

BanglaMM-Disaster: A Multimodal Transformer-Based Deep Learning Framework for Multiclass Disaster Classification in Bangla

Ariful Islam

*Department of Computer Science and Engineering
Chittagong University of Engineering and Technology
Pahartoli, Raozan-4349, Chittagong, Bangladesh
arifulislamnayem11@gmail.com*

Md Rifat Hossen

*Department of Computer Science and Engineering
Chittagong University of Engineering and Technology
Pahartoli, Raozan-4349, Chittagong, Bangladesh
rifat8851@gmail.com*

Md. Mahmudul Arif

*Department of Electronics and Telecommunication Engineering
Chittagong University of Engineering and Technology
Pahartoli, Raozan-4349, Chittagong, Bangladesh
mdmahmudularif896@gmail.com*

Abdullah Al Noman

*Wilmington University
320 N Dupont Hwy, New Castle, DE 19720
anoman001@my.wilmu.edu*

Md Arifur Rahman

*College of Graduate and Professional Studies
Trine University
1 University Ave, Angola, IN 46703
Mrahman22@my.trine.edu*

Abstract—Natural disasters remain a major challenge for Bangladesh, so real-time monitoring and quick response systems are essential. In this study, we present BanglaMM-Disaster, an end-to-end deep learning-based multimodal framework for disaster classification in Bangla, using both textual and visual data from social media. We constructed a new dataset of 5,037 Bangla social media posts, each consisting of a caption and a corresponding image, annotated into one of nine disaster-related categories. The proposed model integrates transformer-based text encoders, including BanglaBERT, mBERT, and XLM-RoBERTa, with CNN backbones such as ResNet50, DenseNet169, and MobileNetV2, to process the two modalities. Using early fusion, the best model achieves 83.76% accuracy. This surpasses the best text-only baseline by 3.84% and the image-only baseline by 16.91%. Our analysis also shows reduced misclassification across all classes, with noticeable improvements for ambiguous examples. This work fills a key gap in Bangla multimodal disaster analysis and demonstrates the benefits of combining multiple data types for real-time disaster response in low-resource settings.

Index Terms—multimodal deep learning, transformer models, Bangla language, disaster classification, BanglaBERT, ResNet50, early fusion, low-resource languages, XLM-RoBERTa

IEEE Copyright Notice: © 2025 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Publication: Accepted for publication in IEEE SPICSCON 2025.

I. INTRODUCTION

Bangladesh frequently experiences natural hazards such as floods, cyclones, landslides, and fires with major effects on

communities and economic stability. The increased prevalence and the intensity of such events underscore the pressing necessity for effective monitoring systems. Social networking sites have turned into valuable sources of live disaster news with citizens offering text-based descriptions of incidents and posting them frequently along with images that provide instant context [1]. This explosion of multimodal data offers a unique opportunity for automated disaster monitoring but extracting actionable insights from noisy, multilingual content is still a major technical challenge especially for low-resource languages such as Bangla [2].

Traditional disaster monitoring systems rely on satellite imagery and official reports, which are vulnerable to delays and limited coverage. On the contrary, social media provides real-time, ground-level footage that supplements traditional data sources. Despite these promising developments, processing under-resourced languages like Bangla remains challenging due to the language's inherent complexity and the widespread use of informal expressions.

Most of the research in disaster analytics so far has focused on unimodal techniques, processing text or image data independently. These techniques have reasonable performance for resourced languages, they do not capture the complementary character of multimodal information. Recent advances in deep learning have facilitated sophisticated multimodal fusion methods [3]. These methods are significantly under-explored for Bangla, and annotated multimodal datasets are scarce. Merging multimodal data has particular advantages in classifying



Fig. 1: Multimodal disaster content from social media.

disasters, as text provides clear explanations of events and outcomes while visual information offers immediate evidence of damage and context.

This paper introduces BanglaMM-Disaster, a comprehensive multimodal Bangla disaster classification framework. Our approach systematically incorporates transformer-based text encoders with convolutional neural networks for image analysis. The main contributions are:

- We create a large-scale multimodal disaster dataset of 5,037 annotated social media posts covering nine disaster categories.
- We propose a transformer–CNN fusion framework that achieves 83.76% accuracy and outperforms unimodal baselines.
- We conduct detailed error analysis, showing consistent cross-modal benefits to enable more reliable disaster understanding.

II. LITERATURE REVIEW

Automatic disaster analysis has moved from conventional rule-based systems to cutting-edge deep learning approaches. Early research focused on single-modality analysis, with separate streams for textual and visual content independently. This section reviews current methods, referring to trends for multimodal systems and identification of critical gaps in low-resource language processing. **Unimodal Approaches for Disaster Analytics:** Text-based disaster analysis has been extensively studied, particularly for high-resource languages. Caragea et al. [1] were the first to utilize CNN for differentiating between informational disaster messages, while transformer models like BERT, Devlin et al. [4] demonstrated better performance in contextual semantic comprehension. Recent work by Han et al. [5] introduced QuakeBERT, achieving 84.33% F1-score for earthquake disaster tweet classification, demonstrating the effectiveness of domain-specific transformer models. For image-based analysis, Kaur et al. [6] proved the effectiveness of transfer learning in hurricane damage

detection, and Mouzannar et al. [7] investigated crowdsourced photos for damage assessment. **Multimodal Fusion Techniques:** Recent multimodal disaster research has achieved significant breakthroughs through sophisticated fusion strategies. El-Niss et al. [8] looked into federated learning approaches with F1-score 85.2%, and Khattar et al. [9] demonstrated cross-attention mechanisms with 84.08% F1-score. Dar et al. [10] presented CrisisSpot, a social context-aware graph-based multimodal framework achieving 88.1% F1-score on the CrisisMMD benchmark, highlighting the importance of incorporating social context in disaster classification. **Low-Resource Language Processing:** Most of Bangla text analysis till date has focused on overall sentiment or abuse detection, rather than disaster events. Nabil et al. [11] collected Bangla social media content to categorize emergency posts, with 95.25% F1 based on XLM-RoBERTa. Ghosh et al. [12] mention about 84% macro-F1 for multilingual disaster tweet classification using graph-augmented attention networks with transformers. Farjana et al. [2] demonstrated that a CNN–LSTM model with BanglaBERT layer can reach 97.94% accuracy for gender-abuse detection task in Bangla. **Multimodal Research in Bangla:** Karim et al. [13] achieved 83% F1-score using XLM-RoBERTa + DenseNet-161 for Bengali multimodal hate speech detection from memes and texts, demonstrating the viability of transformer–CNN fusion architectures for Bengali social media content analysis. Taheri et al. [14] obtained 77.5% F1-score on emotion classification, while Alam et al. [15] addressed Bangla multimodal aggressive meme classification with 76% weighted F1-score. However, both used datasets with fewer than 4,000 samples and did not include disaster-specific tasks.

A. Research Gaps and Motivation

Despite being used by over 300 million people in high-risk locations, there is currently no systematic multimodal framework for classifying disasters in Bangla. Second, Bangla multimodal studies have so far only been conducted on non-disaster applications and small-scale datasets. This paper addresses these shortcomings by offering the first thorough multimodal large-scale disaster classification system in Bangla.

III. DATASET

A novel multimodal dataset was constructed for Bangla disaster classification, comprising social media text and corresponding images capturing real-world disaster complexity in Bangladeshi scenarios.

A. Data Collection and Annotation

We gathered data from public Bangla social media platforms and local news websites focusing on disaster-related content. The collection process involved systematic sampling from Facebook posts, Twitter feeds, and online news portals during major disaster events in Bangladesh from 2020–2023. Two trained disaster-management annotators, native speakers of Bangla with expertise in emergency response, independently annotated all samples achieving Cohen’s kappa ($\kappa = 0.82$),

ensuring high mutual agreement. Our samples consist of single Bangla sentences and their corresponding images, annotated into one of nine disaster classes revealing disaster impacts in Bangladeshi scenarios.

B. Dataset Statistics and Characteristics

Our dataset includes 5,037 annotated samples well-balanced across different disaster classes. Text samples average 15.3 words with 12,847 distinct tokens, covering diverse linguistic patterns typical of Bangla social media communication. Image samples have standard size of 224×224 pixels and exhibit a broad range of visual contexts including infrastructure damage, natural landscapes, human activities, and weather conditions. We divided our dataset using strategic sampling: 70% training (3,526 samples), 10% validation (504 samples), and 20% test (1,007 samples), ensuring balanced class distribution.

TABLE I: Class Distribution in the Bangla Disaster Dataset

Class	Count	Percentage
Agricultural Damage (AD)	800	16%
Non Damage (ND)	650	13%
Infrastructural Damage (ID)	450	9%
Landslides (LS)	400	8%
Damage to Natural Landscape (DNL)	600	12%
Floods (FL)	500	10%
Fires (FR)	300	6%
Economic Loss (EL)	400	8%
Others (OT)	937	18%
Total	5,037	100%

IV. METHODOLOGY

Our approach integrates Bangla text and image data by a single multimodal deep learning pipeline as in Figure 2. The framework consists of four main components: data preprocessing, feature extraction, early fusion, and classification. The architectural choices were motivated by computational efficiency requirements for real-time disaster response and proven effectiveness in low-resource language scenarios. We selected BanglaBERT for native Bangla understanding, mBERT for multilingual robustness, and XLM-RoBERTa for superior cross-lingual transfer learning capabilities [16]. ResNet50, DenseNet169, and MobileNetV2 were chosen over Vision Transformers for computational efficiency in real-time disaster scenarios [10].

A. Data Preprocessing Pipeline

The preprocessing involves making inputs of uniform format for both modalities. **Text Preprocessing:** Input text undergoes deep cleaning by removing punctuation, unwanted whitespace, and non-textual aspects. English and Banglish text are translated into Bangla using Google Translate API, and frequent misspelling errors are corrected to achieve standard preprocessing for Bangla models. Special characters and emojis are normalized to preserve semantic meaning. **Image Preprocessing:** All input images are rescaled to 224×224 pixels according to ImageNet standards, normalized to [0, 1] pixel values and standardized. Data augmentation techniques such as random horizontal flips, rotations ($\pm 15^\circ$), and zooms (0.8-1.2x) are applied to improve model generalization.

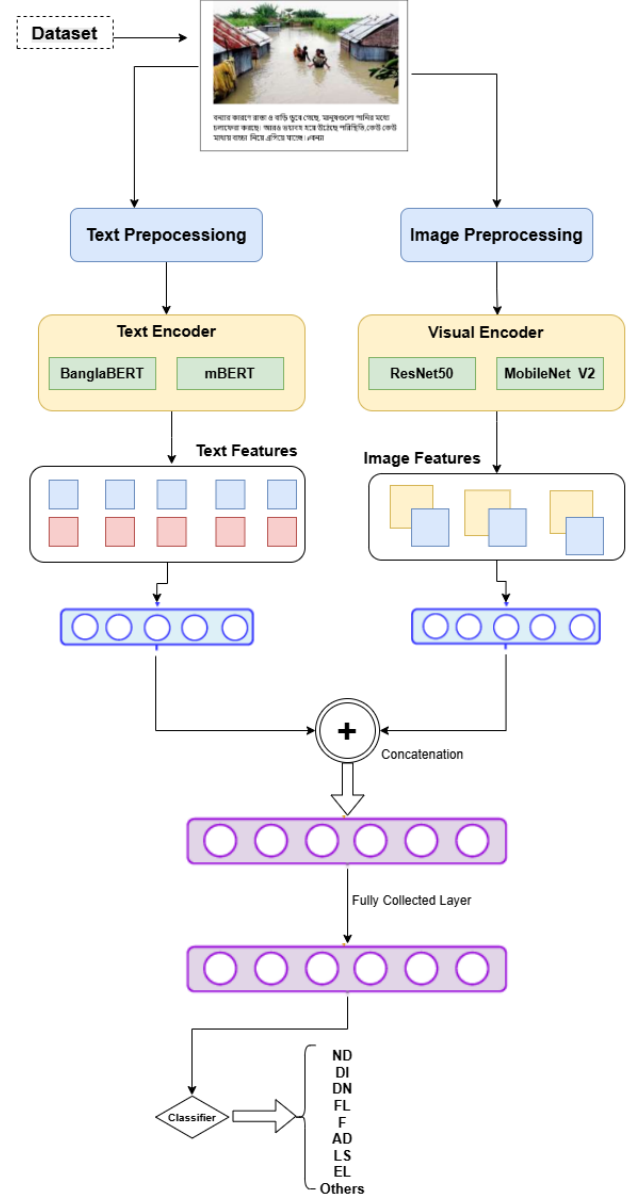


Fig. 2: Overview of the proposed multimodal disaster classification framework.

B. Textual Feature Encoding

The text encoder processes Bangla disaster-related text using transformer-based models (BanglaBERT, mBERT, XLM-RoBERTa) following an effective approach to low-resource language understanding. **Tokenization and Embedding:** All sentences are tokenized into sub-words based on WordPiece tokenization, handling out-of-vocabulary words characteristic of social media content. Tokens are embedded using positional encodings:

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (1)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (2)$$

where pos is the position, i is the dimension index, and d_{model} is the embedding dimension. **Feature Extraction:** Contextual representations are produced by multi-head self-attention mechanisms:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (3)$$

where Q, K, V represent query, key, and value matrices, and d_k is the key dimension. The final text representation is obtained from the [CLS] token embedding.

C. Visual Feature Encoding

Visual features of pre-trained CNNs (ResNet50, DenseNet169 and MobileNetV2) are obtained by eliminating top classification layers and adding global average pooling layers. These models trained on ImageNet offer stable feature representations for disaster images despite domain mismatch. Hierarchical feature extraction computes the visual feature representation:

$$F_{visual} = \text{GlobalAvgPool}(\text{CNN}(I)) \quad (4)$$

where I represents the input image and CNN denotes the convolutional feature extractor. This approach captures low-level patterns (textures, edges) and high-level semantic concepts (damaged buildings, flooding).

D. Early Fusion and Classification

We employ early fusion through feature concatenation over late fusion due to computational efficiency and superior performance for disaster classification tasks where textual and visual information are complementary [8]. The multimodal fusion concatenates textual and visual information:

$$F_{joint} = [F_{text}; F_{visual}] \in \mathbb{R}^{d_{text}+d_{visual}} \quad (5)$$

where $d_{text} = 768$ and $d_{visual} = 2048$ for our best configuration. This early fusion concatenation enables cross-modal feature interactions at the representation level, allowing the model to learn joint patterns leveraging both semantic descriptions and visual evidence simultaneously. This joint representation passes through fully connected layers with dropout (rate 0.1) for regularization, preventing overfitting on the limited disaster dataset. The final classification is performed using softmax activation:

$$\hat{y} = \text{softmax}(W_f F_{joint} + b_f) \quad (6)$$

where $W_f \in \mathbb{R}^{C \times (d_{text}+d_{visual})}$ and $b_f \in \mathbb{R}^C$ are learned fusion parameters, and $C = 9$ is the number of disaster classes.

E. Training Configuration

The training process is designed to manage multimodal training with constrained disaster data. We use Adam optimizer with different learning rates: 1×10^{-5} for pre-trained text encoders and 3×10^{-5} for fusion layers respectively. Batch

size is 32 with early stopping based on validation loss. The optimization follows categorical cross-entropy loss:

$$\mathcal{L} = - \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic}) \quad (7)$$

V. EXPERIMENTAL RESULTS

This section validates comprehensive evaluation of the BanglaMM-Disaster framework comparing unimodal and multimodal approaches. We examine performance indicators, error patterns and show the capability of multimodal fusion on Bangla disaster classification.

A. Unimodal Performance Analysis

1) *Visual-Only Models:* Table II summarizes image-only model performance. ResNet50 achieved the best accuracy of 66.85%, outperforming DenseNet169 (65.47%) and ResNet101 (64.32%). Figure 3 shows that visual models effectively separate distinct classes like Non Damage (ND) and Floods (FL), but struggle with ambiguous ones like Fires (FR) and Economic Loss (EL).

TABLE II: Performance of Image-Only Models

Model	Acc.	Prec.	Recall	F1
ResNet50	66.85	66.62	67.08	66.85
DenseNet169	65.47	65.23	65.71	65.47
ResNet101	64.32	64.15	64.49	64.32
MobileNetV2	62.41	62.17	62.65	62.41

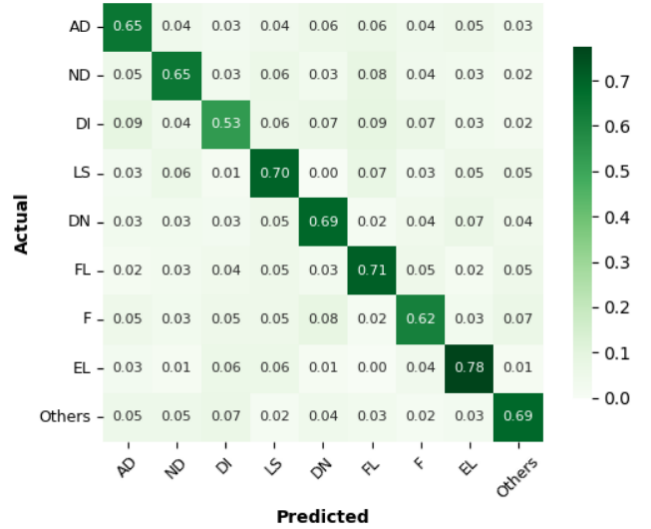


Fig. 3: Confusion matrix for best visual model (ResNet50).

2) *Text-Only Models:* Text-only models performed significantly better than visual models as shown in Table III. XLM-RoBERTa demonstrates superior multilingual support for Bangla text with 79.92% accuracy. The model effectively distinguishes between Agricultural Damage (AD) and Non Damage (ND) but encounters challenges with textual complexity in Economic Loss (EL) and Others categories.

TABLE III: Performance of Text-Only Models

Model	Acc.	Prec.	Recall	F1
XLM-RoBERTa	79.92	79.65	80.21	79.93
mBERT	78.15	77.82	78.48	78.15
BanglaBERT	76.73	76.29	77.18	76.73

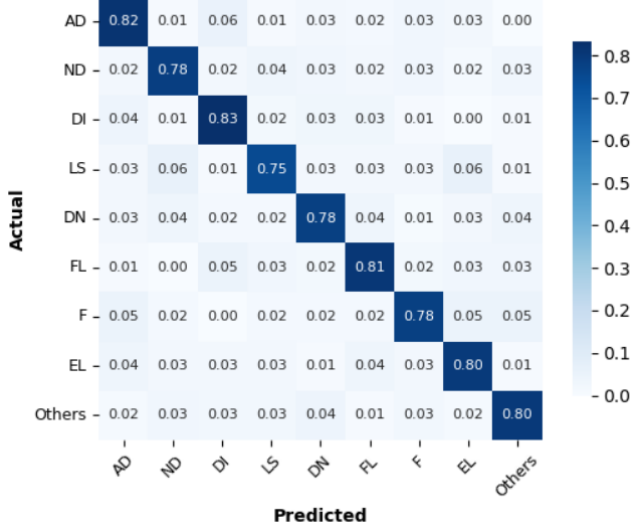


Fig. 4: Confusion matrix for best text model (XLM-RoBERTa).

B. Multimodal Fusion Results

Table IV presents comprehensive multimodal performance. The optimized mBERT + ResNet50 fusion achieved the best accuracy of 83.76%, representing a 3.84% improvement over the best text-only model and 16.91% over the best visual-only model. The early fusion concatenation mechanism enables this improvement by allowing cross-modal feature interactions that disambiguate ambiguous textual descriptions using visual confirmation. Even the worst multimodal combination outperformed the best visual-only baseline, demonstrating consistent multimodal improvements.

C. Comparative Analysis

Table V compares our results with related work. Due to the lack of Bangla multimodal disaster datasets, direct comparison is challenging as different datasets and languages are used across studies. However, our approach achieves competitive performance despite working with a morphologically complex low-resource language and challenging 9-class taxonomy. For real-time deployment, our framework achieves 0.45 seconds average inference time with 1.8GB memory footprint on standard GPU hardware.

TABLE V: Comparison with Related Work

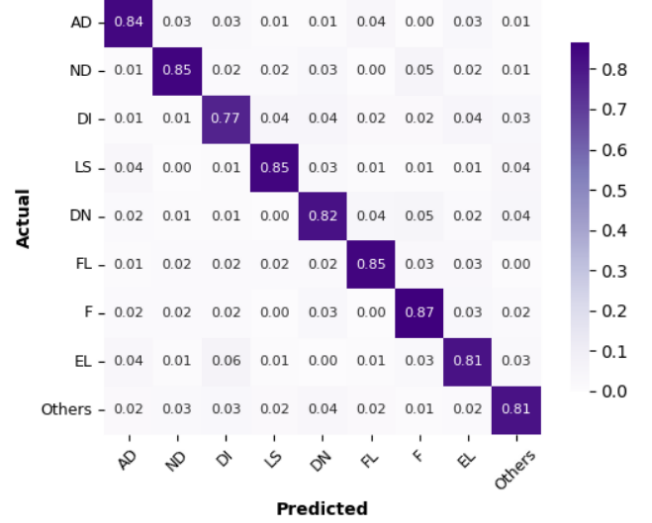


Fig. 5: Confusion matrix for best multimodal model (mBERT+ResNet50).

Method	Domain	Classes	F1 (%)
<i>Bengali Multimodal</i>			
Karim et al. [13]	Hate Speech	2	83.0
Taheri et al. [14]	Emotion	4	77.5
Alam et al. [15]	Meme	2	76.0
<i>English Multimodal Disaster (Context)</i>			
CrisisSpot [10]	Disaster	6	88.1
El-Niss et al. [8]	Disaster	6	85.2
Khattar et al. [9]	Disaster	5	84.1
<i>Bengali Text-Only Disaster</i>			
Nabil et al. [11]	Emergency	2	95.3
Ghosh et al. [12]	Disaster	Multi	84.0
Our Work	Disaster	9	83.76

D. Error Analysis and Cross-Modal Benefits

Figure 6 demonstrates consistent superiority of multimodal fusion across all disaster classes. The multimodal model achieves the lowest error rates for each disaster type, with complementary information from both modalities contributing to performance gains. For example, the Fires (FR) class error drops from 45.3% (visual-only) and 28.4% (text-only) to 22.7% with multimodal fusion, a 50% improvement over visual baseline. Similarly, Economic Loss (EL) errors decrease from 39.2% (visual) and 24.6% (text) to 19.4% (multimodal). These consistent improvements demonstrate that the early fusion approach effectively captures cross-modal dependencies.

VI. CONCLUSION

This research introduces BanglaMM-Disaster, a comprehensive multimodal deep learning framework for disaster classification in Bangla, addressing critical gaps in disaster analytics for low-resource languages. We developed a novel dataset of 5,037 annotated Bangla social media posts across 9 disaster

TABLE IV: Performance of Multimodal Model Combinations

Model	Acc.	Prec.	Recall	F1
mBERT + ResNet50 (Optimized)	83.76	83.54	83.98	83.76
mBERT + ResNet50	82.34	82.15	82.53	82.34
XLM-RoBERTa + DenseNet169	81.87	81.64	82.11	81.87
BanglaBERT + ResNet101	81.42	81.18	81.66	81.42
mBERT + DenseNet201	80.95	80.73	81.17	80.95

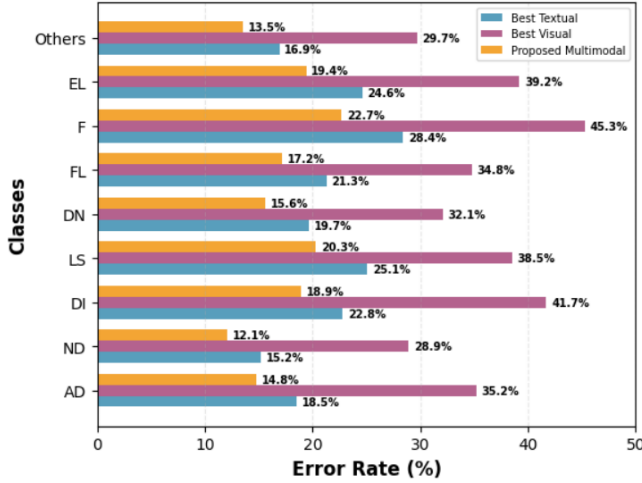


Fig. 6: Class-wise error rates comparing best visual, text, and multimodal models.

categories and demonstrated that early fusion of transformer-based text encoders with CNN visual features achieves 83.76% accuracy, with significant improvements of 3.84% over text-only and 16.91% over image-only baselines. Our error analysis confirms consistent cross-modal benefits across all disaster types. The framework's computational efficiency (0.45s inference, 1.8GB memory) enables practical deployment for real-time disaster monitoring systems. This work establishes new benchmarks for Bangla multimodal disaster classification and demonstrates the effectiveness of multimodal learning in low-resource linguistic settings. Future directions include exploring attention-based fusion mechanisms and graph neural networks [10] to capture more complex cross-modal relationships for enhanced disaster understanding.

REFERENCES

- [1] C. Caragea, A. Silvescu, and A. H. Tapia, "Identifying informative messages in disaster events using convolutional neural networks," in *Proceedings of the 13th International Conference on Information Systems for Crisis Response and Management*. ISCRAM Association, 2016, pp. 137–147.
- [2] M. Farjana, S. Afroge, and A. Y. Srizon, "Gender abusive bengali text classification using enhanced cnn-lstm model with bangla bert base preprocessing," in *2024 27th International Conference on Computer and Information Technology (ICCIT)*. IEEE, 2024, pp. 1093–1098.
- [3] S. Kamoji, M. Kalla, and C. Joshi, "Fusion of multimodal textual and visual descriptors for analyzing disaster response," in *2023 5th International Conference on Smart Systems and Inventive Technology (ICSSIT)*. IEEE, 2023, pp. 1614–1619.
- [4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, 2019, pp. 4171–4186.
- [5] J. Han, W. Zhang, A. Morrison, H. Huang, L. Hou, and R. Ranjan, "Quakebert: Accurate classification of social media texts for emergency management using deep learning transformer models," *International Journal of Disaster Risk Reduction*, vol. 104, p. 104354, 2024.
- [6] S. Kaur, S. Gupta, S. Singh, V. T. Hoang, S. Almakdi, T. Alelyani, and A. Shaikh, "Transfer learning-based automatic hurricane damage detection using satellite images," *Electronics*, vol. 11, no. 9, p. 1448, 2022.
- [7] H. Mouzannar, Y. Rizk, and M. Awad, "Damage identification in social media posts using multimodal deep learning," in *Proceedings of the 15th International Conference on Information Systems for Crisis Response and Management*. ISCRAM Association, 2018, pp. 529–543.
- [8] A. El-Niss, A. Alzu'Bi, and A. Abuarqoub, "Multimodal fusion for disaster event classification on social media: A deep federated learning approach," in *Proceedings of the 7th International Conference on Future Networks and Distributed Systems*, 2023, pp. 758–763.
- [9] A. Khattar and S. Quadri, "Camm: cross-attention multimodal classification of disaster-related tweets," *IEEE Access*, vol. 10, pp. 92 889–92 902, 2022.
- [10] H. Dar, F. Alam, and C. Castillo, "A social context-aware graph-based multimodal attentive framework for disaster content classification during emergencies," *Expert Systems with Applications*, vol. 252, p. 124154, 2024.
- [11] A. A. Nabil, D. Das, M. S. Salim, S. Arifeen, and H. A. Fattah, "Bangla emergency post classification on social media using transformer based bert models," in *2023 6th International Conference on Electrical Information and Communication Technology (EICT)*. IEEE, 2023, pp. 1–6.
- [12] S. Ghosh, S. Maji, and M. S. Desarkar, "Gnom: graph neural network enhanced language models for disaster related multilingual text classification," in *Proceedings of the 14th ACM Web Science Conference 2022*, 2022, pp. 55–65.
- [13] M. R. Karim, S. K. Dey, T. Islam, M. Shajalal, and B. R. Chakravarthi, "Multimodal hate speech detection from bengali memes and texts," *arXiv preprint arXiv:2204.10196*, 2022.
- [14] Z. S. Taheri, A. C. Roy, and A. Kabir, "Bemofusionnet: A deep learning approach for multimodal emotion classification in bangla social media posts," in *2023 26th International Conference on Computer and Information Technology (ICCIT)*. IEEE, 2023, pp. 1–6.
- [15] M. A. Alam, J. Hossain, S. Ahsan, and M. M. Hoque, "Multimodal aggressive meme classification using bidirectional encoder representations from transformers," in *2024 27th International Conference on Computer and Information Technology (ICCIT)*. IEEE, 2024, pp. 3542–3547.
- [16] A. Bhattacharjee, T. Hasan, W. U. Ahmad, K. S. Mubasshir, M. S. Islam, A. Iqbal, M. S. Rahman, and R. Shahriyar, "Banglabert: Language model pretraining and benchmarks for low-resource language understanding evaluation in bangla," in *Findings of the Association for Computational Linguistics: NAACL 2022*. Association for Computational Linguistics, 2022, pp. 1318–1327.