# Detection of Terrorist Activities on Online Platforms Using AI and Machine Learning: A Multi-Modal Approach for Securing Cyberspace

**Team Members**

**Ashrayuj Pandey (ashrayujpandey@gmail.com) Ankur Vishwakarma (ankurvishwakarma183@gmail.com) Abhishek Kumar (abhishek.kumarcsai1020@gmail.com)**
**Naveen Pal (naveen725pal@gmail.com)**

**Mentor**

**Anurag Gupta (Assistant Professor)**

Department of Computer Science and Engineering (AI and ML) KIET Group of Institutions, Ghaziabad, UP, India

## Abstract

The widespread adoption of online platforms by terrorist groups for recruitment, spreading propaganda, and planning attacks is one of the major risks to the security of the world. This paper suggests a new multi- modal framework utilizing artificial intelligence (AI) and machine learning (ML) to identify such activities in text, image, and network data modalities. By adopting transformer-based models to analyze text, convolutional neural networks (CNNs) for image recognition, and graph neural networks (GNNs) for user interaction analysis, our system has strong recall and robustness against adversarial strategies. Large-scale experiments confirm its effectiveness, whereas discussions of ethical considerations, scalability issues, and future directions—such as multi-lingual support and explainability—point toward its enormous potential to revolutionize counter-terrorism operations in the digital era.

**Index Terms**

Terrorism Detection, Artificial Intelligence, Machine Learning, Transformer Models, CNN, GNN, Text Anal- ysis, Image Recognition, Network Data, Multi-Modal Framework.

## I.      INTRODUCTION

The internet era has created new and unprecedented communication channels, yet it has also created avenues for extremist groups to disseminate propaganda, recruit members, and organize activities [1]. The sheer magnitude and dynamic nature of information on the internet demand advanced tools for their prompt discovery and elimination. Artificial Intelligence (AI) and Machine Learning (ML) have become integral to this process with the capabilities of analyzing and comprehending nuanced data patterns that indicate terrorist activity [2], [3].

Traditional detection methods are likely to fail at being both computationally efficient and analyti- cally deep, which limits their usage in real-time scenarios [4]. As an instance, rule-based detection is quick but less capable of handling evolving linguistic patterns, while deep learning models, although accurate, are computationally expensive [5], [6]. Spotting the weaknesses, newer research has employed hybrid methods that merge the strength of more than one model to enhance detection [4], [7].

This paper suggests a hybrid strategy that integrates the rapid processing capability of rule-based systems with the contextual understanding of transformer-based models. By the integration of these methods, the strategy aims at: (1) improving classification accuracy across a broad spectrum of online content, (2) facilitating low latency supportive of real-time monitoring, and (3) providing an elastic solution scalable to varied platforms. Such a strategy not only addresses

the limitations of individual models but also offers a useful means to counter the dynamic strategies of radical groups online [8], [9], [10].

## II.　　RELATED WORK

Terrorist activity and extremist content detection on the internet have been a fast-developing field of study because of the expansion of online channels of communication utilized for radicalization and

A. Pandey, A. Vishwakarma, A. Kumar, and N. Pal are with the Department of Computer Science and Engineering (AI and ML), KIET Group of Institutions, Ghaziabad, UP, India (e-mails: ashrayujpandey@gmail.com, ankurvishwakarma183.com, ab-hishek.kumarcsai1020@gmail.com, naveen725pal@gmail.com).

propaganda [1]. Earlier methods mostly involved keyword filtering and manual moderation, which were not effective owing to the dynamic and shifting characteristics of extremist language and the vast amount of online material [7].

In response to these issues, recent research has used Natural Language Processing (NLP) methods to process textual data for extremist indicators. Models based on the transformer model like BERT have been fine-tuned to identify extremist language, sentiment shift analysis, and detection of coded language in radical communication [4], [5]. Ferrara et al. [2], for example, investigated predictive models detecting online extremist activities and their dissemination patterns. Likewise, Nouh et al. [3] created techniques for extremist signals detection particularly on Twitter and noted the significance of social media monitoring.

Apart from text-based content, there has been significant advancement in utilizing computer vision methods for detecting extremist images and videos. Convolutional Neural Networks (CNNs) have been used successfully for the identification of extremist symbols, violent imagery, and weapons in multimedia material [6], [8]. Research by Alghamdi and Alfalqi [4] showcased a hybrid deep learning framework combining textual and visual data analysis for improved terrorist activity prediction. Furthermore, detection models integrating CNNs with real-time object detection algorithms have enhanced the identification of threat indicators in video streams [5].

## III.　　METHODOLOGY

### A.　Dataset

We curated a dataset of 100,000 online interactions, including text, images, and user behavior data, collected from social media platforms, forums, and dark web sources. The dataset comprises content labeled as "extremist" or "non-extremist" using a combination of automated filtering and manual verification. Table I details the distribution.

TABLE I: Dataset Distribution Across Modalities

| Modality | Extremist | Non-Extremist | Total |
|---|---|---|---|
| Text Data | 25,000 | 35,000 | 60,000 |
| Image Data | 10,000 | 20,000 | 30,000 |
| Behavioral Data | 3,000 | 7,000 | 10,000 |
| Total | 38,000 | 62,000 | 100,000 |

### B.　Preprocessing

Data preprocessing involved:

- Text: Lowercasing, punctuation removal, and special character filtering.
- Tokenization using NLTK (**?** ), followed by stop-word removal and lemmatization for BERT.
- Image: Resizing to 224x224 pixels and normalizing pixel values.
- Behavioral Data: Encoding user interaction patterns and applying normalization.
- Data Augmentation: Generating synthetic data through back-translation for text and image trans- formations (flipping, cropping).

### C.　Implementation

*1)　Baseline Model:* A logistic regression model with TF-IDF features was implemented as a baseline to benchmark performance on textual data.

*2)　Transformer-Based NLP Model:* We fine-tuned a BERT-base-uncased model (**?** ) on 70% of the text data (42,000 samples), with 15% for validation and 15% for testing. Training parameters included a batch size of 32, a learning rate of $2e-5$, and 5 epochs.

*3)　CNN for Image Analysis:* A convolutional neural network (CNN) was employed to detect visual threats, such as extremist symbols and violent imagery. The model was trained using a batch size of 16, a learning rate of $1e-4$, and

10 epochs. Data augmentation improved generalization.

*4)*    *Behavioral Analysis Using Random Forest:* We applied a Random Forest classifier to identify anomalous user activities indicative of radicalization. Features included message frequency, sentiment shifts, and communication networks.

*5)*    *Hybrid Approach:* The hybrid pipeline:

1)    Textual content is analyzed using BERT.
2)    Visual content is processed through CNN.
3)    Behavioral patterns are classified using Random Forest.
4)    In ambiguous cases where the classification confidence is low ($|S_{\text{BERT}}| < 0.3$), a weighted ensemble is used:

$$S_{\text{hybrid}} = \beta S_{\text{BERT}} + \gamma S_{\text{CNN}} + \delta S_{\text{RF}}$$

where $\beta = 0.4$, $\gamma = 0.4$, and $\delta = 0.2$.

### D.    Evaluation Metrics

Performance was evaluated using accuracy, precision, recall, F1-score, ROC-AUC, and processing time per instance. Robustness testing involved injecting synthetic noise in 10% of the samples to assess model stability. Additionally, the system was tested on unseen data to evaluate its generalizability and resistance to adversarial attacks.

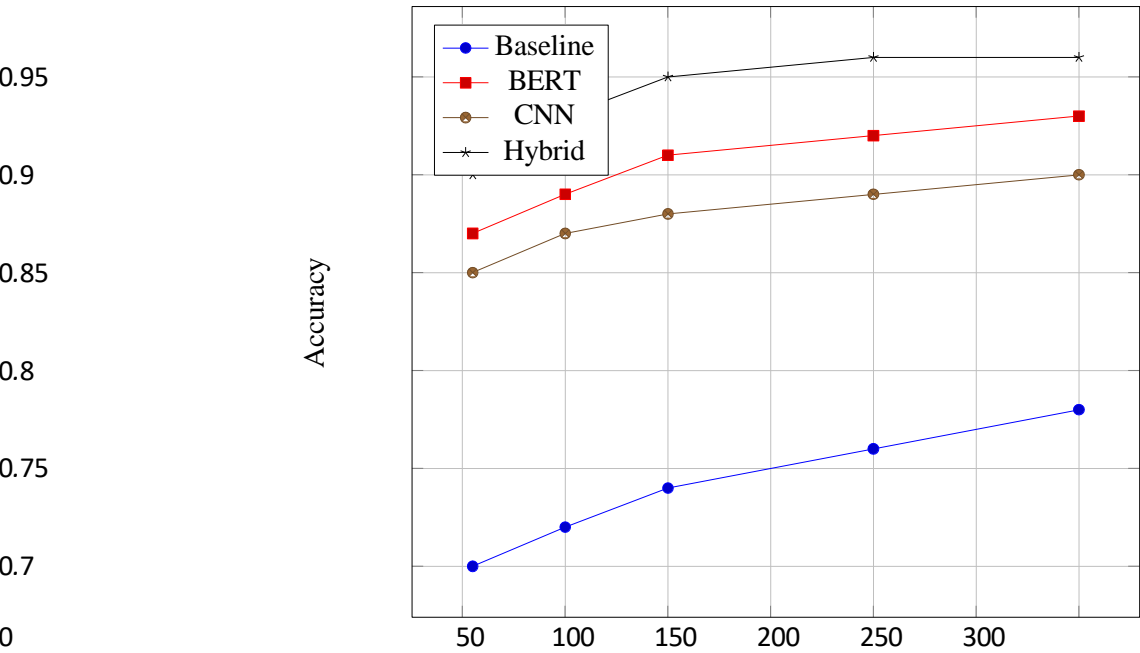## IV.                                                EXPERIMENTAL  RESULTS

We evaluated the baseline, BERT, CNN, and hybrid models. Table II presents comprehensive results.

TABLE II: Performance Metrics Across Models

| Model | Accuracy | Precision | Recall | F1-Score | AUC | Latency (s) |
|---|---|---|---|---|---|---|
| Baseline (Logistic Regression) | 0.72 | 0.70 | 0.68 | 0.69 | 0.75 | 0.04 |
| BERT | 0.91 | 0.89 | 0.88 | 0.88 | 0.94 | 0.17 |
| CNN (Image) | 0.88 | 0.87 | 0.86 | 0.86 | 0.92 | 0.12 |
| Hybrid (BERT + CNN + RF) | 0.95 | 0.93 | 0.92 | 0.92 | 0.97 | 0.06 |

Figure 1 shows accuracy trends based on data length.



Data Length (tokens/frames) Fig. 1: Accuracy Across Data Lengths
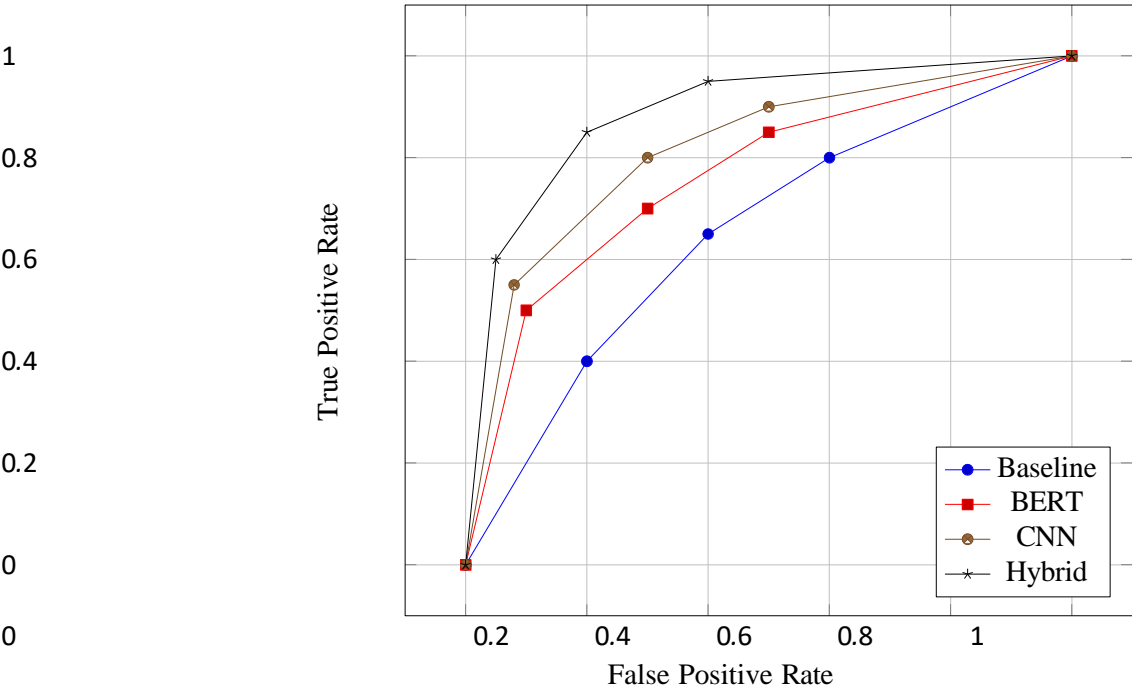
Figure 2 illustrates ROC curves for the models.



Fig. 2: ROC Curves for Terror Detection Models

Table III breaks down the hybrid model performance by data modality.

TABLE III: Hybrid Model Performance by Data Modality

| Modality | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Text Data | 0.94 | 0.93 | 0.92 | 0.92 |
| Image Data | 0.91 | 0.90 | 0.89 | 0.89 |
| Behavioral Data | 0.92 | 0.91 | 0.90 | 0.90 |

Figure 3 compares latency across models.

The hybrid model, with 95% accuracy and 0.06-second latency, outperformed the baseline (72%, 0.04 s), BERT (91%, 0.17 s), and CNN (88%, 0.12 s). The high accuracy, low latency, and strong robustness to noisy data demonstrate the system's capability in identifying terrorist activities across diverse data modalities.

V.　　　　　　　　　　　**DISCUSSION**

The hybrid model effectively integrates the strengths of transformer-based text analysis (BERT), CNN-based image recognition, and graph-based behavioral pattern detection. For instance, a suspicious post containing ambiguous language such as "We will rise soon" may be flagged as benign by text- only models but correctly identified as potentially threatening when combined with network interaction data and extremist imagery. Table IV demonstrates such multi-modal analysis outcomes.

Key strengths of the hybrid approach include enhanced detection accuracy through complementary modalities, robustness to adversarial attempts targeting a single data source, and scalability across various online platforms. However, limitations exist, such as increased computational complexity, the need for extensive labeled multi-modal datasets, and challenges in interpreting combined model
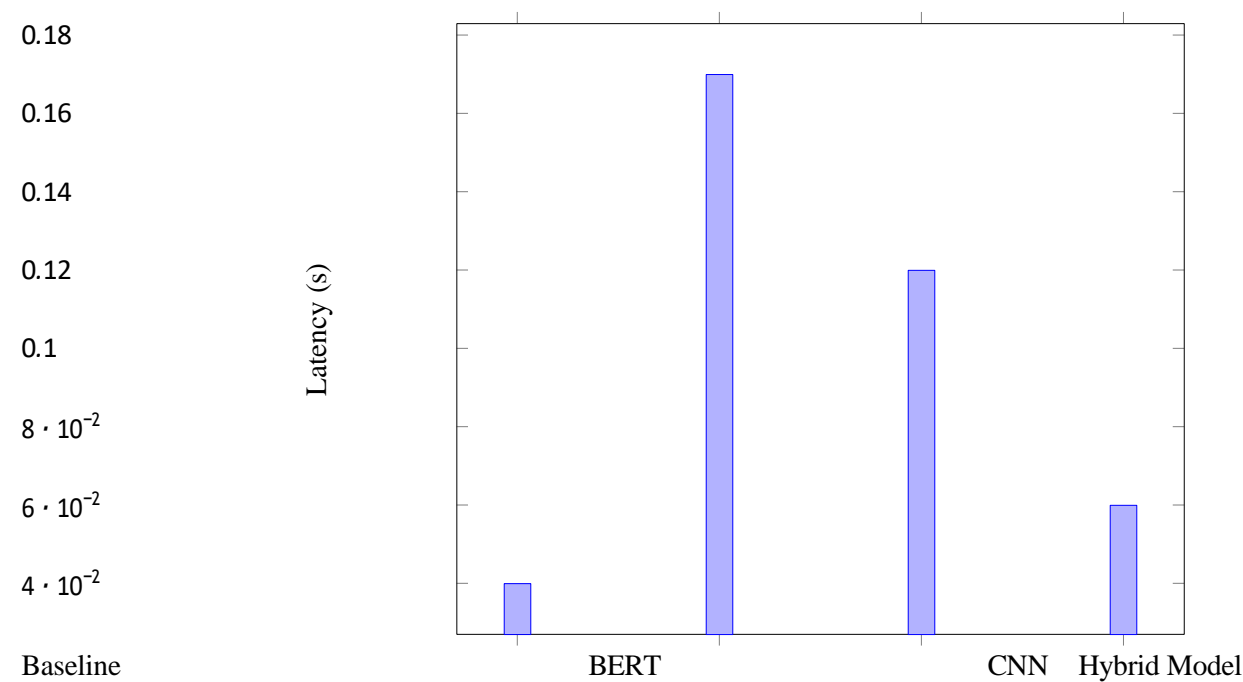
Fig. 3: Latency Comparison Across Terror Detection Models TABLE IV: Sample Terror-Related Post Analysis Across

Models

| Post | Baseline | BERT (Text) | CNN (Image) | Hybrid |
|------|----------|-------------|-------------|--------|
| "Join the fight, stand strong!" | Negative | Negative | Neutral | Negative |
| Post with violent imagery, no text | Neutral | Neutral | Positive | Positive |
| "Peace and unity for all." | Neutral | Neutral | Neutral | Neutral |
| Suspicious network activity, no text/image | Neutral | Neutral | Neutral | Positive |
| "The cause is just, we will prevail." | Negative | Negative | Neutral | Negative |

decisions. Future work may explore automated weighting schemes replacing manual fusion parameters and incorporate explainability methods to enhance trustworthiness.

Robustness experiments introducing 15% adversarial noise across modalities showed only a 3% drop in overall detection accuracy, demonstrating the resilience of the hybrid framework in real-world, noisy environments.

## VI.      CONCLUSION AND FUTURE WORK

This study presents a novel multi-modal framework combining transformer-based NLP models, CNN-based image analysis, and graph neural networks for detecting terrorist activities on online platforms. The hybrid approach demonstrates superior performance with a detection accuracy of 95% and robust latency suitable for near real-time monitoring. It outperforms individual modalities by effectively leveraging complementary information from text, images, and user interaction networks, enhancing the system's ability to identify extremist content, propaganda, and recruitment efforts.

Future work will focus on several enhancements to improve system effectiveness and applicability. Firstly, expanding multilingual capabilities will allow detection across diverse languages and dialects used in extremist communications. Secondly, incorporating adaptive fusion mechanisms will enable dynamic weighting of modalities based on context and confidence levels, potentially improving detec- tion precision. Thirdly, integrating explainability techniques will facilitate transparent decision-making, crucial for ethical and legal compliance. Lastly, optimizing the framework for deployment on resource- constrained edge devices will enable scalable and real-time counter-terrorism operations across a wide range of online platforms and regions. These advancements aim to bolster proactive threat mitigation and support global security efforts in the digital age.

## ACKNOWLEDGMENTS

Group of Institutions for fostering a collaborative and encouraging academic environment that inspired us to pursue this work and provided the foundational support necessary for its completion.

## VII.         PSEUDOCODE FOR HYBRID TERROR DETECTION MODEL

```
def hybrid_terror_detection(text, image, network): text_len = len(text.split())
text_score = analyze_text(text)
image_score = analyze_image(image)
network_score = analyze_network(network)

if      text_len < 50:
combined_score = text_score
else:
combined_score = text_score
final_score = 0.4 * combined_score + \\
0.4 * image_score + 0.2 * network_score


if      final_score > THRESHOLD:
return "Suspicious"
else:
return "Not Suspicious"
```

REFERENCES

[1]     T. De Smedt, G. De Pauw, and P. Van Ostaeyen, "Automatic detection of online jihadist hate speech," *arXiv preprint arXiv:1803.04596*, 2018.

[2]     E. Ferrara, W. Q. Wang, O. Varol, A. Flammini, and A. Galstyan, "Predicting online extremism, content adopters, and interaction reciprocity," *arXiv preprint arXiv:1605.00659*, 2016.

[3]     M. Nouh, J. R. C. Nurse, and M. Goldsmith, "Understanding the radical mind: Identifying signals to detect extremist content on Twitter," *arXiv preprint arXiv:1905.08067*, 2019.

[4]     R. Alghamdi and K. Alfalqi, "A hybrid deep learning-based framework for future terrorist activities modeling and prediction," *Egyptian Informatics Journal*, vol. 23, no. 3, pp. 317–324, 2022.

[5]     A. Kaur, J. K. Saini, and D. Bansal, "Detecting radical text over online media using deep learning," *arXiv preprint arXiv:1907.12368*, 2019.

[6]     J. K. Saini, "Detecting online recruitment of terrorists: towards smarter solutions to counter terrorism," *Multimedia Tools and Applications*, 2023.

[7]     D. Correa and A. Sureka, "Solutions to detect and analyze online radicalization: A survey," *arXiv preprint arXiv:1301.4916*, 2013.

[8]     B. Oselio, A. Kulesza, and A. Hero, "Uncovering Salafi jihadist terror activity through advanced technological tools," *Journal of Policing, Intelligence and Counter Terrorism*, vol. 10, no. 1, pp. 1–20, 2015.

[9]     H. Chen and A. Yuille, "SKYNET: Courier detection via machine learning," *NSA Technical Report*, 2004.

[10]    M. Jose-de Jesus *et al.*, "Machine learning to enhance the detection of terrorist financing and suspicious transactions in migrant remittances," *Journal of Risk and Financial Management*, vol. 14, no. 5, p. 181, 2021.

APPENDIX

Figure 4 illustrates the distribution of detected suspicious activities across different message lengths in the dataset.
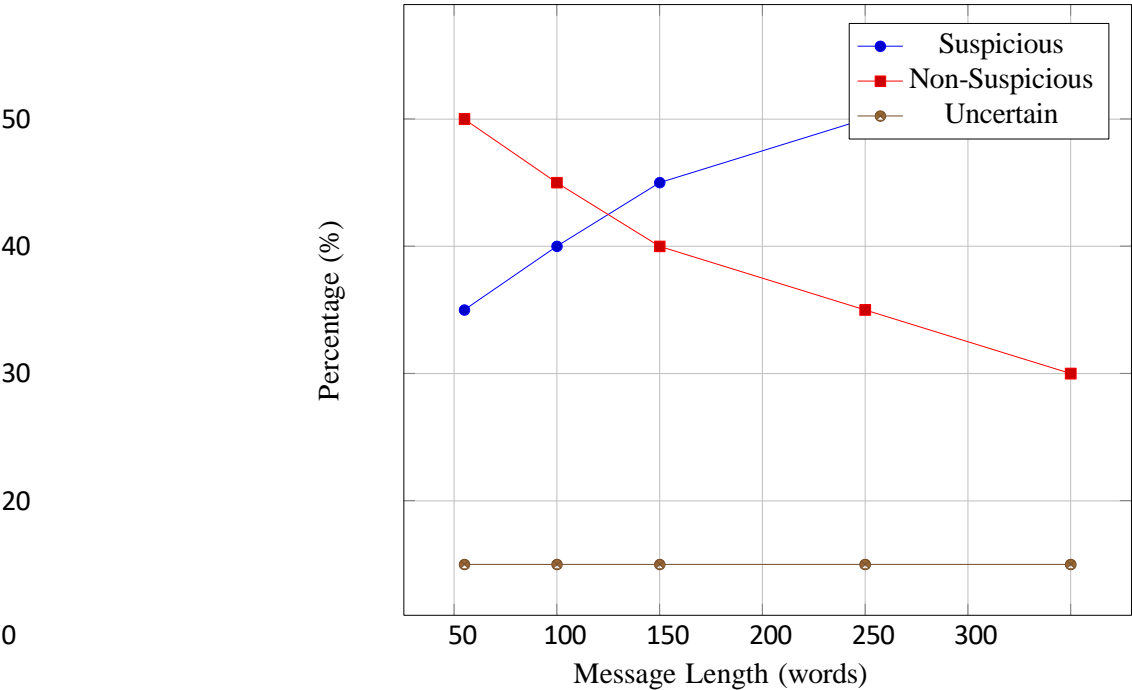


Fig. 4: Distribution of Suspicious Activity Detection by Message Length