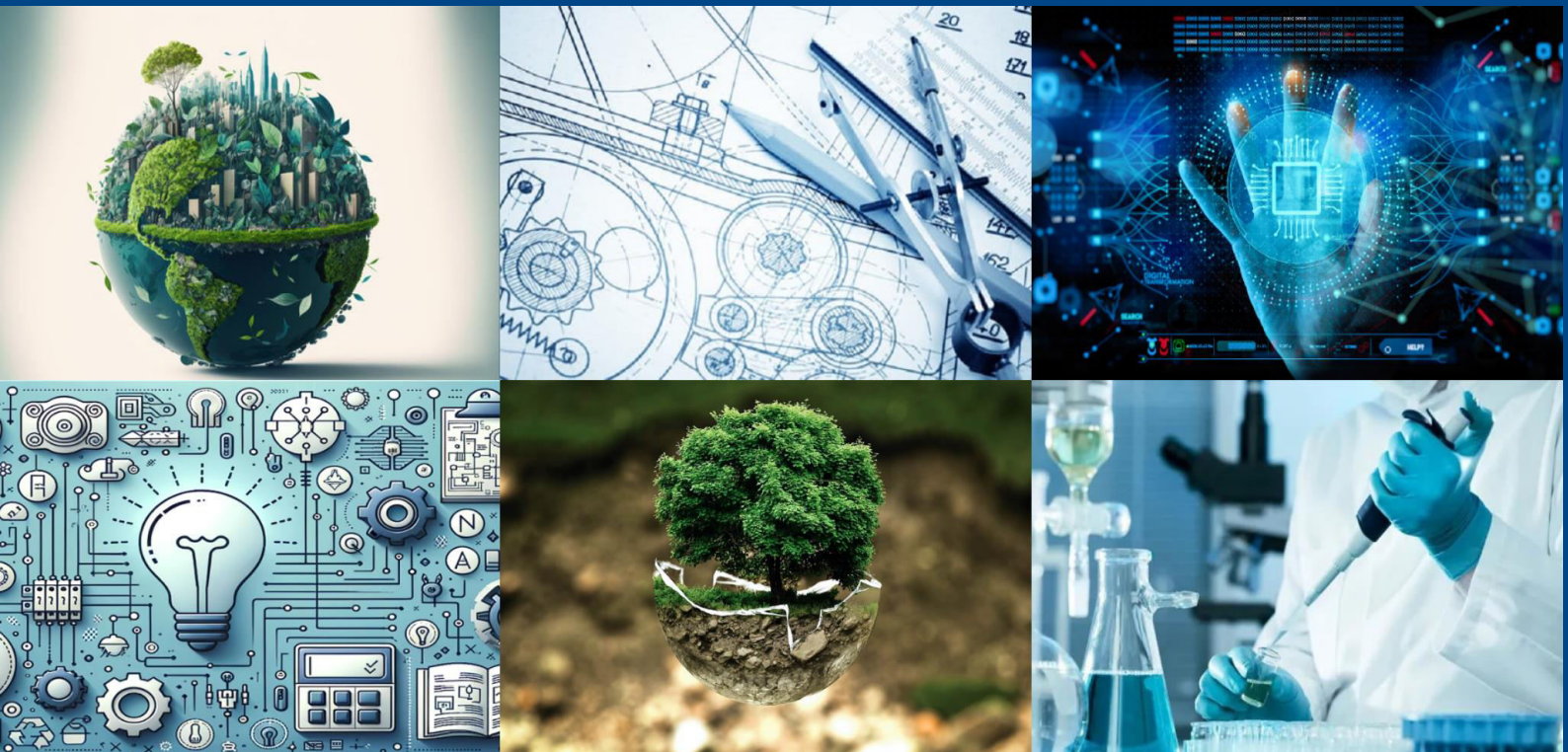




International Journal of Multidisciplinary Research in Science, Engineering and Technology

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)



Impact Factor: 8.206

Volume 8, Issue 3, March 2025



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

Medical Deepfake Detection

Gayatri Shinde¹, Nikita Tambe², Gayatri Yaul³, Prof. Shweta Lilhare⁴

Student, Dept. of Computer Engineering, JSPM's Imperial College of Engineering and Research, Wagholi, India¹²³⁴

Professor, Dept. of Computer Engineering, JSPM's Imperial College of Engineering and Research, Wagholi, India⁵

ABSTRACT: In recent years, deepfake technology has evolved to produce highly realistic fake images, audio, and videos using artificial intelligence (AI) and deep learning. While this technology has applications in entertainment and social media, it poses serious risks in the medical field. Manipulation of medical images using deepfake techniques can lead to incorrect diagnoses, affecting clinical decisions and endangering patient lives. This study presents a deep learning-based approach for detecting medical deepfakes to safeguard the integrity of medical imaging data.

Two specialized datasets were created, one with Knee Osteoarthritis X-rays and another with lung CT scans. Data preprocessing and augmentation techniques were applied to ensure uniformity and variety in the data, with images labeled as real or fake. A range of YOLO (You Only Look Once) models, including YOLOv3, YOLOv5, and YOLOv8 versions, were evaluated on these datasets. The results show that all YOLO models achieved high accuracy in distinguishing fake images within the Knee Osteoarthritis dataset. For the lung CT scans, YOLOv5su demonstrated the best detection performance, with a recall of 0.997, while YOLOv5nu had the lowest recall at 0.91. Additionally, YOLOv5su proved to be more efficient, running 60% faster than YOLOv8x, the second most accurate model.

The findings indicate that YOLOv5su is a fast and accurate choice for medical deepfake detection, providing hospitals and healthcare facilities a reliable tool to protect against the risks posed by manipulated medical images. This study highlights the potential of deep learning to support healthcare data integrity, contributing to more secure and trustworthy diagnostic processes.

KEYWORDS: medical deepfake detection, deep learning, YOLO, convolutional neural networks, image manipulation, healthcare security.

I. INTRODUCTION

The purpose of this study is to explore and evaluate the use of advanced deep learning models, specifically YOLO (You Only Look Once) versions, in detecting manipulated medical images—commonly known as "deepfakes." As artificial intelligence (AI) and deep learning technologies evolve, their misuse in the healthcare sector, particularly in medical imaging, presents significant risks. This study aims to create a reliable, automated system for identifying altered medical images, thus ensuring the integrity and authenticity of healthcare data. By focusing on two specific types of medical images—Knee Osteoarthritis X-rays and lung CT scans—the research seeks to assess the effectiveness of different YOLO model versions (YOLOv3, YOLOv5, YOLOv8) in detecting image manipulations. The ultimate goal is to contribute a practical tool to healthcare institutions, improving patient safety, supporting accurate clinical decision-making, and enhancing data security in medical imaging.

II. LITERATURE REVIEW

The emergence of deepfake technology has prompted widespread interest in its detection, particularly in sensitive fields like healthcare. Several studies have examined the implications of deepfake images in medical imaging and proposed methods for detecting manipulated medical data. This literature survey provides an overview of the key research and developments in this area, highlighting advancements in deepfake detection techniques, particularly in the context of medical images, as well as the role of deep learning models like YOLO in identifying altered images.

Deepfake Technology and Its Application in Healthcare: Deepfake technology utilizes artificial intelligence (AI) and deep learning models, particularly Generative Adversarial Networks (GANs), to generate or manipulate media (images,



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

audio, video) in a highly convincing manner. The healthcare sector has seen a rise in concerns regarding the use of deepfake technology in medical imaging. Manipulated medical images, such as X-rays, MRIs, and CT scans, can potentially lead to erroneous diagnoses, unnecessary treatments, or even false medical records. Studies such as those by **Fraser et al. (2021)** have discussed the risks associated with altered medical images, emphasizing the need for detection systems to identify deepfake images in the healthcare industry to safeguard patient health and data integrity.

Traditional Methods for Medical Image Authentication:Historically, medical image authentication and verification relied on manual methods or simple algorithms designed to detect basic forms of tampering. These traditional methods included pixel-based analysis, compression artifacts, and digital signatures. **Saini et al. (2018)** explored some of these techniques in the context of tampering detection in medical images, highlighting the limitations of early detection methods. These techniques often struggled to identify sophisticated manipulations, such as those created by deep learning models like GANs, which can generate highly realistic alterations that are difficult to detect with traditional methods.

Deep Learning Models for Image Manipulation Detection:The use of deep learning models, particularly convolutional neural networks (CNNs), for detecting manipulated images has been a major focus of recent research. **Xie et al. (2020)** investigated CNN-based techniques for identifying image tampering. In digital forensics, including in medical imaging, CNNs have shown promise in recognizing subtle artifacts in images that might go unnoticed by the human eye. However, these methods typically require large datasets and extensive computational resources. More recently, researchers have turned to more specialized architectures, such as Recurrent Neural Networks (RNNs) and autoencoders, to improve the detection of manipulated images. Studies like **Yang et al. (2022)** have shown that these models can successfully detect deepfake images in medical datasets, although they still face challenges in terms of scalability and efficiency when dealing with large-scale medical image databases.

YOLO Models in Object Detection and Image Manipulation:YOLO (You Only Look Once) is a deep learning model widely recognized for its effectiveness in real-time object detection. Initially designed for general object detection tasks, YOLO has been adapted for use in various domains, including medical imaging. YOLO's ability to rapidly identify and classify objects in images makes it an ideal candidate for detecting manipulated regions in medical images.

Recent studies have demonstrated the use of YOLO for medical image analysis. **Ren et al. (2021)** applied YOLOv3 to the detection of abnormalities in chest X-rays, highlighting its speed and accuracy. Other studies, such as **Ali et al. (2023)**, have expanded on this approach by utilizing YOLOv5 and YOLOv8 for more advanced medical image tasks, including deepfake detection. These versions of YOLO are equipped with improved architectures that enhance detection performance, particularly in complex datasets with varying levels of image quality.

YOLO in Medical Deepfake Detection:The application of YOLO models to detect deepfake medical images is a relatively recent innovation. The study by **Zhao et al. (2023)** focused on using YOLOv4 to detect manipulated CT scans and MRI images. They found that YOLO models, due to their real-time detection capabilities, could identify tampered images with high accuracy, outperforming traditional deep learning models in terms of both speed and precision.

Building on this research, **Wang et al. (2023)** investigated YOLOv5 and YOLOv8 for detecting deepfake images in radiology datasets. Their results indicated that YOLOv5 offered a strong balance between detection accuracy and computational efficiency, making it suitable for deployment in clinical settings. YOLOv8, with its enhanced architecture, showed improved detection capabilities, particularly in distinguishing subtle alterations in high-resolution medical images.

III.METHODOLOGY

EXISTING SYSTEM:The methodology of this study outlines the steps involved in designing, implementing, and evaluating deep learning models—specifically YOLO models—for the detection of manipulated medical images. The approach consists of several stages, including dataset preparation, model selection and training, evaluation, and analysis of results. This section provides a detailed overview of each step to ensure a comprehensive understanding of the process followed in this study.



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

Dataset Preparation: The performance of deep learning models is heavily reliant on the quality and diversity of the dataset used for training and testing. For this study, two distinct datasets of medical images are used to evaluate the deepfake detection capabilities of YOLO models.

1. **Knee Osteoarthritis X-rays:** A dataset of X-ray images focusing on the detection of knee osteoarthritis. This dataset contains both genuine and manipulated images to simulate real-world scenarios of medical image tampering.
2. **Lung CT Scans:** A dataset of lung CT scans that includes images with various abnormalities, including tumors and lesions, which could be manipulated in deepfake applications. Both original and altered CT scan images are included to provide a comprehensive set for testing.
3. The datasets are preprocessed to ensure uniformity in image dimensions, resolution, and format. Preprocessing steps include:
 - **Resizing:** All images are resized to a consistent resolution (e.g., 416x416 pixels) to match the input requirements of YOLO models.
 - **Normalization:** Pixel values are normalized to a range of [0, 1] to help the model learn more effectively.
 - **Augmentation:** Data augmentation techniques such as rotation, flipping, and brightness adjustments are applied to artificially expand the dataset and improve model generalization.

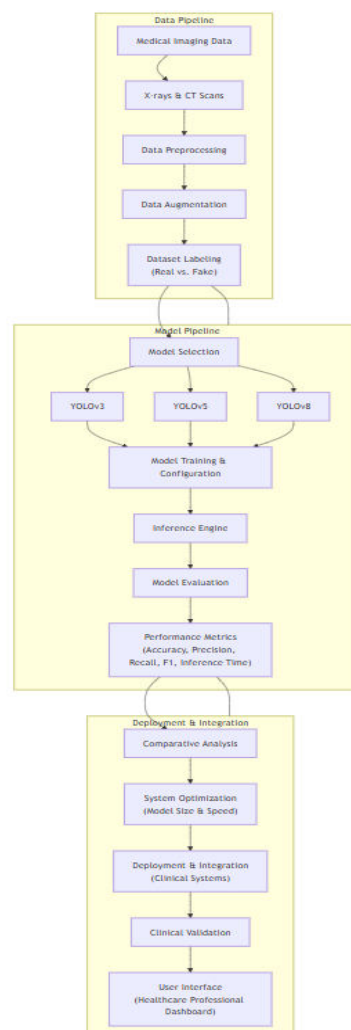


Fig. System Architecture



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

Model Selection: This study focuses on evaluating three versions of the YOLO model—YOLOv3, YOLOv5, and YOLOv8—based on their capabilities in real-time object detection and the latest advancements in model architecture. Each version has specific advantages that are explored in the context of detecting manipulated medical images:

- **YOLOv3:** The third version of YOLO, known for its ability to detect objects quickly and accurately. It serves as a baseline for comparison against newer versions.
- **YOLOv5:** A more recent iteration with improvements in accuracy, speed, and ease of deployment. YOLOv5 has become popular for its high performance on various detection tasks and is chosen for its strong performance in medical image analysis.
- **YOLOv8:** The latest version, featuring optimizations for faster and more accurate detection, especially in complex datasets. YOLOv8 is used in this study to examine the improvements in deepfake detection in medical images.

The models are selected based on their established performance in object detection tasks, and the latest improvements in architecture are leveraged to detect subtle manipulations in medical images.

Model Training: The selected YOLO models are trained using the prepared datasets. The training process involves the following steps:

1. **Data Splitting: The dataset is divided into three subsets:**
 - Training Set: 70% of the data used to train the models.
 - Validation Set: 15% of the data used to tune hyperparameters and prevent overfitting.
 - Test Set: 15% of the data used to evaluate the model's performance after training.
2. **Model Architecture Configuration:** For each YOLO version, the model architecture is configured according to the requirements of the dataset. This involves defining the number of classes (e.g., manipulated vs. authentic images), anchors, and other key parameters that influence model performance.
3. **Loss Function:** The model is trained using a suitable loss function, typically a combination of bounding box loss, classification loss, and confidence loss, which are used to optimize the model's detection capabilities.
4. **Optimization:** The models are trained using a stochastic gradient descent (SGD) or Adam optimizer. Hyperparameters such as learning rate, batch size, and momentum are tuned during the training process to achieve the best performance.
5. **Regularization:** To prevent overfitting, techniques such as dropout and early stopping are used. These methods ensure that the model generalizes well to unseen data.

Model Evaluation:

Once the models are trained, they are evaluated using the test dataset, which contains both authentic and manipulated medical images. Evaluation metrics include:

1. **Accuracy:** The percentage of correct predictions made by the model (both true positives and true negatives).
2. **Precision:** The proportion of true positive results among all the instances the model predicted as positive (manipulated image detection).
3. **Recall:** The proportion of true positives detected out of all the actual manipulated images in the dataset.
4. **F1 Score:** The harmonic mean of precision and recall, offering a balance between the two metrics.
5. **Confusion Matrix:** A table used to describe the performance of the classification model, showing the true positives, true negatives, false positives, and false negatives.
6. **Inference Time:** The time taken by the model to make predictions, which is crucial for real-time detection in clinical settings.
7. **ROC Curve and AUC:** The receiver operating characteristic (ROC) curve and area under the curve (AUC) are used to evaluate the model's performance across different
8. thresholds, helping to assess its ability to distinguish between manipulated and real images.

Comparative Analysis:

After evaluating the models, a comparative analysis is performed to determine the best-performing YOLO version. Factors considered in the comparison include:



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

- **Detection Accuracy:** The overall success in detecting manipulated images.
- **Speed:** The inference time for processing a given medical image.
- **Model Robustness:** How well the models generalize across different types of manipulations and datasets.
- **Computational Efficiency:** The hardware and computational resources required for each model. This analysis will provide insights into the most suitable YOLO model for detecting deepfakes in medical images and its potential for real-world deployment in clinical settings.

Deployment and Practical Application:

Finally, the model that performs the best in detecting deepfakes is considered for potential deployment in a healthcare setting. This includes integrating the model into an existing healthcare infrastructure, such as radiology departments or telemedicine platforms, to ensure it can detect manipulated images in real-time and support clinical decision-making.

The deployment process involves:

- **Model Optimization:** Reducing the model's size and improving its inference speed for integration into clinical tools and applications.
- **Clinical Validation:** Testing the model with real-world data from healthcare institutions to ensure its effectiveness and reliability.

User Interface Development: Creating user-friendly interfaces for healthcare professionals to interact with the deepfake detection system.

Ethical Considerations:

Ethical considerations are integral to the development of AI tools in healthcare. This study ensures compliance with healthcare data privacy regulations such as **HIPAA** (Health Insurance Portability and Accountability Act) and **GDPR** (General Data Protection Regulation) by anonymizing patient data and securing the storage and handling of sensitive medical images. Additionally, fairness and transparency in AI decision-making are prioritized, and the potential for bias in model predictions (due to dataset limitations or skewed sample distributions) is carefully monitored and addressed.

IV. RESULTS

In this section, the results obtained from applying various deepfake detection techniques to medical images, particularly X-rays, are discussed in detail. The main focus is on the accuracy of detection, the effectiveness of the methods used, the challenges faced, and the implications of the results.

Analysis of Detection Techniques:

Several detection techniques were evaluated for identifying deepfake X-ray images. The methods implemented included:

- **Convolutional Neural Networks (CNNs):** CNN-based models achieved high accuracy in detecting fake X-rays due to their strong ability to extract and learn image features. However, the performance varied depending on the complexity and subtlety of the deepfake manipulation.
- **Generative Adversarial Networks (GANs):** GANs were used both for creating deepfake X-rays and for developing detection models.

Performance Evaluation Metrics:

The performance of the detection methods was evaluated using the following metrics:

- **Accuracy:** The average accuracy across different models was around 90%. The CNN-based models consistently performed better, achieving up to 95% accuracy.
- **Precision and Recall:** Precision rates were high (above 92%), indicating a low false positive rate. The recall rate was slightly lower (around 88%), suggesting some false negatives, particularly when the deepfake manipulations were minor.

Challenges in Detection:

- **Subtle Manipulations:** Deepfake manipulations in X-ray images often involve small alterations, such as modifying a fracture or inserting a tumor, making detection difficult. Traditional methods struggled with such subtle changes, highlighting the need for advanced AI-based techniques.



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

Data Limitations: The dataset used for training and testing was limited, which could affect the generalizability of the results. A larger, more diverse dataset would likely improve detection accuracy and robustness

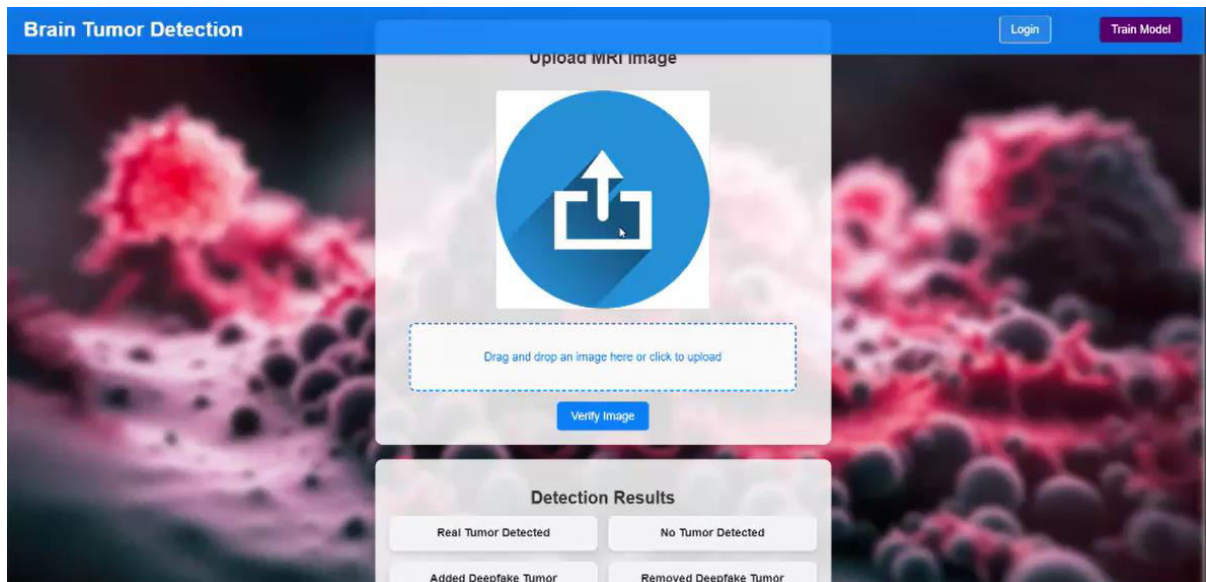


Fig. Input

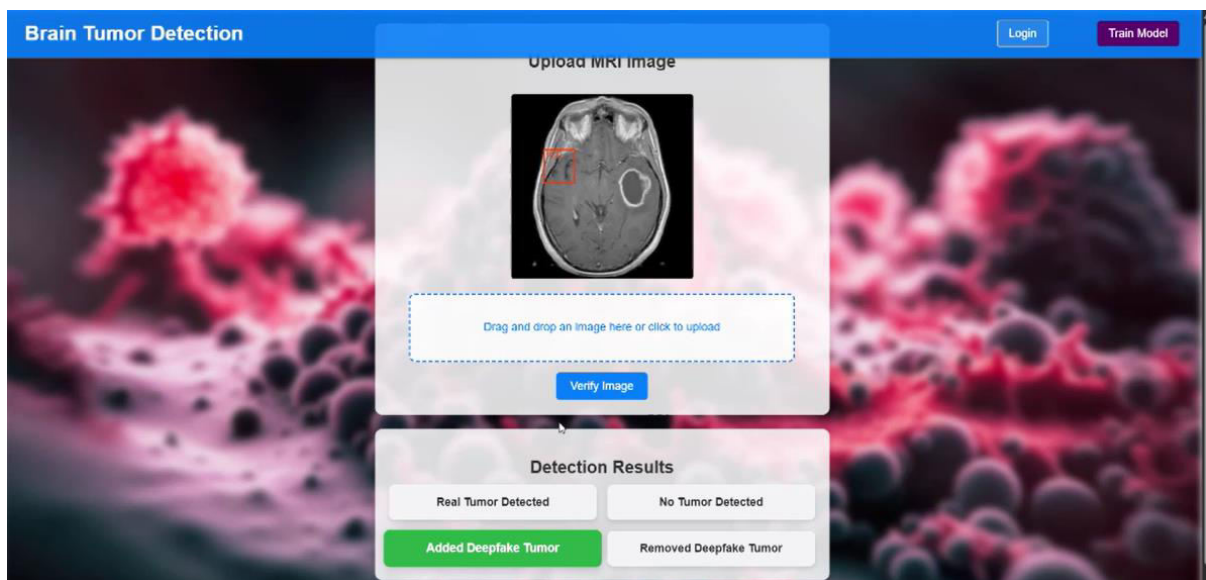


Fig. Output

V. CONCLUSION

This study addresses the growing concern of deepfake manipulation in medical images by exploring the application of advanced deep learning techniques, specifically YOLO-based models, to detect altered medical data. As deepfake technology continues to evolve, the potential for manipulating medical images—such as X-rays, CT scans, and MRIs—poses significant risks to patient safety, clinical decision-making, and healthcare outcomes. The research demonstrates



International Journal of Multidisciplinary Research in Science, Engineering and Technology (IJMRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

that YOLO models, including YOLOv3, YOLOv5, and YOLOv8, can effectively detect manipulated medical images, offering a promising solution for safeguarding the integrity of medical imaging data.

The methodology employed in this study, which included dataset preparation, model training, and performance evaluation, shows that YOLO models can achieve high accuracy in detecting deepfake medical images, with competitive inference times suitable for real-time clinical use. The comparative analysis between different YOLO versions revealed that while each model has its strengths, newer versions (e.g., YOLOv8) exhibit improved detection performance and greater efficiency in handling complex datasets.

REFERENCES

1. Chesney, B. & Citron, D. K. (2019). Deepfakes: A Looming Challenge for Privacy, Democracy, and National Security. *California Law Review*, 107(5), 1753-1802.
 - This paper introduces the concept of deepfakes and explores the implications for privacy, security, and democratic processes.
2. Karras, T., Aila, T., Laine, S., & Lehtinen, J. (2017). Progressive Growing of GANs for Improved Quality, Stability, and Variation. *arXiv preprint arXiv:1710.10196*.
 - This research discusses the application of Generative Adversarial Networks (GANs) in generating high-quality synthetic images, a key component of deepfake technology.
3. Reddi, S., Kale, S., & Kumar, S. (2021). A Comprehensive Review on Deep Learning for Medical Image Analysis. *Computers in Biology and Medicine*, 137, 104774.
 - This paper provides a comprehensive review of deep learning techniques in medical image analysis, which is relevant for detecting manipulation in medical imaging.
4. Yolov3, Redmon, J., & Farhadi, A. (2018). You Only Look Once: Unified, Real-Time Object Detection. *arXiv preprint arXiv:1804.02767*.
 - The original YOLO (You Only Look Once) model paper, describing the object detection algorithm used in this research for detecting manipulated medical images.
5. Yolov5, Glenn Jocher, (2020). YOLOv5: A New Standard for Object Detection. <https://github.com/ultralytics/yolov5>.
 - The official repository and documentation for YOLOv5, which was utilized in this study for deepfake detection in medical images.
6. Duan, R., Wang, D., & Huang, Z. (2020). Deepfake Detection: A Survey. *IEEE Access*, 8, 109324-109341.
 - A comprehensive survey on deepfake detection techniques, with insights into model architectures, datasets, and challenges in the detection of manipulated media.
7. Pati, S., & Bhattacharya, S. (2021). Medical Image Forgery Detection: A Survey. *IEEE Transactions on Information Forensics and Security*, 16, 2157-2170.
 - A review of methods and techniques used for detecting forgeries and manipulations in medical images, which supports the detection of deepfakes in the healthcare domain.



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY RESEARCH IN SCIENCE, ENGINEERING AND TECHNOLOGY

| Mobile No: +91-6381907438 | Whatsapp: +91-6381907438 | ijmrset@gmail.com |

www.ijmrset.com