



International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



# **Digital Purifier**

Shaima Shahul<sup>1</sup>, Renjini L R<sup>2</sup>, Harikrishnan S R<sup>3</sup>

Student, MCA, CHMM College for Advanced Studies, Trivandrum, India<sup>1</sup> Associate Professor, MCA, CHMM College for Advanced Studies, Trivandrum, India<sup>2</sup> Associate Professor, MCA, CHMM College for Advanced Studies, Trivandrum, India<sup>3</sup>

Abstract: The digital age is increasingly challenged by the spread of manipulated media (deep fakes) and the pervasive issue of online hate speech and toxic comments. This project proposes an integrated system, deployed as a Flask web application, to address these critical issues. The system will comprise two primary modules: a deep fake detection module and a toxic comment analysis and removal module. The deepfake detection module will employ Generative Adversarial Networks (GANs) to accurately classify uploaded images as either fake or real, analyzing subtle inconsistencies indicative of manipulation. The toxic comment module will utilize BERT (Bidirectional Encoder Representations from Transformers) to identify and categorize hate speech and toxic language within user-generated text. Upon detection of toxic comments, the system will automatically remove the offending content. Furthermore, all detected toxic comments, along with associated metadata, will be systematically recorded and saved to an Excel spreadsheet for future analysis and moderation purposes. This comprehensive platform aims to enhance online safety and information integrity by providing real-time detection and mitigation of both deepfakes and toxic content, while maintaining a detailed record for ongoing review and improvement

Keywords: GAN, BERT, NLP, Pytorch, Machine Learning, Deep Learning

### I. INTRODUCTION

In today's digital landscape, the authenticity of online content and the quality of user interactions are increasingly under threat. The rise of deepfakes-highly realistic yet artificially generated images and videos-has made it difficult to distinguish real media from manipulated content, fueling misinformation and eroding public trust. At the same time, social platforms are plagued by toxic comments and hate speech, contributing to harassment, mental health issues, and a hostile online environment. To address these dual challenges, this project presents Digital Purifier, a unified AI-driven web application that integrates advanced machine learning techniques for content moderation. The system includes two main modules: a deepfake detection module that leverages Generative Adversarial Networks (GANs) to identify and flag fake images by analyzing visual inconsistencies, and a toxic comment analysis module that uses BERT (Bidirectional Encoder Representations from Transformers) to detect, classify, and automatically remove harmful textual content. Furthermore, all flagged toxic comments and related metadata are logged in an Excel file for future review and moderation enhancement. By combining cutting-edge image and text analysis tools, Digital Purifier aims to foster a safer, more credible, and respectful digital ecosystem. This project not only enhances the security and authenticity of digital content but also promotes healthier online interactions by automating the moderation of harmful material. Its modular architecture allows for scalable deployment across various platforms cloud and local and its resilience mechanisms ensure robust performance even under high content volume or adversarial input. The ultimate goal of the Digital Purifier is to create a safer digital environment by proactively identifying and mitigating the effects of deceptive visual media and offensive language, making it a powerful tool for digital content moderation, education platforms, social networks, and media verification services.

### **II. LITERATURE REVIEW**

The project proposes an integrated system to combat manipulated media and online hate speech through a Flask web application. It includes two main modules: a deep fake detection module using Generative Adversarial Networks

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

#### Volume 5, Issue 1, July 2025



(GANs) to classify images as real or fake, and a toxic comment analysis module utilizing BERT to identify and remove hate speech and toxic language. Detected toxic comments will be logged with metadata in an Excel spreadsheet for future analysis. The platform aims to enhance online safety and information integrity by providing real-time detection and mitigation of deepfakes and toxic content. paragraphs must be indented. All paragraphs must be justified, i.e. both left-justified and right-justified.

### **III. PROPOSED METHOD**

The proposed method in the Digital Purifier project is centered on the integration of advanced Artificial Intelligence (AI) techniques to detect and mitigate two major digital threats: deepfake images and toxic user comments. The system adopts a modular architecture that combines image forensics and natural language processing (NLP) within a single web-based interface, thereby offering a holistic solution for automated content moderation. For image analysis, the method utilizes Generative Adversarial Networks (GANs), a deep learning framework where a generator and discriminator are trained simultaneously. The discriminator, after proper training, becomes proficient at identifying forged images by detecting subtle irregularities such as unnatural lighting, texture anomalies, or inconsistencies in facial expressions-typical indicators of deepfakes. On the textual front, the system incorporates BERT (Bidirectional Encoder Representations from Transformers), a state-of-the-art language model developed by Google. BERT enables the system to understand the contextual meaning of words and sentences, making it highly effective for detecting various forms of toxic content, including hate speech, insults, threats, and implicit abusive language.Once a user uploads an image or enters a comment, the data is routed through respective detection pipelines. The deepfake detector preprocesses the image and feeds it into a trained CNN or GAN discriminator model, while the text input undergoes tokenization and contextual encoding before being classified by the BERT model. If an image is flagged as fake or a comment is marked toxic, the system triggers automated moderation actions-such as blocking the comment or flagging the image—while logging the incident for future analysis. This ensures not only immediate moderation but also long-term tracking and accountability. Additionally, all processed data is recorded in a structured log (e.g., Excel), allowing for transparency and retrospective examination of user behaviour. The entire method is encapsulated within a lightweight Flask-based web application, which provides a simple, intuitive interface for users to interact with. Designed to operate seamlessly across both cloud environments (for training, e.g., Google Colab with GPU support) and local servers (for deployment and user interaction), the proposed method ensures platform independence, ease of use, and scalability. In essence, the Digital Purifier's proposed method represents a unified, intelligent, and automated approach to managing digital content, combining the strengths of GANs and BERT to create a system that not only protects users from misinformation and online abuse but also enhances the integrity of digital communication platforms.

#### **IV. ALGORITHM**

### **Convolutional Neural Networks (CNNs)**

In the Digital Purifier system, Convolutional Neural Networks (CNNs) play a critical role in analyzing image-based content for deepfake detection. A CNN is a class of deep learning models specifically designed for processing visual data. It operates by extracting hierarchical patterns from input images through a series of layers, such as convolutional layers, pooling layers, and activation functions. The convolutional layers apply learnable filters across the image to detect features like edges, textures, and facial contours. In the context of deepfake detection, CNNs are trained on large datasets containing both real and manipulated images. Through training, the network learns to identify subtle inconsistencies typically found in synthetic media, such as unnatural skin textures, irregular lighting, mismatched facial landmarks, or compression artifacts. These are features often introduced during the generation or tampering process using deepfake algorithms. Once trained, the CNN can accurately classify an input image as either "real" or "fake" by analyzing these extracted features. This classification result is then passed to the system's moderation module, which either flags the image or allows it to proceed, based on its authenticity. The use of CNNs ensures a high degree of precision and speed in image analysis, making them an ideal choice for the deepfake detection task in the Digital Purifier.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



#### **Generative Adversarial Networks (GANs)**

Generative Adversarial Networks (GANs) form the core of the deepfake detection algorithm due to their powerful ability to model and understand the structure of synthetic images. A GAN consists of two neural networks—the Generator and the Discriminator—that are trained simultaneously in a competitive setup. The Generator's task is to create fake images that closely resemble real ones, while the Discriminator's role is to evaluate and distinguish between real and generated (fake) images. During training, the Generator learns to improve the quality of its fake outputs by trying to "fool" the Discriminator, and the Discriminator becomes increasingly skilled at identifying forgeries by learning the subtle artifacts, distortions, and irregular patterns that typically appear in manipulated visuals. In the Digital Purifier, once this adversarial training is complete, the Discriminator network is utilized as the deepfake detector. It receives input images from users and analyzes them for signs of tampering, such as unnatural facial features, inconsistent lighting, or pixel-level abnormalities. Based on this analysis, it classifies the image as authentic or fake. The use of GANs in this context is particularly effective because the Discriminator has been trained on a wide range of realistic forgeries, making it adept at catching even highly convincing deepfakes. This ensures the Digital Purifier delivers reliable results in real-time, helping prevent the spread of misleading or deceptive media content.

#### **BERT (Bidirectional Encoder Representations from Transformers)**

BERT is employed as the core algorithm for detecting toxic or harmful text-based content. Developed by Google, BERT is a powerful language representation model that uses a transformer-based architecture to deeply understand the meaning of words in context. Unlike traditional NLP models that read text from left to right or right to left, BERT reads text bidirectionally, meaning it considers both the preceding and following words simultaneously. This enables it to grasp the full context of a sentence, even when the language used is sarcastic, indirect, or emotionally nuanced. When a user submits a comment, the text is first tokenized and converted into embeddings, which are then fed into the BERT model. Through multiple transformer layers, BERT processes these embeddings and captures relationships between words to determine whether the input is toxic. It can identify a range of toxic behaviors such as hate speech, abusive language, threats, and discriminatory remarks. Once the text is classified, the system automatically takes moderation actions, such as blocking the comment or logging it for review. BERT's deep contextual understanding and ability to handle complex sentence structures make it an ideal solution for detecting harmful content that might otherwise go unnoticed by simpler keyword-based filters. As a result, the Digital Purifier can offer real-time and highly accurate content moderation across diverse forms of user input.

#### V. PACKAGES

### OpenCV

(Open Source Computer Vision Library) is a library of programming functions mainly for real-time computer vision. Originally developed by Intel, it was later supported by Willow Garage, then Itseez (which was later acquired by Intel). The library is cross-platform and licensed as free and open-source software under Apache License 2. Starting in 2011, OpenCV features GPU acceleration for real-time operations. When it is integrated with various libraries, such as NumPy, python is capable of processing the opencv array structure for analysis. To Identify an image pattern and its various features we use vector space and perform mathematical operations on these features.

#### PyTorch

PyTorch is a deep learning library built on Python and Torch (a Lua-based framework). It provides GPU acceleration, dynamic computation graphs, and an intuitive interface for deep learning researchers and developers. PyTorch follows a "define-by-run" approach, meaning that its computational graphs are constructed on the fly, allowing for better debugging and model customization.

PyTorch uses dynamic graphs, allowing flexibility in model execution and debugging and provides an automatic differentiation engine that simplifies gradient computation for deep learning. It supports CUDA, allowing computations to be performed efficiently on GPUs.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



#### Natural Language Processing (NLP)

NLP algorithms are crucial for analyzing user-generated text and detecting toxic or harmful comments. These algorithms work by transforming raw text into structured data that can be understood by machine learning models. Techniques such as tokenization, stemming, lemmatization, and vectorization are used in the preprocessing phase to prepare the text. Advanced models like BERT are then employed to capture the contextual meaning of words in a sentence, allowing the system to accurately detect subtle forms of toxicity, including sarcasm, threats, and abusive language. Unlike traditional rule-based methods, NLP algorithms powered by deep learning can learn from vast datasets and adapt to complex linguistic patterns, making them highly effective for real-time moderation. As a result, NLP serves as a foundational component in ensuring that the platform remains safe, respectful, and free from harmful digital communication.

### VI. EXPERIMENTAL RESULTS&PERFORMANCE EVALUATION

The Digital Purifier system was evaluated based on the performance of its two core modules—Deepfake Detection and Toxic Comment Detection—using standard classification metrics including accuracy, precision, recall, F1-score, and latency (response time). These experiments validate the robustness and efficiency of the proposed system in real-world content moderation scenarios.

#### **Deepfake Detection Module**

The deepfake detection component, built using a CNN architecture enhanced with a GAN discriminator, was trained on a balanced dataset containing real and forged facial images. The system achieved the following performance:

- Accuracy: 94.2%
- Precision: 93.5%
- Recall: 92.8%
- F1-Score: 93.1%
- Average Response Time: ~1.3 seconds per image

The high precision and recall values indicate the system's strong ability to correctly identify fake images while minimizing false positives and false negatives. The model was particularly effective in detecting inconsistencies in facial geometry, unnatural skin textures, and lighting mismatches, which are commonly present in synthetic images generated by deepfake techniques.

### **Toxic Comment Detection Module**

The textual analysis module leverages the BERT language model, fine-tuned on datasets such as the Jigsaw Toxic Comment Classification Challenge. It was tested against a diverse range of inputs including offensive, sarcastic, and subtly toxic content. Performance metrics include:

- Accuracy: 92.6%
- Precision: 91.8%
- Recall: 90.4%
- F1-Score: 91.1%
- Average Response Time: ~1.5 seconds per comment

This module accurately detected various forms of toxic language such as insults, threats, identity-based hate, and profanity. The use of BERT allowed for contextual understanding, making the system capable of identifying toxicity even when phrased indirectly or humorously.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 1, July 2025



#### System Architecture



#### VII. LIMITATION

While The Digital Purifier offers an innovative and effective approach to combating digital misinformation and toxicity, it is not without its limitations. One major limitation lies in its dependency on pre-trained models such as BERT and GANs, which, although powerful, can struggle with domain-specific nuances, newer slang, or unseen manipulation techniques not present in the training data. As a result, false positives and false negatives may still occur, particularly in edge cases involving sarcasm, regional dialects, or subtle deepfake techniques. Additionally, the system currently supports only English-language text, which limits its effectiveness in multilingual or culturally diverse environments. Real-time processing is another area of constraint; the system currently handles user inputs in a batch or on-request format rather than live streams, making it less suited for high-speed, dynamic platforms like live chats or video feeds. Furthermore, while the system performs well on images, it does not yet support video deepfake detection or audio moderation, which are increasingly relevant in today's content landscape. Computational resource requirements for deep learning inference—especially GAN-based image analysis—may also present challenges for deployment on low-powered or edge devices. Lastly, despite its automated nature, the absence of a human-in-the-loop

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

#### Volume 5, Issue 1, July 2025



review system can reduce contextual judgment in borderline cases, potentially leading to user dissatisfaction. These limitations, while important to acknowledge, also offer valuable opportunities for future development and enhancement.

## VIII. FUTURE SCOPE

In the future, the Digital Purifier has strong potential for expansion beyond its current capabilities, laying the foundation for a highly adaptable and intelligent content moderation system. In the future, the system can be extended to support multilingual toxic comment detection using advanced models like mBERT or XLM-RoBERTa, enabling broader applicability in global and regional contexts. Similarly, video-based deepfake detection can be incorporated, allowing the system to analyze frame sequences and detect tampered video content, which is a growing concern on social media platforms. Another major enhancement involves real-time processing for both image and text moderation, enabling seamless integration into live chats, video conferencing, and interactive streams. With optimization techniques such as model pruning and quantization, the system can be made lightweight and deployable on mobile and edge devices, increasing accessibility. Moreover, integrating emotional tone and sentiment analysis could improve classification accuracy by distinguishing between genuine criticism, sarcasm, and actual abuse. Features like a moderator dashboard, user feedback system, appeal mechanism, and blockchain-based tamper-proof logging could improve transparency, user trust, and administrative control. Eventually, API integration with third-party platforms could make the Digital Purifier a plug-and-play tool for digital platforms seeking robust, AI-powered content protection. These future enhancements will transform the Digital Purifier into a scalable, multilingual, and ethically aware solution that addresses the evolving threats in the digital landscape.

#### **IX. CONCLUSION**

In an era where digital communication is ubiquitous, the integrity and safety of online content have become crucial. The Digital Purifier addresses these modern challenges by leveraging advanced AI models—GAN for detecting deepfakes and BERT for identifying toxic language—within a unified, efficient platform. Through the successful integration of these technologies into a Flask-based web application, the system effectively detects manipulated media and harmful text, enabling proactive moderation and promoting healthier digital interactions. The project demonstrates the feasibility of deploying cutting-edge machine learning models for real-world content verification and moderation tasks, while also emphasizing ethical responsibility and user privacy. Although current limitations exist, the Digital Purifier lays a solid foundation for future enhancements such as real-time moderation, multilingual support, mobile deployment, and improved explainability. Ultimately, this system not only showcases the power of AI in combating misinformation and hate speech but also serves as a blueprint for future innovations aimed at creating a safer, more trustworthy digital environment.

#### REFERENCES

[1]. Ian Goodfellow, Yoshua Bengio, and Aaron Courville.Deep Learning. MIT Press, 2016.

[2]. Vaswani, A., Shazeer, N., Parmar, N., et al. "Attention is All You Need." Advances in Neural Information Processing Systems (NeurIPS), 2017.

[3]. Devlin, J., Chang, M.W., Lee, K., & Toutanova, K. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv preprint arXiv:1810.04805, 2018.

[4]. Tschannen, M., Bachem, O., & Lucic, M. "Recent Advances in Autoencoder-Based Representation Learning." arXiv preprint arXiv:1812.05069, 2018.

[5]. Goodfellow, Ian, Bengio, Yoshua, and Courville, Aaron. Deep Learning. MIT Press, 2016. A foundational text on deep learning covering neural networks, optimization, unsupervised learning, and key concepts used in GANs.

[6]. Devlin,Jacob,Chang,Ming-Wei,Lee,Kenton, and Toutanova, Kristina."BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv:1810.04805, 2018. The original paper introduced BERT, a transformer-based model pre-trained on a large corpus, now widely used for text classification tasks like toxic comment detection.

Copyright to IJARSCT www.ijarsct.co.in



DOI: 10.48175/568





International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

#### Volume 5, Issue 1, July 2025



[7]. Vaswani, Ashish, et al. "Attention is All You Need." Advances in Neural Information Processing Systems (NeurIPS), 2017. This seminal paper introduces the transformer architecture which underpins BERT and revolutionized NLP.

[8]. Radford, Alec, et al. "Language Models are Unsupervised Multitask Learners." OpenAI GPT-2 paper, 2019. Describes GPT, another large-scale language model. Relevant for comparative context when discussing BERT.

[9]. O'Reilly, Sebastien Raschka. Machine Learning with PyTorch and Scikit-Learn: Develop machine learning and deep learning models with Python. Packt Publishing, 2022. Practical applications and code examples for implementing models like BERT and GANs in PyTorch.

[10].Chollet, François.Deep Learning with Python. Manning Publications, 2017



