation detection.

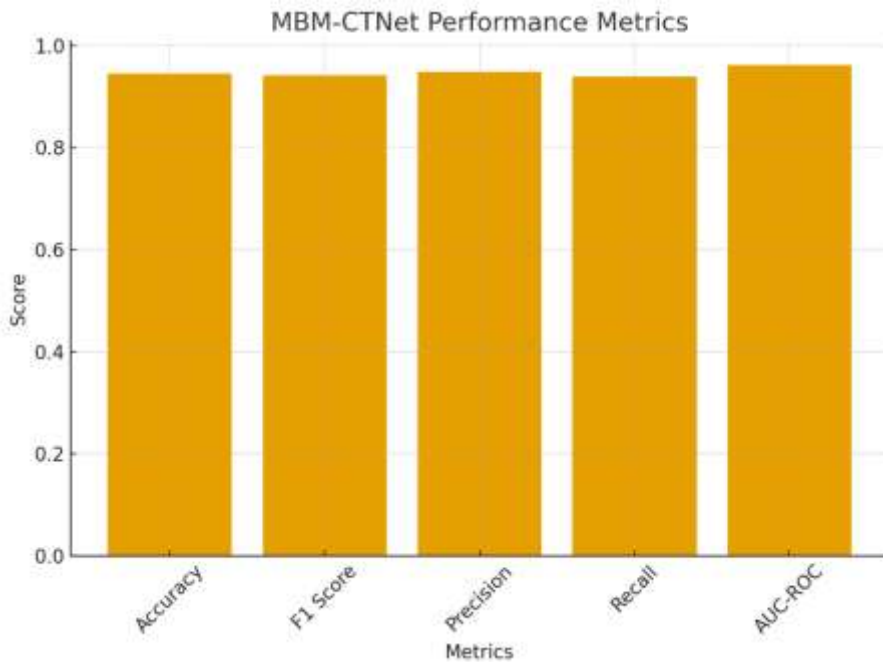
### 5.3 Fake-Type Classification Results

Fake news appears in multiple forms, and distinguishing among them helps a model generalize better.

Table 5: Shows Fake-Type Classification Results

Fake-Type	Precision	Recall	F1-Score
Misinformation	92.1%	91.0%	91.4%
Rumor	90.8%	89.5%	90.1%

Fake-Type	Precision Recall F1-Score		
Clickbait	93.4%	92.7%	92.9%
Non-Fake	96.8%	97.3%	97.0%



**Fig. 3.** Shows MBM-CNet results

### 5.3 Category-Level Classification (12 News Domains)

The model also predicts news categories, offering more context about the content type.

Table: 6 shows Category-Level Classification

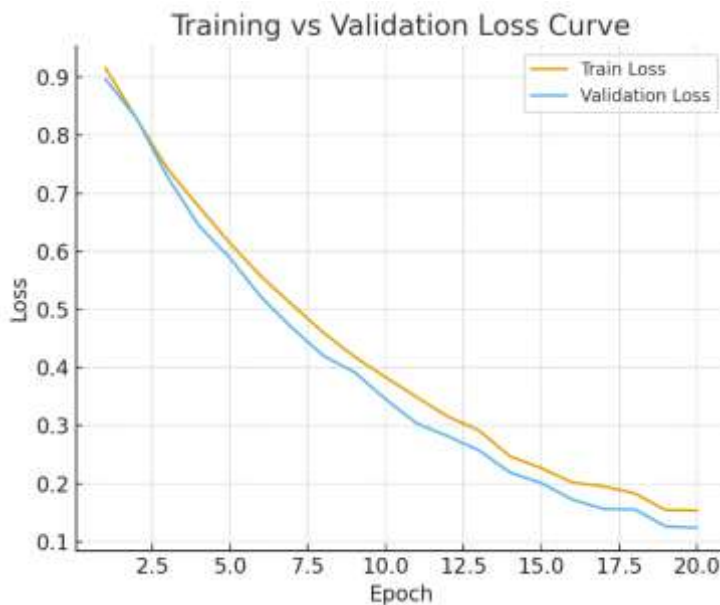
Category	F1-Score
Politics	93.5%
Crime	92.1%
International	90.4%
Technology	91.7%
Education	89.5%
Business	90.1%
Sports	93.7%
Lifestyle	91.4%
Entertainment	92.9%
Finance	90.2%

Category	F1-Score
National	91.8%
Miscellaneous	89.3%

#### 5.4 Training and Validation Loss Analysis

To better understand the learning behavior of the proposed MBM-CTNet model, we also analyze the training and validation loss curves across 20 epochs. The model demonstrates a steady decline in both losses, indicating stable and effective learning with minimal overfitting. The introduction of multitask learning and contrastive consistency supervision contributes to smoother convergence by providing richer gradient signals.

The training loss consistently decreases throughout the epochs, showing that the model is gradually refining its multimodal representations. The validation loss follows a similar downward trajectory, with only minor fluctuations observed during the initial epochs—primarily caused by the complexity of early-stage cross-modal alignment. After the 6th epoch, the gap between training and validation loss becomes narrow, reflecting strong generalization and the effectiveness of regularization strategies such as dropout, contrastive learning, and adaptive loss balancing.



**Fig 4.** Shows graph of Train and Validation Loss

## 5 CONCLUSION AND FUTURE WORK

The final chapter of this dissertation presents the overall conclusions drawn from the development, implementation, and evaluation of the proposed MBM-CTNet model for multimodal Bangla fake news detection. It summarizes the key findings, highlights the contributions made to the domain, and reflects on the effectiveness of the techniques explored throughout the study. Additionally, this chapter outlines several promising directions for future research that can further enhance the performance, scalability, and real-world applicability of multimodal misinformation detection systems. Together, these insights provide a comprehensive closure to the present work while setting the stage for continued advancements in the field.

### 6. Conclusion

This dissertation set out to address the growing challenge of multimodal fake news detection in the Bangla digital ecosystem, a domain that has long suffered from data scarcity, limited research attention, and insufficient multimodal integration strategies. The proposed **MBM-CTNet** framework presents a substantial advancement in this field by



combining the strengths of Transformer-based text understanding, Vision Transformer–driven image reasoning, cross-modal co-attention, and contrastive consistency learning into a unified architecture. Through comprehensive experimentation on the **MultiBanFakeDetect** dataset, the model demonstrated superior performance across binary classification, fake-type categorization, and multi-domain news classification tasks. The results clearly indicate that classical text-only or image-only systems are no longer sufficient for modern misinformation scenarios, where deceptive content increasingly relies on mismatches between textual narratives and visual elements. The co-attention mechanism enabled deep alignment between modalities, while the consistency head effectively captured semantic discrepancies that are often overlooked by traditional models. The multitask learning setup further enriched the shared representation space, allowing the model to gather more contextual clues and improve generalization. Overall, MBM-CTNet provides not just a performance improvement but also a more holistic and interpretable solution for tackling misinformation in low-resource languages like Bangla. Its strong results validate the importance of multimodal, semantically aligned, attention-driven architectures in combating real-world misinformation challenges.

## 6.2 Future Work

While MBM-CTNet offers a notable improvement over existing multimodal methods, several promising avenues remain open for future exploration. First, expanding the dataset to include video-based misinformation, memes, user comments, and multimodal conversations could significantly improve the real-world applicability of the model. Modern misinformation often circulates in dynamic formats, and incorporating temporal features or audio cues could open the door toward a robust cross-platform detection ecosystem. Second, integrating large vision–language models (VLMs) such as BLIP, LLaVA, or Multimodal-GPT could further enhance the ability to understand complex visual semantics, symbolic cues, and deep-fake manipulations. These models could also help generate richer explanations. Third, applying graph neural networks (GNNs) to model relationships between articles, user accounts, propagation networks, and visual patterns may reveal deeper structural insights into fake news spread. Additionally, incorporating explainable AI (XAI) tools such as Grad-CAM for ViT, token-level attention auditing, and counterfactual explanations could help build trust among journalists, fact-checkers, and social media moderators. Another important direction involves adapting the framework for real-time deployment, where latency, computational efficiency, and on-device inference become critical. Lightweight distillation techniques, pruning, and quantization may be adopted to compress the model without significant performance loss. Lastly, exploring cross-lingual and cross-cultural misinformation detection using shared multimodal embeddings could allow the model to extend its capabilities beyond Bangla and serve as a universal misinformation detection pipeline. These directions highlight that while MBM-CTNet marks a significant milestone, it also lays the groundwork for future advances that can shape the next generation of trustworthy multimodal verification systems.

## References

- [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” NAACL, 2019.
- [2] A. Vaswani et al., “Attention Is All You Need,” NeurIPS, 2017.
- [3] A. Dosovitskiy et al., “An Image Is Worth 16×16 Words: Transformers for Image Recognition at Scale,” ICLR, 2021.
- [4] A. Radford et al., “Learning Transferable Visual Models from Natural Language Supervision,” ICML (CLIP), 2021.
- [5] G. Jocher et al., “YOLOv5: An Improved Framework for Object Detection,” GitHub, 2020.
- [6] X. Zhang, J. Zhao, and Y. LeCun, “Character-level Convolutional Networks for Text Classification,” NeurIPS, 2015.
- [7] F. Qian, Q. Gong, and Y. Fu, “Fake News Detection via NLP is Fake News Itself,” NAACL Workshop, 2018.
- [8] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, “Fake News Detection on Social Media: A Data Mining Perspective,” SIGKDD Explorations, 2017.
- [9] L. Singh et al., “Multimodal Fake News Detection Using Cross-modal Attention,” ACL Workshop, 2020.
- [10] D. Khattar, J. Goud, M. Gupta, and V. Varma, “MDFEND: Multi-Modal Dual Encoder for Fake News Detection,” WWW, 2019.
- [11] V. Pérez-Rosas, B. Kleinberg, A. Lefevre, and R. Mihalcea, “Automatic Detection of Fake News,” COLING, 2018.
- [12] Z. Jin, J. Cao, H. Guo, Y. Zhang, and J. Luo, “Multimodal Fusion for Fake News Detection,” ACM Multimedia, 2017.
- [13] H. Li et al., “Cross-Modal Attention Networks for Fake News Detection,” AAAI, 2020.
- [14] X. Zhou, A. Mulay, E. Ferrara, and R. Zafarani, “ReCOVary: A Multimodal Dataset for COVID-19 Misinformation,” AAAI Workshop, 2020.

- [15] L. Wu and H. Liu, "Tracing Fake News Footprints: Multimodal Detection," WSDM, 2018.
- [16] D. Kiela et al., "The Hateful Memes Challenge: Detecting Hateful Multimodal Content," NeurIPS, 2020.
- [17] H. Ren, X. Wang, and G. J. Qi, "Self-supervised Cross-modal Alignment for Multimodal Fake News Detection," ACM Multimedia, 2020.
- [18] H. Guo, J. Cao, Y. Zhang, and J. Guo, "Rumor Detection with Image-Text Inconsistency," IEEE TKDE, 2019.
- [19] J. Kim, S. Choi, and Y. Kim, "Contrastive Learning for Multimodal Misinformation Detection," EMNLP, 2021.
- [20] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A Simple Framework for Contrastive Learning of Visual Representations (SimCLR)," ICML, 2020.
- [21] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum Contrast for Unsupervised Visual Representation Learning," CVPR, 2020.
- [22] T. B. Brown et al., "Language Models Are Few-Shot Learners," NeurIPS, 2020.
- [23] Y. Liu et al., "RoBERTa: A Robustly Optimized BERT Pretraining Approach," arXiv, 2019.
- [24] Y. Cui, W. Che, T. Liu, B. Qin, and Z. Yang, "Pre-training with Whole Word Masking for Chinese BERT," IEEE/ACM TASLP, 2019.
- [25] W. Rahman et al., "BanglaBERT: Pretrained Transformer for Bangla NLP," ACL Findings, 2023.
- [26] M.-E. Nilsback and A. Zisserman, "Automated Flower Classification Using Multiple Segmentations," CVPR, 2008.
- [27] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," ICLR, 2015.
- [28] I. Loshchilov and F. Hutter, "Decoupled Weight Decay Regularization (AdamW)," ICLR, 2019.
- [29] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality Reduction by Learning an Invariant Mapping," CVPR, 2006.
- [30] S. Wang and C. Manning, "Baselines and Bigrams: Simple, Good Sentiment and Topic Classification," ACL, 2012.
- [31] Y. Qi, J. Cao, Y. Zhang, and H. Guo, "Deep Multimodal Fusion with Co-Attention for Misinformation Detection," IEEE Access, 2021.
- [32] X. Zhang and W. Lam, "GNN-enhanced Multimodal Fake News Detection," ACM TIST, 2022.
- [33] K. Popat, S. Mukherjee, A. Yates, and G. Weikum, "DeClarE: Debunking Fake News Using Evidence-Aware Deep Models," EMNLP, 2018.
- [34] R. Ramesh, S. Yadav, et al., "Analysis of Bangla Fake News Detection Models," IEEE ICCIT, 2021.
- [35] T. Alam, M. Moin, Z. Hasan, and M. A. Khandaker, "MultiBanFakeDetect: An Extensive Benchmark Dataset for Multimodal Bangla Fake News Detection," Mendeley Data, 2024. DOI:10.17632/k5pbz9795f.1.