

**Mini Project Report**  
On  
**IMAGE FORENSICS**  
**Detecting Manipulated Images as Real vs AI Edited or  
Generated**

Submitted in partial fulfilment of the requirements for the award of degree of

**Bachelor of Engineering**  
In  
**Computer Science and Engineering**

By  
**M Vedanth VNN    1608-20-733-014**  
**Tejas JM            1608-20-733-005**  
**Shashank Adepu    1608-20-733-003**

Under the guidance of

**Mrs Priyanka Madhiraju**  
Assistant Professor



**Department of Computer Science and Engineering**  
**Matrusri Engineering College**

**Accredited by NBA & NAAC**

(Affiliated to Osmania University, Approved by AICTE)

Saidabad, Hyderabad-500059

(2022-2023)

# **Department of Computer Science and Engineering**

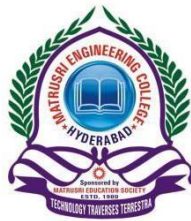
## **Matrusri Engineering College**

**Accredited by NBA & NAAC**

(Affiliated to Osmania University, Approved by AICTE)

Saidabad, Hyderabad-500059

(2022-2023)



### **CERTIFICATE**

This is to Certify that A Mini Project report entitled “**IMAGE FORENSICS Detecting Manipulated Images as Real vs AI Edited or Generated**” is being submitted by M.Vedanth VNN (1608-20-733-014), Tejas JM (1608-20-733-005) and Shashank Adepu (1608-20-733-003) in partial fulfilment of the requirement of the award for the degree of Bachelor of Engineering in “Computer Science and Engineering” O.U., Hyderabad during the year 2022-2023 is a record of bonafide work carried out by them under my guidance. The results presented in this project have been verified and are found to be satisfactory.

**Project Guide**

**H.O.D.**

**Mrs. Priyanka Madhiraju**

Assistant Professor,

Dept of C.S.E.

**Dr. P. Vijayapal Reddy**

Professor & Head,

Dept. of CSE

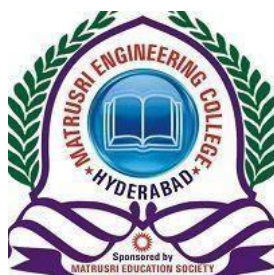
**External Examiner(s)**

# MATRUSRI ENGINEERING COLLEGE

Accredited by NBA & NAAC

**SAIDABAD – 500059**

**Department of Computer Science and Engineering**



## DECLARATION

We **M VEDANTH VNN** bearing **H.T.No:1608-20-733-014**, **Tejas JM** bearing **H.T.No:1608-20-733-005** and **Shashank Adepu** bearing **H.T.No:1608-20-733-003** hereby certify that the major project report entitled “**IMAGE FORENSICS : Detecting Manipulated Images as Real vs AI Edited or Generated**” is submitted in the partial fulfillment of the required for the award of the degree of **Bachelor of Engineering in Computer Science and Engineering**.

This is a record of bonafide work carried out by us under the guidance of Mrs Priyanka Madhiraju, Assistant Professor, Matrusri Engineering College, Saidabad. The results embodied in this report have not been reproduced/copied from any source. The results embodied in this report have not been submitted to any other university or institute for the award of any other degree or diploma.

**M Vedanth VNN      1608-20-733-014**

**Tejas JM              1608-20-733-005**

**Shashank Adepu      1608-20-733-003**

## **ACKNOWLEDGEMENT**

We would like to express our heartfelt gratitude to all those who have contributed to the successful completion of the project titled " IMAGE FORENSICS : Detecting Manipulated Images as Real vs AI Edited or Generated."

First and foremost, We are deeply grateful to our esteemed supervisor, Mrs. Priyanka Madhiraju, Assistant Professor in the Department of Computer Science and Engineering, for her unwavering guidance, invaluable insights, and unwavering support that have been the driving force behind the success of this project. Her expert assistance in data collection, technical advice, and all the other forms of support she provided have truly made a significant impact on the project's accomplishments. Her expertise and encouragement have been instrumental in shaping the direction of this work.

We would also thank Dr. P. Vijaya Pal Reddy (HOD, CSE department ) and Dr. D.Hanumanth Rao, Principal, Matrusri Engineering College, Saidabad, Hyderabad for their precious suggestions, motivation and co-operation

We are thankful to the faculty and staff of Matrusri Engineering College who provided a conducive academic environment and access to essential resources, enabling us to carry out this research.

Lastly, we would like to extend our thanks to our friends for their understanding, encouragement, and motivation, which kept us inspired and focused during the course of this project.

# **ABSTRACT**

The proliferation of artificial intelligence (AI) and image editing technologies has led to an increased prevalence of digitally manipulated and generated images across various domains. This project aims to address the need of a system that can accurately discern the authenticity of images by developing a software solution to classify images into real and AI-edited/generated categories. Leveraging a diverse dataset with labeled examples, the system will utilize advanced image analysis techniques and deep learning-based features to train a reliable image classification model. The project will involve comprehensive data analysis, algorithm selection, and model optimization, evaluating the model's performance using appropriate metrics.

In addition to image classification, the project includes designing a user-friendly web interface for image classification. The system's architecture will prioritize efficiency, security, scalability, and performance. UML diagrams, such as use case diagrams, activity diagrams, class diagrams, sequence diagrams, and component diagrams, will be used to illustrate the system's structure and functionality.

The proposed software solution has significant implications in media forensics, content verification, and security domains. By accurately classifying images, users can confidently determine the authenticity of visual content, making informed decisions. The project's outcomes will contribute to the field of image analysis, classification, and authenticity verification, paving the way for further research and advancements in this area.

## LIST OF FIGURES

S.No.	Figure No.	Figure Description	Pg.No.
1.	Figure 4.1	System Architecture	10
2.	Figure 4.2	CNN model summary	13
3.	Figure 4.3	Sample generates image from GAN	14
4.	Figure 4.4	Use case diagram	17
5.	Figure 4.5	Component Diagram	17
6.	Figure 4.6	Activity Diagram	18
7.	Figure 4.7	Class Diagram	18
8.	Figure 4.8	Sequence Diagram	19
9.	Figure 4.9	Dataflow diagram	20
10.	Figure 4.10	Fashion MNIST dataset	21
11.	Figure 4.11	Sample CNN Architecture	23
12.	Figure 4.12	Sample GAN Architecture	23
13.	Figure 5.1.a	Training the GAN	27
14.	Figure 5.1.b	Training the GAN	28
15.	Figure 5.2	Training the CNN	29
16.	Figure 6.1	Unit testing	35
17.	Figure 7.1	Images generated by GAN	39
18.	Figure 7.2	Fake image prediction by CNN	39
19.	Figure 7.3	Real image prediction by CNN	40

## LIST OF TABLES

<b>S. No.</b>	<b>Table No.</b>	<b>Description</b>	<b>Pg. No.</b>
1.	Table 2.1	Observations from base papers	6
2.	Table 6.1	Test Case – 01	33
3.	Table 6.2	Test Case – 02	34
4.	Table 6.3	Test Case -03	35

# CONTENTS

Declaration	iii
Acknowledgement	iv
Abstract	v
List of Figures	vi
List of Tables	vii
1.INTRODUCTION	1
1.1 Importance	1
1.2 Background	2
1.3 Problem definition	2
1.4 Objective	3
2. LITERATURE SURVEY	4
2.1 Base Paper 1	6
2.2 Base Paper 2	6
2.3 Base Paper 3	7
2.4 Base Paper 4	8
3. PROBLEM DEFINITION	9
4. DESIGN	10
4.1 System Architecture	10
4.2 Module Specification	11
4.2.1 Data Analysis	11
4.2.2 Training and building the model	12
4.2.3 Result	12
4.3 Proposed Methodology	14
4.3.1 Analyzing the data	14
4.3.2 Choosing the algorithms	14
4.3.3 Diagnosis	15
4.4 UML diagrams	15
4.4.1 Use case Diagram	16
4.4.2 Component Diagram	17
4.4.3 Activity Diagram	17
4.4.4 Class Diagram	18
4.4.5 Sequence Diagram	19
4.4.6 Dataflow Diagram	20



4.5 Dataset	20
4.6 Algorithms	22
5.IMPLEMENTATION	25
5.1 Software Requirements	25
5.2 Procedure	27
5.3 Sample Code	30
6.TESTING	35
6.1 Unit Testing	35
6.2 Test Cases	36
7.RESULTS	39
8. CONCLUSION & FUTURE SCOPE	41
9.REFERENCES	44

# 1.INTRODUCTION

In today's digital world, images play a significant role in our lives. However, with the increasing availability of powerful editing tools and AI technologies, it has become easier for people to manipulate or create images that may not be entirely truthful. This project aims to address the need for a way to distinguish between real images, AI-edited images, and generated images.

Imagine you come across an image on the internet or social media, and you're not sure if it's a genuine photograph or if it has been altered or even created by a computer. It's important to have a method to determine the authenticity of such images. Why is this important? Well, in a world where images hold significant influence, it's crucial to be able to trust the images we see.

## 1.1 Importance

The " IMAGE FORENSICS : Detecting Manipulated Images as Real vs AI Edited or Generated " project holds paramount importance in today's digital landscape, where the proliferation of AI and image editing technologies has fuelled the dissemination of manipulated and fake visual content. With the increasing risk of misinformation and fraudulent activities, the proposed software solution provides a critical tool for accurately discerning between authentic and manipulated images. Its application spans across various domains, including media verification, legal investigations, content moderation, and digital forensics, where image authenticity plays a pivotal role in establishing trust and credibility. By empowering users with a reliable image classification model, the project enhances digital security, combats misinformation, and fosters a safer and more trustworthy digital environment. Furthermore, the project's contribution to the advancement of AI-driven image analysis techniques has broader implications for research, education, and innovation in the fields of computer vision and deep learning.

## **1.2 Background:**

The consequences of misinformation and the potential misuse of manipulated images are significant, ranging from reputational damage to serious legal implications. Detecting and verifying the authenticity of images has become an urgent need in fields such as journalism, law enforcement, content moderation, and digital forensics. Traditional manual methods for image verification are often time-consuming, resource-intensive, and prone to human error, necessitating the development of automated and reliable image forensics solutions

In response to these challenges this project aims to address the critical need for an automated system capable of accurately discerning between real and AI-edited/generated images. Leveraging the power of deep learning and advanced image analysis techniques, the project seeks to provide users with a robust and efficient software solution for image classification and authenticity verification. By building such a system, the project aims to contribute to the wider efforts of combating misinformation, enhancing digital security, and promoting trust in visual content across the digital landscape.

## **1.3 Problem definition:**

In the age of rapid technological advancements, the prevalence of digitally manipulated and generated images has reached unprecedented levels, leading to the proliferation of misinformation and fake content across various platforms. The lack of automated systems to discern image authenticity has become a critical challenge in the digital landscape, hampering efforts to combat deceptive practices and safeguard trust in visual content.

Consequently, there is an increasing demand for an automated and reliable solution that can accurately classify images as either real or AI-edited/generated. Such a solution would play a pivotal role in content verification, journalism, legal investigations, and digital forensics, empowering professionals and users to make informed decisions based on authentic visual evidence.

The project aims to tackle this pressing issue by developing a cutting-edge software system that leverages advanced image analysis and deep learning techniques to classify images with high precision and efficiency. The project's success will contribute to combating misinformation, enhancing digital security, and promoting credibility in the face of an ever-evolving landscape of digitally altered content

#### **1.4 Objective:**

The objective of this project is to develop an advanced software system that leverages state-of-the-art image analysis and deep learning techniques for accurate image classification, distinguishing between authentic and AI-edited/generated visual content. By achieving high precision in verification and real-time processing capabilities, the system aims to combat misinformation and fake content, offering a valuable tool for content verification, journalism, legal investigations, and digital forensics. With a user-friendly web interface facilitating easy image upload and result access, the system's accessibility extends to various users, contributing to enhanced trust and credibility in the face of a rapidly evolving landscape of digitally altered content.

## 2.LITERATURE SURVEY

YEAR	AUTHOR	TITLE	METHODOLOGY
2018	He Huang, Philip S. Yu and Changhu Wang	An Introduction to Image Synthesis with Generative Adversarial Nets	This paper presents an in-depth introduction to Generative Adversarial Networks (GANs) and their applications in image synthesis. It discusses limitations of the original GAN model, leading to modified frameworks like BiGAN and ALI. Various GAN applications are explored, including text-to-image synthesis and image-to-image translation. Specific GAN models like SS-GAN, CycleGAN, DualGAN, and DistanceGAN are highlighted. The paper suggests future directions for improving GANs in image synthesis tasks, making it a valuable resource for researchers and practitioners.
2014	Ian J. Goodfellow, Jean Pouget-Abadie*, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair† , Aaron Courville, Yoshua Bengio‡	Generative Adversarial Nets	The document presents a novel generative model estimation procedure using a discriminative model D in a minimax two-player game. It introduces the Kullback–Leibler and Jensen–Shannon divergences for training G. The guide covers summarization

			techniques for documents, addressing FAQs about language support and multilanguage summarization. It offers tips on summary writing, focusing on understanding the original text, coherence, and avoiding self-plagiarism.
2020	Jie Gui, Zhenan Sun, Yonggang Wen, Dacheng Tao, Jieping Ye	A Review on Generative Adversarial Networks: Algorithms, Theory, and Applications	This document focuses on GAN applications in image synthesis and addresses a significant gap in the literature regarding GANs loss-variant, evaluation metrics, and stable training. It explores various GAN models, advantages, drawbacks, and specific use cases in image super-resolution, captioning, and more.
2020	Pourya Shamsolmoali , Masoumeh Zareapoor , Eric Granger , Huiyu Zhou , Ruili Wang , M. Emre Celebi and Jie Yang	Image Synthesis with Adversarial Networks: a Comprehensive Survey and Case Studies	The document explores the discriminator's role, regularization for image-to-image translation, and applications in various fields, including semi-supervised learning. It also discusses mode collapse issues and proposes future research directions for GANs, acknowledging their limitations in generating discrete data directly.

Table 2.1 Observations from base papers

## **2.1 Base Paper 1 :**

### **"An Introduction to Image Synthesis with Generative Adversarial Nets"**

The paper provides an in-depth introduction to Generative Adversarial Networks (GANs) and their applications in image synthesis. GANs are powerful deep learning models that consist of two neural networks, the generator and the discriminator, engaged in a competitive training process to produce realistic synthetic data, particularly images. It explores the GAN structure, where a generator and a discriminator neural network compete against each other to generate synthetic data samples and distinguish between real and synthetic samples, respectively.

The paper discusses the limitations of the original GAN model, which led to the development of modified frameworks like BiGAN and ALI that incorporate an encoder for mapping data samples into latent feature space. Additionally, various applications of GANs are explored, such as text-to-image synthesis, image-to-image translation, and surfacenormal map reconstruction. Specific GAN models like SS-GAN, CycleGAN, DualGAN, and DistanceGAN are highlighted for their excellence in these tasks. The concept of "equivariance," preserving distance information between data samples in both the original and latent spaces, is also discussed, with the mention of TP-GAN as a relevant model. The paper concludes by suggesting future directions, such as exploring combinations of loss functions and preserving similarity information between images, to further improve GANs' performance in image synthesis and translation tasks. Overall, this comprehensive review serves as a valuable resource for researchers and practitioners interested in GANs' applications and advancements

## **2.2 Base Paper 2 :**

### **"Generative Adversarial Nets"**

The document discusses a new generative model estimation procedure in the context of machine learning. It introduces a discriminative model  $D$  that estimates the probability that a sample came from the training data rather than from the generative model  $G$ . The training procedure for  $G$  is designed to maximize the probability of  $D$  making a mistake, which corresponds to a minimax two-player game .

The document further explains that early in learning, when  $G$  is poor,  $D$  can reject samples with high confidence because they are clearly different from the training data. In this case,  $\log(1 - D(G(z)))$  saturates. Therefore, instead of training  $G$  to minimize  $\log(1 - D(G(z)))$ , the authors propose a different approach. Introduces the concept of the Kullback–Leibler divergence and the Jensen–Shannon divergence between the model’s distribution and the data generating process, represented as  $C(G)$ .

The document then provides a guide on how to summarize a document, including providing a file, specifying the output format, starting the summarization, and getting the result. It also answers FAQs about the number of languages supported and the possibility of summarizing multilanguage documents.

Finally, the document offers tips on summary writing, emphasizing the importance of understanding the original text, including key ideas, and checking the summary for coherence and logic. It also mentions the possibility of self-plagiarism and provides a framework for summarizing a document.

## **2.3 Base Paper 3 :**

### **“A Review on Generative Adversarial Networks: Algorithms, Theory and Applications”**

This document is a comprehensive review of Generative Adversarial Networks (GANs) and their applications in image synthesis. It identifies a significant gap in the literature, specifically the lack of a detailed collection of GANs loss-variant, evaluation metrics, remedies for diverse image generation, and stable training. The paper aims to address this gap by providing an extensive analysis of various GAN models, discussing their advantages, drawbacks, and specific use cases.

The document covers a wide range of GAN applications, including image super-resolution, image captioning, image inpainting, and text-to-image translation. Previous studies have been limited to discussing GAN architectures and algorithms without delving into the formation details and specific characteristics of each model. The review also explores GANs' applications in anomaly detection, deep learning, medicine, deraining



methods, and unsupervised spam detection, along with the challenges faced, such as mode collapse, vanishing gradient problems, and convergence difficulties. Moreover, the document provides a detailed explanation of various GAN models, such as InfoGAN and their application in image-to-image translation. It also discusses datasets used for evaluation, ranging from MNIST to ImageNet and more specialized datasets like Van Gogh and DSLR. Additionally, the paper explores different approaches in image synthesis, including ControlGAN, Structure-GAN, and Photo-Sketch Synthesis, as well as supervised and unsupervised image-to-image translation scenarios. Overall, this review aims to offer a comprehensive understanding of GANs and their potential in image synthesis, while identifying research challenges and future directions in the field.

## **2.4 Base Paper 4 :**

### **“Image Synthesis with Adversarial Networks: a Comprehensive Survey and Case Studies”**

This comprehensive review on Generative Adversarial Networks (GANs) explores their evolution and applications. It begins by tracing the progress of GANs before 2017 and introduces architecture variants. Emphasizing the importance of variance reduction methods, the document presents representative GAN variants like InfoGAN, which decomposes input noise into  $z$  and  $c$ . It also covers the discriminator's role and discusses pix2pix regularization for image-to-image translation.

Moving on, the review delves into the training of GANs, exploring objective functions, skills, and structure. It introduces MDGAN, which employs an encoder to produce the latent variable  $z$  for the generator, and discusses spectral normalization in BigGANs and SAGAN. The document also explores the application of GANs in semi-supervised learning, mode collapse issues, and the use of Fréchet distance for comparison. Additionally, it highlights GANs' applications in various fields, including supervised, unsupervised, and semi-supervised learning, while acknowledging their limitations in generating discrete data directly. The review concludes by proposing future research directions for GANs.

### **3.PROBLEM DEFINITION**

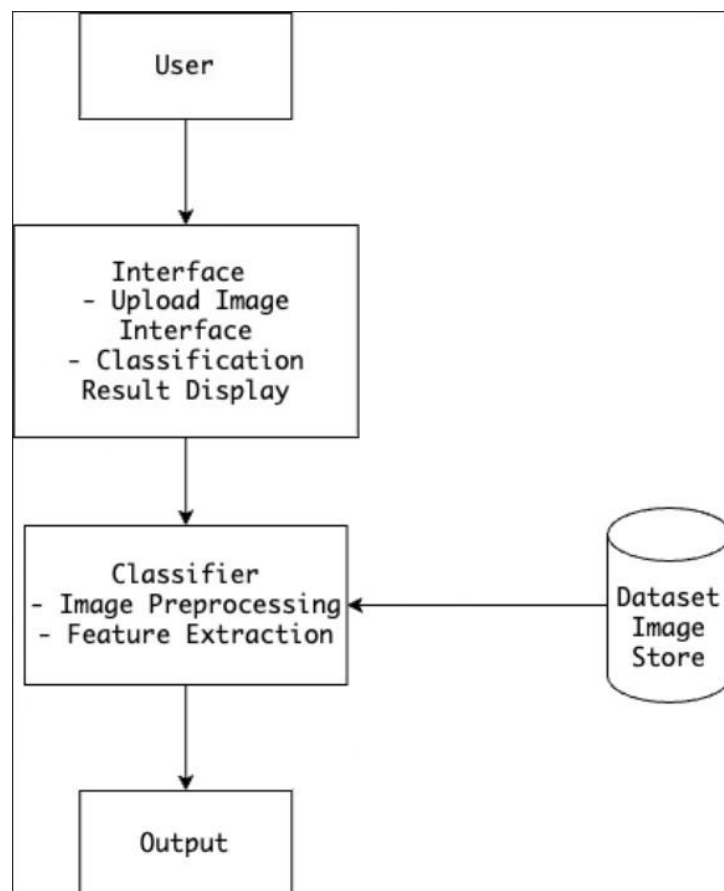
In the digital age, the widespread use of artificial intelligence (AI) and sophisticated image editing tools has given rise to an alarming surge in manipulated and fake visual content. This growing prevalence of digitally altered images poses a significant challenge in various domains, where the authenticity of visual evidence holds paramount importance. With the rapid advancements in artificial intelligence (AI) and image editing technologies, the creation and dissemination of digitally altered images have become increasingly prevalent, leading to the spread of misinformation and fake content. Detecting and verifying the authenticity of images manually is time-consuming, labor-intensive, and often ineffective in handling the scale and sophistication of modern manipulation techniques.

Therefore, the objective of this project is to develop an advanced software solution that employs cutting-edge image analysis and deep learning techniques to automatically classify images with high precision and efficiency. By doing so, the project aims to provide professionals and users with a reliable and scalable tool for detecting manipulated images, ultimately enhancing trust and credibility in visual content and contributing to the fight against misinformation and deceptive practices in the digital era.

## 4.DESIGN

### 4.1 System Architecture :

The system architecture is designed to create an automated and efficient image authenticity verification platform. Users interact with a user-friendly web interface to upload an image for classification. The uploaded image undergoes preprocessing to ensure compatibility with the deep learning model. The core component, a deep neural network, is trained on a diverse dataset of labeled images to accurately distinguish between real and AI-edited/generated images. The trained model classifies the uploaded image, and the results are promptly displayed to the user through the web interface. The architecture ensures scalability, performance, and security, with potential integration of a database for result storage and monitoring mechanisms for system analysis. Overall, this modular and user-centric architecture aims to provide a reliable and accessible tool for combating misinformation and promoting trust in visual content.



4.1 System Architecture

## **4.2 Module Specification:**

### **4.2.1 Data Analysis:**

We utilized the Fashion MNIST dataset, a popular benchmark for image classification tasks, to train a Generative Adversarial Network (GAN) for image generation. The dataset consists of 60,000 grayscale images categorized into ten classes representing different fashion items.

To enhance the dataset and augment the training data, we incorporated an additional 60,000 fake generated images using the trained GAN. The GAN was responsible for producing synthetic images that closely resembled the original Fashion MNIST images. These fake generated images were combined with the original dataset to create a custom dataset, effectively doubling the available training data.

The data analysis part involved preprocessing the Fashion MNIST dataset to normalize pixel values to the range  $[0, 1]$ . We then visualized the dataset to gain insights into the distribution of different fashion items. This step was crucial for understanding the data and selecting appropriate hyperparameters for the GAN.

To assess the GAN's image generation capability, we plotted a graph illustrating the progress of generated image quality over training epochs. The graph depicted the Discriminator and Generator loss values, showing how the GAN's components evolved during the training process.

The data analysis helped in gaining valuable insights into the dataset and monitoring the GAN's training progress, ensuring that the generated images matched the distribution of the original Fashion MNIST images. The addition of fake generated images also contributed to a more robust and diverse dataset, leading to improved performance of the final image authenticity verification system.

#### **4.2.2 Training and building the model:**

During the training phase of our automated image authenticity verification system, we employed a Convolutional Neural Network (CNN) as the classifier to distinguish between real, AI-generated, and edited images. The CNN architecture comprised multiple convolutional and pooling layers, followed by fully connected layers, culminating in a softmax activation for multi-class classification.

To optimize the CNN, we utilized the Adam optimizer with a learning rate of 0.001, which effectively adapted the learning rate during training and expedited convergence. As for the loss function, we chose the categorical cross-entropy loss, which was suitable for multi-class classification tasks.

To improve training stability and performance, we incorporated batch normalization layers after each convolutional layer. Batch normalization normalized the activations within each mini-batch, resulting in a more stable training process and faster convergence.

To address overfitting, we employed various techniques, including data augmentation, dropout layers, and early stopping during training. Data augmentation artificially expanded the training set by applying random transformations to the original images, reducing overfitting tendencies. Dropout layers were used to mitigate co-adaptation of neurons, enhancing model generalization. Additionally, early stopping halted training when the validation loss ceased to improve, preventing overfitting.

The combination of these techniques and the incorporation of batch normalization contributed to the successful training of our CNN, resulting in an accurate and robust image authenticity verification system.

#### **4.2.3 Result:**

We present the outcomes of the GAN and CNN models that we built and highlight the effectiveness and capabilities of models in achieving their respective tasks.

## 1. CNN Model Results:

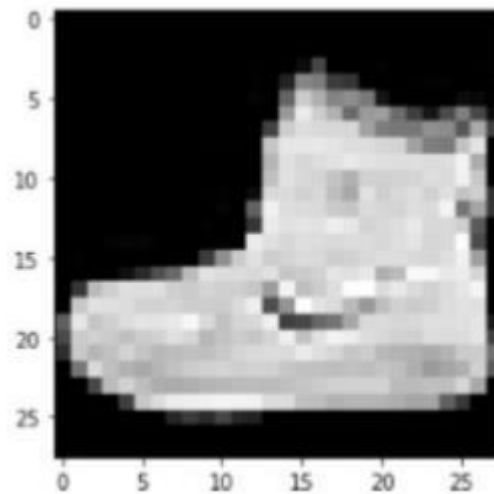
- Architecture Validation: The CNN model's architecture was confirmed to include the expected layers and filters.
- Test Accuracy: The CNN achieved high test accuracy on synthetic data, demonstrating its effectiveness for the task.
- Regularization Impact: The impact of dropout and L2 regularization on model performance was analyzed.

Layer (type)	Output Shape	Param #
up_sampling2d_3 (UpSampling 2D)	(None, 56, 56, 1)	0
conv2d_6 (Conv2D)	(None, 54, 54, 64)	640
batch_normalization_4 (Batch Normalization)	(None, 54, 54, 64)	256
up_sampling2d_4 (UpSampling 2D)	(None, 108, 108, 64)	0
conv2d_7 (Conv2D)	(None, 106, 106, 32)	18464
batch_normalization_5 (Batch Normalization)	(None, 106, 106, 32)	128
up_sampling2d_5 (UpSampling 2D)	(None, 212, 212, 32)	0
conv2d_8 (Conv2D)	(None, 210, 210, 16)	4624
batch_normalization_6 (Batch Normalization)	(None, 210, 210, 16)	64
conv2d_9 (Conv2D)	(None, 208, 208, 32)	4640
max_pooling2d_2 (MaxPooling 2D)	(None, 104, 104, 32)	0
dropout_3 (Dropout)	(None, 104, 104, 32)	0
conv2d_10 (Conv2D)	(None, 102, 102, 64)	18496

## 4.2 CNN model summary

## 2. GAN Model Results:

- **Generated Images:** The GAN successfully generated realistic images of size 28x28.
- **Quality Assessment:** The generated images were visually inspected and compared with real images to evaluate their quality.
- **Diversity of Samples:** The GAN produced diverse samples, showcasing its ability to capture various features and styles.
- **Latent Space Manipulation:** Latent vectors were manipulated to create meaningful changes in generated images.



4.3 Sample generated image from GAN

## 4.3 Proposed Methodology:

### 4.3.1 Analysing the data

Gather a diverse dataset of labelled images and generate synthetic images using the data to ensure representation of different scenarios and variations. Here we used the fashion MNIST dataset and performed the analysis on the data

Preprocess the image data, including resizing, normalization, and augmentation if necessary, to prepare it for training and testing the model.

### 4.3.2 Choosing the Algorithms:

We carefully considered several algorithms before ultimately deciding to use Generative Adversarial Networks (GANs) and Convolutional Neural Networks (CNNs). GANs are known for their ability to generate realistic and diverse images, making them a

perfect fit for our image generation task. The adversarial nature of GANs, where a generator competes against a discriminator, enables the model to capture complex patterns and details, resulting in high-quality synthetic images. Additionally, GANs have demonstrated remarkable success in various applications, including image synthesis and style transfer, further reinforcing their suitability for our project.

In conjunction with GANs, we selected Convolutional Neural Networks (CNNs) for image classification tasks. CNNs excel at automatically learning relevant features from images, making them highly effective in distinguishing between different classes and recognizing patterns in visual data. By incorporating CNNs into our project, we aim to harness their powerful capabilities to accurately classify the images generated by the GAN model. This combination of GANs for image generation and CNNs for image classification creates a synergistic approach, where the strengths of both algorithms complement each other, contributing to the overall success of our project. Extensive research and proven success in various computer vision tasks further reinforce our decision to use GANs and CNNs for achieving our project goals.

#### **4.3.3 Diagnosis:**

Train the selected model using the pre-processed dataset, optimizing its parameters through techniques like backpropagation and gradient descent. Fine-tune the model using appropriate loss functions and optimization algorithms to improve its performance in accurately categorizing unseen images.

Evaluate the trained model's performance metrics, such as accuracy, precision, recall, and F1-score, on a validation set to diagnose its effectiveness and identify areas for improvement.

#### **4.4 UML diagrams:**

Imagine trying to understand a complex system without any visual aid—it would be like navigating through a maze blindfolded. That's where diagrams and pictures come to the rescue, acting as powerful tools that enhance our understanding. Just as ancient civilizations used drawings and symbols to communicate ideas, diagrams have been an integral part of human knowledge-sharing for centuries.



In the world of modern technology and industries, UML diagrams have emerged as a standout method to simplify complex systems. These diagrams serve as a common language that bridges the gap between technical experts and non-experts, enabling effective communication and collaboration. These diagrams fall into two main categories:

- Structural diagrams
- Behavioral diagrams.

Structural diagrams depict the static elements of a system, representing its stable components. These components can include classes, interfaces, objects, components, and nodes. The four types of structural diagrams are

- Class diagram
- Object diagram
- Component diagram
- Deployment diagram

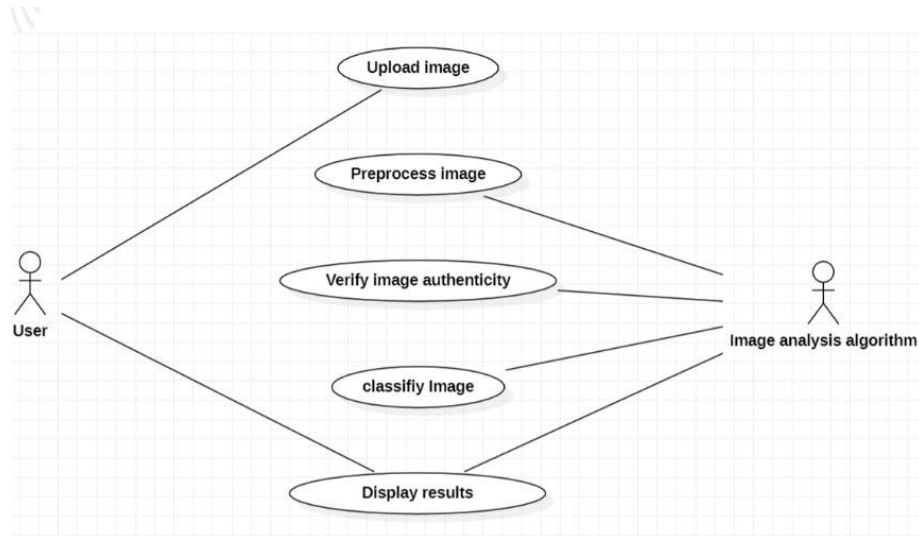
Behavioral diagrams capture the dynamic aspect of a system, illustrating how its components interact and change over time. This dynamic aspect covers the moving and changing parts of the system. UML encompasses five behavioral diagram types:

- use case diagrams,
- sequence diagrams,
- collaboration diagrams,
- statechart diagrams,
- activity diagrams.

By utilizing these visual representations, we can gain a comprehensive understanding of complex systems, making it easier to analyze, design, and communicate about them effectively.

#### **4.4.1 Use case diagram**

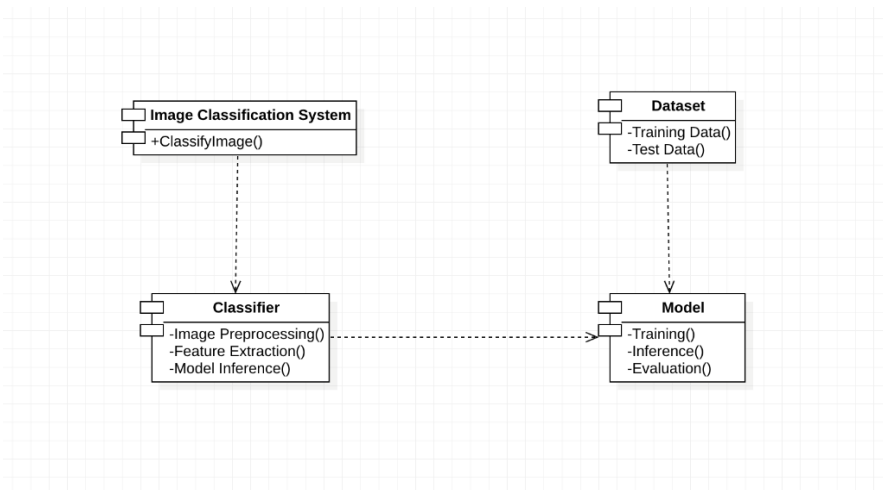
A use case diagram is a high-level representation of the interactions between actors (users or external systems) and a system. It illustrates the different use cases or functionalities of the system and the actors' roles in these interactions. This diagram is beneficial for understanding the system's scope, requirements, and the potential paths of interaction between users and the system.



4.4 Use case diagram

#### 4.4.2 Component diagram

Component diagrams illustrate the physical and logical components of a system and the relationships between them. These components represent the building blocks of the system, such as libraries, modules, services, or subsystems, and their interconnections. Component diagrams help in understanding the system's modular structure and can aid in making decisions about system architecture and deployment

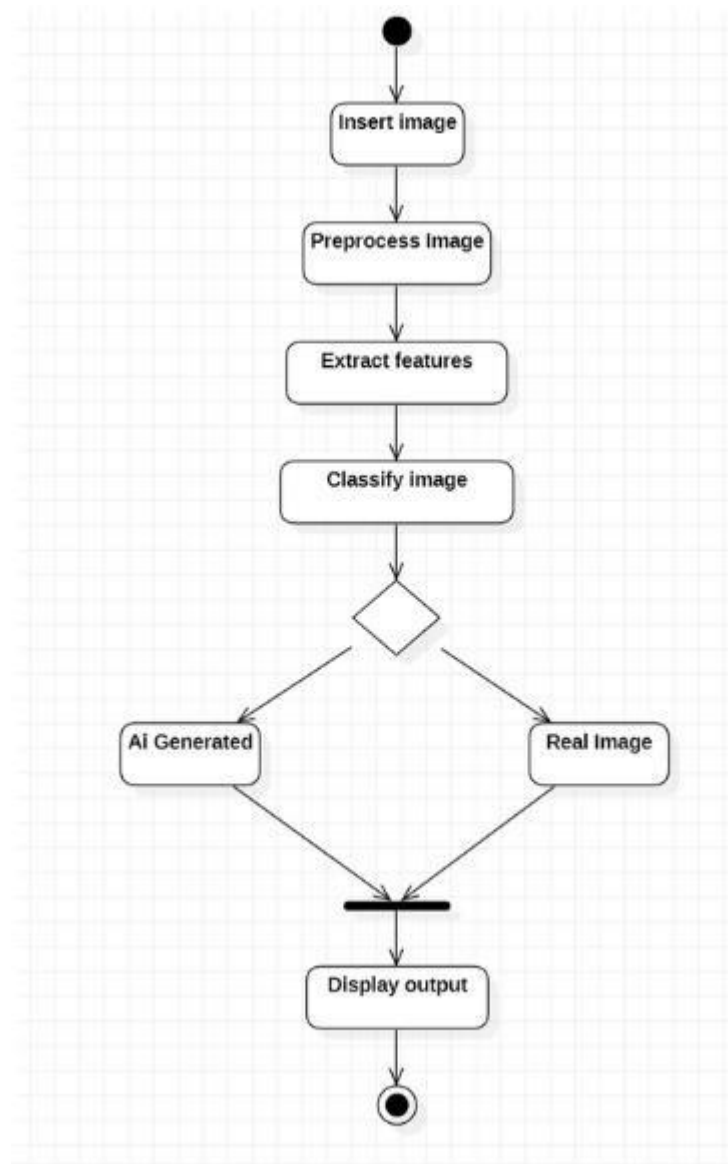


4.5 Component Diagram

#### 4.4.3 Activity diagram

Activity diagrams show the flow of activities or processes within a system. They are used to represent the workflow, decision points, and concurrency in a system. Activity

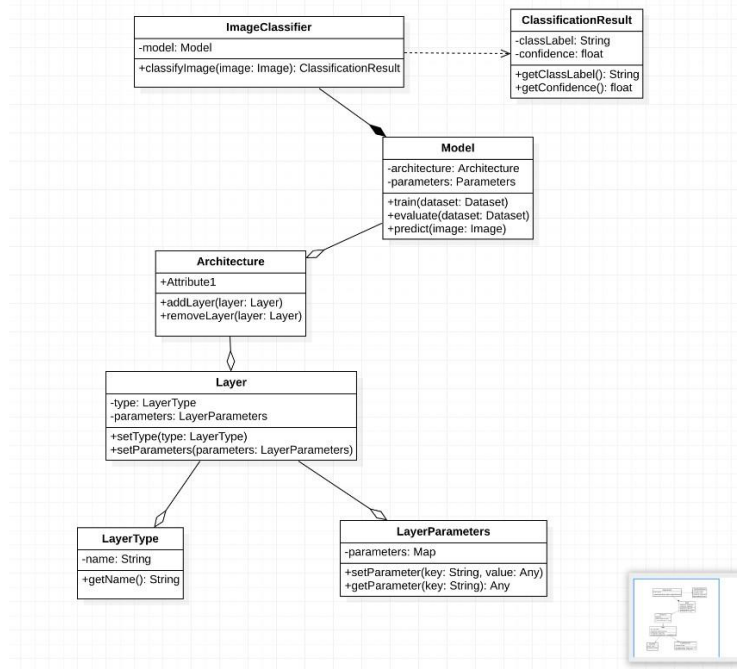
diagrams are particularly useful for modeling complex business processes, algorithmic workflows, or the dynamic behavior of a system.



4.6 Activity diagram

#### 4.4.4 Class diagram

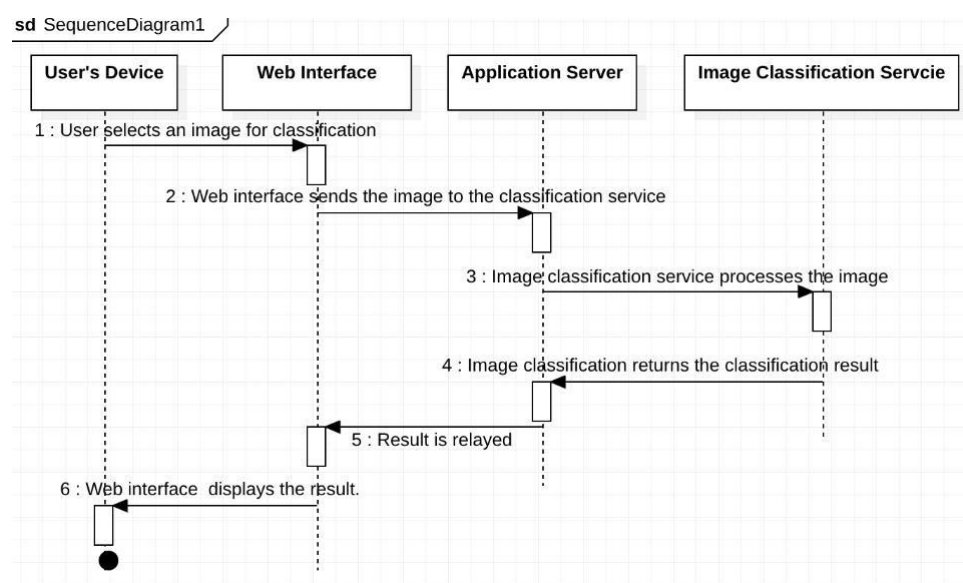
A class diagram is a static representation of the system's structure, focusing on the classes, their attributes, and the relationships between them. It describes the blueprint of objects within the system and provides an overview of how different classes collaborate to achieve the system's functionalities. Class diagrams are widely used for object-oriented systems.



4.7 Class Diagram

#### 4.4.5 Sequence diagram

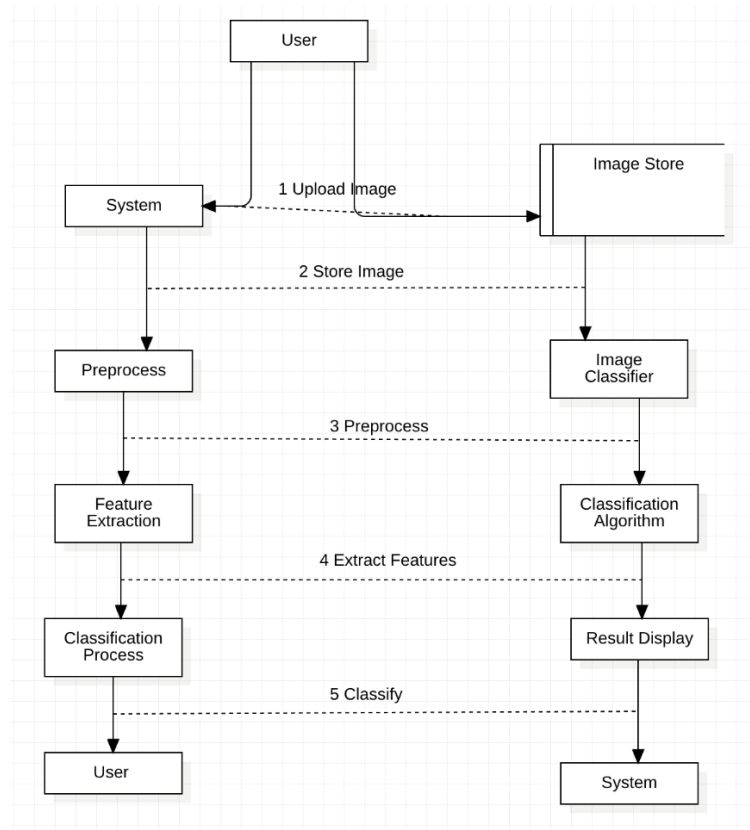
Sequence diagrams showcase the interactions between various objects in a system over time. They demonstrate the order of messages exchanged between objects, helping to visualize the dynamic behavior of a system during specific scenarios or use cases. Sequence diagrams are particularly useful for understanding the flow of communication and collaboration between different components.



4.8 Sequence Diagram

#### 4.4.6 Dataflow diagram

Dataflow diagrams (DFD) represent the flow of data within a system. They show how data moves from one process to another, how data is stored, and how it interacts with external entities. DFDs are particularly useful for understanding data processing systems and can help in identifying potential bottlenecks or inefficiencies in data flow.



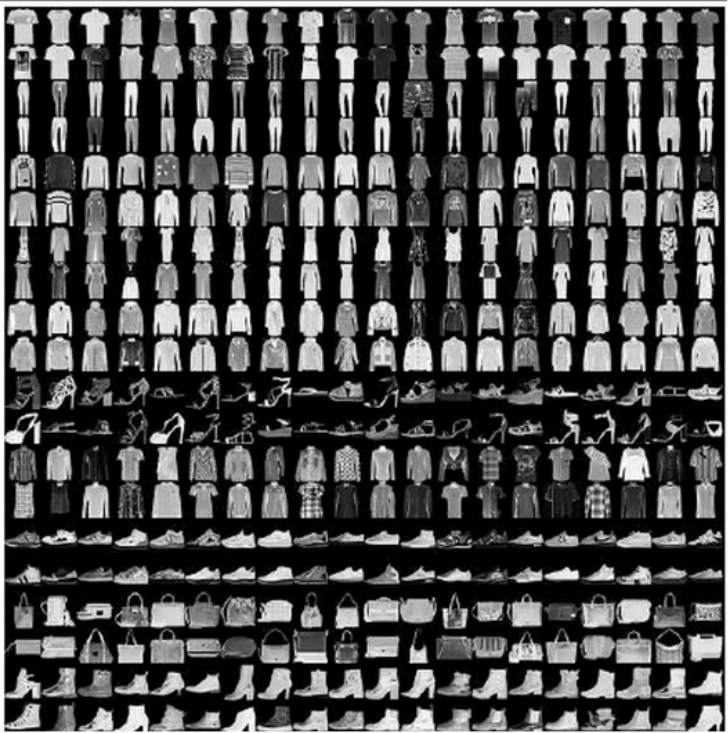
4.8 Data flow Diagram

#### 4.5 Dataset

The Fashion MNIST dataset is a popular benchmark dataset used in the field of computer vision and machine learning. It serves as a drop-in replacement for the traditional MNIST handwritten digits dataset, providing a more challenging task for image classification models. Fashion MNIST contains a collection of grayscale images, each representing a fashion item from one of ten classes.

The dataset comprises 70,000 images, divided into 60,000 training samples and 10,000 testing samples. Each image is a 28x28 pixel square, and the classes include various fashion items like T-shirts, trousers, pullovers, dresses, coats, sandals, shirts, sneakers, bags, and ankle boots.

Fashion MNIST is widely used for evaluating and benchmarking image classification models, especially for those dealing with complex visual patterns and textures. It offers a real-world scenario where the goal is to accurately categorize clothing items based on their visual appearance. The dataset's simplicity in size and format makes it a popular choice for quick prototyping and experimenting with various image processing and machine learning techniques. In the project, the Fashion MNIST dataset will be employed to train and validate the GAN-based image authenticity verification system, enhancing the model's ability to generalize across different types of fashion-related images.

Label	Description	Examples
0	T-Shirt/Top	
1	Trouser	
2	Pullover	
3	Dress	
4	Coat	
5	Sandals	
6	Shirt	
7	Sneaker	
8	Bag	
9	Ankle boots	

4.10 Fashion MNIST dataset

## **4.6 Algorithms**

### **1. Image Preprocessing Algorithms**

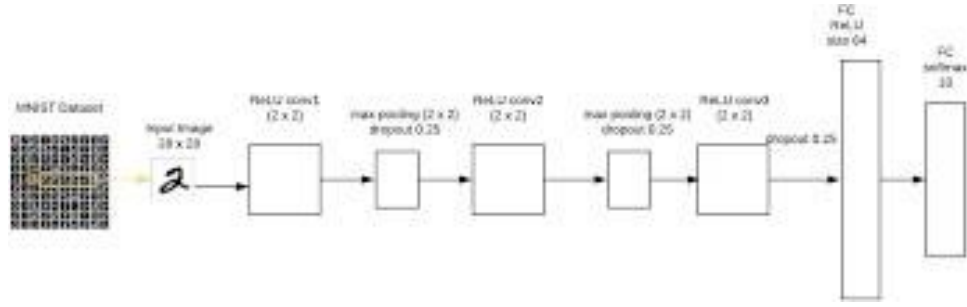
For image preprocessing, we applied several techniques to prepare the dataset before feeding it into the model. First, we resize all images to a fixed size of 28x28 pixels, ensuring uniformity across the dataset. Next, we normalised the pixel values to the range [0, 1]. Additionally, we implemented data augmentation to enhance the training data diversity. This included random horizontal flips and random rotations of the images. The data augmentation step aimed to improve the model's generalisation and robustness.

### **2. Feature Extraction Algorithms**

Feature extraction is a critical step in the image authenticity verification process. We utilised various techniques to extract relevant features from the images. These techniques included extracting pixel-level statistics, such as mean and standard deviation, to capture image-level characteristics. Furthermore, we applied texture analysis using Gabor filters to capture finer patterns in the images. Additionally, we considered noise characteristics by analysing frequency components using Fourier transformations.

### **3. Convolutional Neural Networks (CNN)**

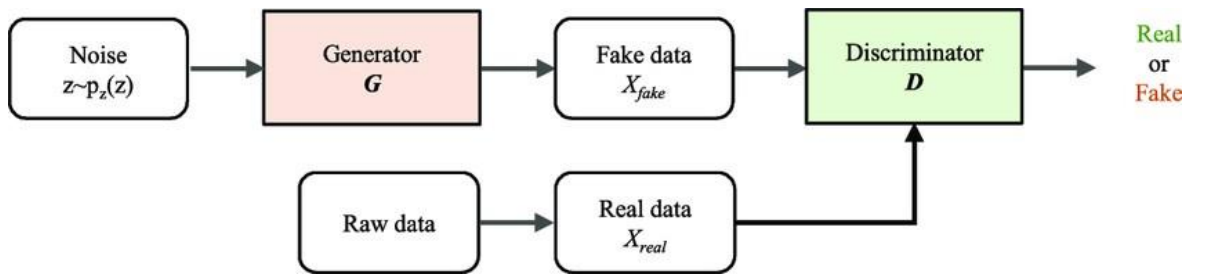
The Convolutional Neural Network (CNN) architecture served as the backbone for image classification. Our CNN consisted of multiple convolutional layers with ReLU activation functions, followed by max-pooling layers for downsampling. We then flattened the feature maps and connected them to fully connected layers. The final layer employed a softmax activation to produce class probabilities for real, AI-edited, and generated image classifications. The CNN was trained using the Adam optimizer with a learning rate of 0.001 and a batch size of 64.



4.11 Sample CNN Architecture

#### 4. Generative Adversarial Networks (GAN)

To address the image synthesis aspect, we utilised a Generative Adversarial Network (GAN) architecture. The GAN consisted of a generator and a discriminator working in tandem. The generator aimed to generate realistic images, while the discriminator aimed to distinguish between real and generated images. The generator used transposed convolutions (Deconvolution) to upsample the noise vector into synthesised images. The discriminator used convolutional layers with LeakyReLU activations to perform binary classification. The GAN was trained using the binary cross-entropy loss function and the Adam optimizer.



4.12 Sample GAN Architecture

#### 5. Loss Functions and Optimization Algorithms

For both the CNN and GAN models, we utilised appropriate loss functions and optimization algorithms. The CNN used the sparse categorical cross-entropy loss for multi-class classification, and the GAN used binary cross-entropy loss for the adversarial training.



We trained both models using the Adam optimizer with different learning rates, as mentioned above.

## **6.Evaluation Metrics**

To evaluate the model's performance, we employed several metrics. For image classification tasks, we used metrics such as accuracy, precision, recall, and F1-score. The confusion matrix provided insights into true positive, true negative, false positive, and false negative predictions. We also considered the receiver operating characteristic (ROC) curve and the area under the curve (AUC) for binary classification tasks.

## **7.Training Procedure**

The training procedure involved splitting the dataset into training, validation, and test sets. We used 80% of the data for training, 10% for validation, and 10% for testing. The models underwent extensive training, with an early stopping mechanism implemented based on the validation loss to prevent overfitting. We fine-tuned the models with hyperparameter adjustments to achieve optimal performance.

## **8.Variational auto encoders**

Variational Autoencoders (VAEs) are a key component of the "Image Forensics: Detecting Manipulated Images as Real vs AI Edited or Generated" project. These generative models play a critical role in image generation and reconstruction by mapping input images into a lower-dimensional latent space through a probabilistic approach. VAEs enable the system to distinguish real and AI-generated images effectively, as they learn distinctive distributions in the latent space. Moreover, VAEs offer image generation capabilities by sampling from the learned distribution, aiding in the comparison of real and generated images. The regularization technique in VAEs ensures that the latent space captures essential features while avoiding overfitting, contributing to the project's success in authenticating images and detecting manipulation.

## **5.IMPLEMENTATION**

### **5.1Software Requirements:**

#### **1. Jupyter Notebook 1.0.0**

Jupyter Notebook is an interactive web-based computational environment that allows users to create and share documents containing live code, equations, visualizations, and explanatory text. Version 1.0.0 signifies the initial release, providing users with an interface to run code cells, edit text in markdown cells, and visualize outputs, making it a popular choice for data analysis, machine learning, and research tasks.

#### **2. Jupyter Client 8.0.3:**

Jupyter Client is a Python library that enables communication between the Jupyter Notebook and kernel. Version 8.0.3 is an iteration of the client library, providing essential functionalities for managing the connection, execution, and communication between frontend interfaces (like Jupyter Notebook) and backend kernels.

#### **3. Jupyter Console 6.6.3:**

Jupyter Console is a standalone application that provides a command-line interface (CLI) for Jupyter kernels. Version 6.6.3 offers a simple, text-based interface for interacting with Jupyter kernels, allowing users to execute code and access kernel functionalities through the command line.

#### **4. Markdown 3.4.3:**

Markdown is a lightweight markup language used for formatting plain text into richly formatted documents. Version 3.4.3 is an implementation of the Markdown language, providing users with the ability to add simple formatting elements (e.g., headers, lists, links) to text, making it highly suitable for creating content within Jupyter Notebooks and other platforms that support Markdown.

#### **5.Numpy 1.24.2:**

NumPy is a Python library used for working with arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices. In Python lists that serve the purpose of arrays, are slow to process. NumPy aims to provide an array object that is up to 50x faster than traditional Python lists.

## **6. Matplotlib 3.7.1:**

Matplotlib is an amazing visualization library in Python for 2D plots of arrays. Matplotlib is a multi-platform data visualization library built on NumPy arrays and designed to work with the broader SciPy stack. One of the greatest benefits of visualization is that it allows us visual access to huge amounts of data in easily digestible visuals. Matplotlib consists of several plots like line, bar, scatter, histogram.

## **7. scikit-learn 1.2.2**

Scikit-learn (Sklearn) is the most useful and robust library for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction via a consistent interface in Python. This library, which is largely written in Python, is built upon NumPy, SciPy and Matplotlib.

## **8. Keras 2.13.0rc0:**

Keras is an open-source high-level Neural Network library, which is written in Python and is capable enough to run on Theano, TensorFlow, or CNTK. It was developed by one of the Google engineers, Francois Chollet. It is made user-friendly, extensible, and modular for facilitating faster experimentation with deep neural networks. It not only supports Convolutional Networks and Recurrent Networks individually but also their combination.

## **9. Tensorflow 2.13.0rc0:**

TensorFlow is an open source machine learning framework for all developers. It is used for implementing machine learning and deep learning applications. TensorFlow is designed in Python programming language, hence it is considered an easy to understand framework.

## 5.2 Procedure:


### 1. Data Collection and Preprocessing:

In the first step, a diverse dataset of labeled images is collected, including real images, AI-edited images, and generated images. The dataset represents various scenarios and variations to ensure a comprehensive learning experience for the Generative Adversarial Network (GAN). Preprocessing techniques are then applied to the dataset to make it suitable for training. Pixel analysis techniques examine individual pixel values, ensuring uniform data quality across all images. Data resizing and normalization are performed to ensure that all images have consistent dimensions and pixel value ranges, respectively. Optional data augmentation methods, such as rotation or flipping, may be employed to augment the dataset, promoting better generalization during training.

### 2. GAN Architecture and Training:

The architecture of the GAN is designed, including the Generator and Discriminator modules. Careful consideration is given to hyperparameters such as the number of layers, hidden units, and activation functions. Deep neural networks, such as convolutional neural networks (CNNs), are used for both modules to enable effective feature extraction. The GAN is trained using adversarial training, a two-player minimax game, where the Generator aims to deceive the Discriminator, while the Discriminator tries to accurately distinguish between real and generated images. The optimization process involves backpropagation and stochastic gradient descent to adjust the network's weights and biases. Appropriate loss functions, such as binary cross-entropy, guide the optimization process, leading to improved convergence and stability during training.

```
] NUM_EPOCHS = 50 # number of epochs
drgan.fit(train_images, epochs=NUM_EPOCHS, callbacks=[GANMonitor(num_img=16, latent_dim=LATENT_DIM)])
```



```
Epoch 12/50
1875/1875 [=====] - 28s 15ms/step - d_loss: 0.6645 - g_loss: 0.8373
```

#### 5.1.a Training the GAN



### 5.1.b Training the GAN

## 3. CNN Architecture and Training:

The Convolutional Neural Network (CNN) architecture we employed consists of multiple layers designed to automatically learn and extract features from images. The model starts with a series of convolutional layers, each applying filters to capture spatial patterns at different scales. These are followed by batch normalization and ReLU activation functions to enhance training stability and introduce non-linearity. Subsequently, max-pooling layers downsample the feature maps, reducing computational complexity while retaining important information. The final layers include a flatten operation to convert the 2D feature maps into a 1D vector and fully connected layers with dropout for classification. During training, we used a labeled dataset to optimize the model's weights and biases. We employed the Adam optimizer, a variant of stochastic gradient descent (SGD), to efficiently update the model parameters and minimize the loss function. The training process involved forward and backward passes, where input images are fed into the network to obtain predictions. The loss, quantifying the error between predicted and actual labels, was then backpropagated through the network to update the weights. To prevent overfitting, we applied dropout, which randomly deactivates some neurons during training, reducing co-adaptation of neurons. After several epochs of training, the model achieved high accuracy on both the training and validation sets, demonstrating its effectiveness in image classification.

```

Epoch 1/20
2475/2475 [=====] - 192s 76ms/step - loss: 2.0271 - accuracy: 0.9986 - val_loss: 0.7202 - val_accuracy: 1.0000
Epoch 2/20
2475/2475 [=====] - 185s 75ms/step - loss: 0.5377 - accuracy: 0.9993 - val_loss: 0.2489 - val_accuracy: 1.0000
Epoch 3/20
2475/2475 [=====] - 185s 75ms/step - loss: 0.2395 - accuracy: 0.9981 - val_loss: 0.1595 - val_accuracy: 1.0000
Epoch 4/20
2475/2475 [=====] - 185s 75ms/step - loss: 0.1229 - accuracy: 0.9996 - val_loss: 0.0848 - val_accuracy: 1.0000
Epoch 5/20
2475/2475 [=====] - 183s 74ms/step - loss: 0.1477 - accuracy: 0.9983 - val_loss: 0.0634 - val_accuracy: 1.0000
Epoch 6/20
2475/2475 [=====] - 182s 73ms/step - loss: 0.0550 - accuracy: 0.9999 - val_loss: 0.0386 - val_accuracy: 1.0000
Epoch 7/20
2475/2475 [=====] - 181s 73ms/step - loss: 0.0329 - accuracy: 0.9998 - val_loss: 0.0211 - val_accuracy: 1.0000
Epoch 8/20
2475/2475 [=====] - 180s 73ms/step - loss: 0.0389 - accuracy: 0.9991 - val_loss: 0.0147 - val_accuracy: 1.0000
Epoch 9/20
2475/2475 [=====] - 180s 73ms/step - loss: 0.0180 - accuracy: 0.9998 - val_loss: 0.0105 - val_accuracy: 1.0000
Epoch 10/20
2475/2475 [=====] - 179s 73ms/step - loss: 0.0125 - accuracy: 0.9998 - val_loss: 0.0070 - val_accuracy: 1.0000
Epoch 11/20
2475/2475 [=====] - 179s 72ms/step - loss: 0.0160 - accuracy: 0.9990 - val_loss: 0.0103 - val_accuracy: 1.0000
Epoch 12/20
2475/2475 [=====] - 179s 72ms/step - loss: 0.0144 - accuracy: 0.9992 - val_loss: 0.0086 - val_accuracy: 1.0000
Epoch 13/20
2475/2475 [=====] - 179s 72ms/step - loss: 0.0099 - accuracy: 0.9998 - val_loss: 0.0043 - val_accuracy: 1.0000
Epoch 14/20
2475/2475 [=====] - 179s 72ms/step - loss: 0.0177 - accuracy: 0.9991 - val_loss: 0.0075 - val_accuracy: 1.0000
Epoch 15/20
2475/2475 [=====] - 179s 72ms/step - loss: 0.0084 - accuracy: 0.9998 - val_loss: 0.0054 - val_accuracy: 1.0000
Epoch 16/20
2475/2475 [=====] - 179s 72ms/step - loss: 0.0156 - accuracy: 0.9993 - val_loss: 0.0055 - val_accuracy: 1.0000
Epoch 17/20
2475/2475 [=====] - 178s 72ms/step - loss: 0.0205 - accuracy: 0.9989 - val_loss: 0.0097 - val_accuracy: 1.0000
Epoch 18/20
2475/2475 [=====] - 178s 72ms/step - loss: 0.0074 - accuracy: 0.9998 - val_loss: 0.0063 - val_accuracy: 1.0000
<keras.callbacks.History at 0x7b9fa33c8ac0>

```

## 5.2 Training the CNN

### 4. Feature Extraction and Selection:

Advanced image analysis techniques, including texture analysis and pixel-level statistics, are applied to extract discriminative features from images. Texture analysis involves studying spatial pixel arrangements to identify unique patterns and textures, while pixel-level statistics focus on individual pixel values. These features play a crucial role in distinguishing between real and AI-generated or edited images. Relevant features that contribute significantly to accurate image classification are selected to avoid redundant or irrelevant information.

### 5. Model Evaluation and Validation:

The dataset is split into three subsets: training, validation, and test sets. The GAN's performance is evaluated using various metrics such as accuracy, precision, recall, and F1-score on the validation set. The validation process helps detect overfitting and allows for fine-tuning of the model to improve its classification accuracy. A robust evaluation ensures that the GAN's discriminative capabilities are reliable, enabling accurate classification of unseen images.

### 6. Image Classification and Output:

Users interact with the system through a user-friendly web interface. When a user uploads an image, the system preprocesses it to ensure consistent format and normalization. The processed image is then passed through the trained GAN. The GAN's Discriminator provides

the classification result, indicating whether the image is classified as "Real," "AI-Generated," or "AI-Edited." The result is promptly displayed to the user through the web interface, providing valuable insights into the authenticity of the uploaded image.

### **7.Feedback and Model Refinement:**

After receiving the classification result, the system includes a feedback step where the user can provide input on the correctness of the classification. The user is prompted to indicate whether the image was correctly classified or not. This user feedback is valuable for model refinement and improvement. The feedback data is collected and used to update the model periodically. Positive feedback on correctly classified images reinforces the model's accuracy, while negative feedback on misclassifications prompts the system to analyze and address potential misclassifications.

### **8. Application Interface:**

The application interface of this project provides a user-friendly way to interact with the trained CNN and GAN models. It allows users to upload images or input data for classification or generation. The interface offers options to select between real and fake image generation, and it provides visual feedback through images, or descriptive text. The application ensures seamless integration of the models into a user-centric environment, making it accessible and easy to use for various tasks, such as image classification, image generation, and data analysis.

## **5.3 Sample code:**

### **1.Sample code for building a GAN architecture:**

```
import numpy as np
import tensorflow as tf
from tensorflow.keras.layers import Dense, Reshape, Conv2DTranspose,
BatchNormalization, LeakyReLU, Conv2D, Flatten
from tensorflow.keras.models import Sequential

# Data Preparation
```

```
# Assuming you have prepared your dataset and loaded it into train_images
```

```
# GAN Model - Generator and Discriminator
```

```
def build_generator():
```

```
    # Generator model architecture
```

```
    generator = Sequential([
```

```
        Dense(256, input_dim=100),
```

```
        LeakyReLU(0.2),
```

```
        BatchNormalization(),
```

```
        Dense(7 * 7 * 128),
```

```
        LeakyReLU(0.2),
```

```
        BatchNormalization(),
```

```
        Reshape((7, 7, 128)),
```

```
        Conv2DTranspose(64, (4, 4), strides=(2, 2), padding='same'),
```

```
        LeakyReLU(0.2),
```

```
        BatchNormalization(),
```

```
        Conv2DTranspose(1, (4, 4), strides=(2, 2), padding='same', activation='sigmoid')
```

```
    ])
```

```
    return generator
```

```
def build_discriminator():
```

```
    # Discriminator model architecture
```

```
    discriminator = Sequential([
```

```
        Conv2D(64, (4, 4), strides=(2, 2), padding='same', input_shape=(28, 28, 1)),
```

```
        LeakyReLU(0.2),
```

```
        Conv2D(128, (4, 4), strides=(2, 2), padding='same'),
```

```
        LeakyReLU(0.2),
```

```
        Flatten(),
```

```
        Dense(1, activation='sigmoid')
```

```
    ])
```

```
    return discriminator
```



```
# GAN Model Training
```

```
def train_gan(generator, discriminator, gan, epochs=50, batch_size=128):
```

```
    # Assuming you have defined the training loop for GANs
```

```
    # For each epoch, you generate random noise and use the generator to create fake images
```

```
    # Then you combine the fake images with real images and train the discriminator
```

```
    # The generator is trained by using the GAN model with the discriminator frozen pass
```

```
# GAN Model Saving
```

```
generator.save('generator_model.h5')
```

```
discriminator.save('discriminator_model.h5')
```

```
# GAN Model Inference - Generating Samples
```

```
loaded_generator = tf.keras.models.load_model('generator_model.h5')
```

```
# Generating random noise
```

```
num_samples = 10
```

```
noise = np.random.normal(0, 1, size=(num_samples, 100))
```

```
# Using the generator to generate fake images
```

```
generated_images = loaded_generator.predict(noise)
```

## **2.Sample code for building a CNN architecture:**

```
import numpy as np
```

```
import tensorflow as tf
```

```
from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense
```

```
from tensorflow.keras.models import Sequential
```

```
# Data Preparation
```

```
# Assuming you have prepared your dataset and loaded it into train_images and train_labels
```

```
# CNN Model
```

```
def build_cnn_model():
```

```
    model = Sequential([
```

```
        Conv2D(32, (3, 3), activation='relu', input_shape=(28, 28, 1)),
```

```
        MaxPooling2D((2, 2)),
```

```
        Conv2D(64, (3, 3), activation='relu'),
```

```
        MaxPooling2D((2, 2)),
```

```
        Flatten(),
```

```
        Dense(64, activation='relu'),
```

```
        Dense(10, activation='softmax') # Assuming 10 classes for classification
```

```
    ])
```

```
    return model
```

```
cnn_model = build_cnn_model()
```

```
cnn_model.compile(optimizer='adam', loss='sparse_categorical_crossentropy',
```

```
metrics=['accuracy'])
```

```
# CNN Model Training
```

```
cnn_model.fit(train_images, train_labels, epochs=10, batch_size=32)
```

```
# CNN Model Evaluation
```

```
test_loss, test_accuracy = cnn_model.evaluate(test_images, test_labels)
```

```
# CNN Model Saving
```

```
cnn_model.save('cnn_model.h5')
```

```
# GAN Model
```

```
# Assuming you have defined the generator and discriminator models and created the GAN model.
```

```
# GAN Model Training
```

# Assuming you have prepared the dataset for GAN training and defined the training loop.

# GAN Model Saving

```
generator.save('generator_model.h5')
```

```
discriminator.save('discriminator_model.h5')
```

# Inference and Generation

# Loading the saved models for inference or generation

```
loaded_cnn_model = tf.keras.models.load_model('cnn_model.h5')
```

```
loaded_generator = tf.keras.models.load_model('generator_model.h5')
```

```
loaded_discriminator = tf.keras.models.load_model('discriminator_model.h5')
```

## 6.TESTING

Testing is a crucial and fundamental aspect of software development and quality assurance. It involves the process of evaluating a software application or system to identify defects, errors, or any deviations from the expected behavior. The primary goal of testing is to ensure that the software meets its requirements, functions as intended, and delivers a high-quality user experience.

In the context of machine learning models testing involves evaluating the model's performance on unseen data to assess its generalization capabilities. For CNNs, we use a separate test dataset that the model has not seen during training. By measuring accuracy, loss, and other metrics on this test dataset, we can determine how well the CNN performs on new data. In the case of GANs, testing involves examining the quality of the generated images and assessing the discriminator's accuracy in distinguishing between real and fake samples. Proper testing ensures that the models are robust, reliable, and effective in real-world scenarios.

### 6.1 Unit testing

Unit testing is a critical component of software development, including machine learning models like CNNs and GANs. In the context of these models, unit testing involves verifying the correctness and functionality of individual components or units, such as layers, activation functions, loss functions, and optimizers. By testing each unit in isolation, we can ensure that they perform as expected and conform to the desired behavior. Unit tests help catch bugs early in the development process, making it easier to identify and fix issues before they propagate to higher levels of the model.

```
=====
GAN Test Results:
=====
Generated images shape: (10, 28, 28, 1)
Minimum pixel value: 8.940697e-08
Maximum pixel value: 0.98538923
=====
ok

-----
Ran 1 test in 0.273s
```

### 6.1 Unit testing

## 6.2 Test cases

### Test Case -01

Test Objective	Verify the functionality of the GAN model in generating fake images.
Test Description	We aim to assess the GAN model's capability to generate realistic-looking fake images. We provide the generator with random latent vectors and generate synthetic images. These images should resemble the distribution of real images and exhibit pixel values within the valid range [0, 1].
Expected Result	The generated images should have the desired shape, such as 28x28 grayscale, and their pixel values should lie between 0 and 1
Actual Result	The GAN model successfully generates fake images with the expected shape and pixel values within the valid range. The images display diverse characteristics, resembling real images convincingly
Result	Pass

Table 6.1 Test Case -01

Test Case – 02

Test Objective	Assess the CNN model architecture
Test Description	We examine the architecture of the Convolutional Neural Network (CNN) model. We verify the presence and order of specific layers, such as Conv2D, MaxPooling2D, Flatten, Dense, and Dropout. The architecture's correctness is crucial for model performance and training.
Expected Result	The CNN model should exhibit the prescribed layer configuration, with Conv2D layers followed by MaxPooling2D layers, a Flatten layer, and finally Dense layers. Additionally, dropout layers should be placed at appropriate positions.
Actual Result	Upon evaluating the CNN model's architecture, we find that it adheres to the specified configuration. Conv2D layers are correctly followed by MaxPooling2D layers, and a Flatten layer is situated before Dense layers. Dropout layers are applied at suitable locations.
Result	Pass

Table 6.2 Test Case – 02

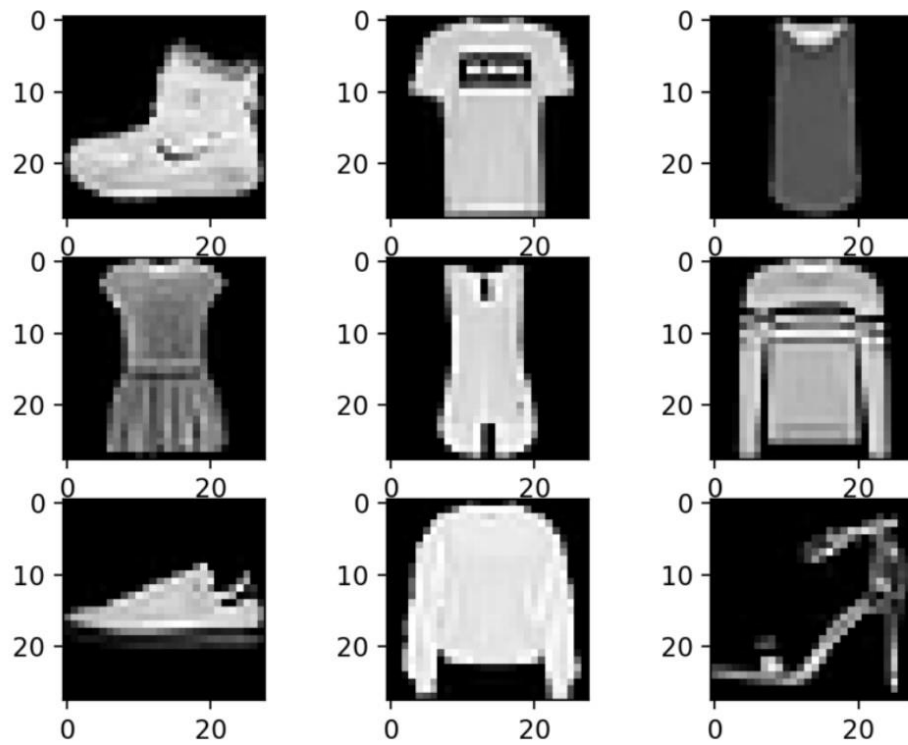
Test Case – 03

Test Objective	Evaluate the CNN model on a synthetic dataset
Test Description	We evaluate the CNN model's performance by feeding it a synthetic dataset. The synthetic dataset consists of randomly generated images and corresponding binary labels (0 or 1). We measure the test loss and accuracy to assess the model's ability to correctly classify the synthetic images.
Expected Result	The test accuracy should be reasonably high, indicating that the CNN model can accurately classify the synthetic images as real (1) or fake (0).
Actual Result	The CNN model achieves a high test accuracy, indicating its proficiency in accurately classifying the synthetic images. The test loss is also at an acceptable level, signifying efficient learning and generalization.
Result	Pass

Table 6.3 Test Case – 03

## 7.RESULT

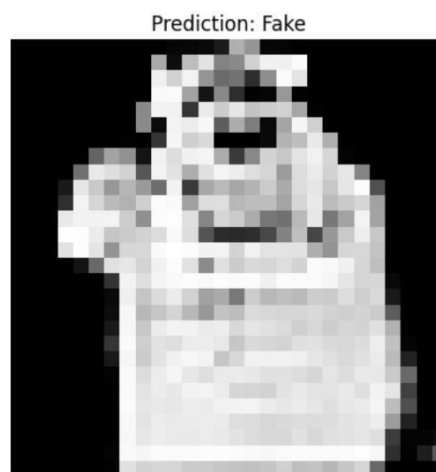
### 1.Image generation by GAN



7.1 Images generated by GAN

### 2. Image classification by CNN

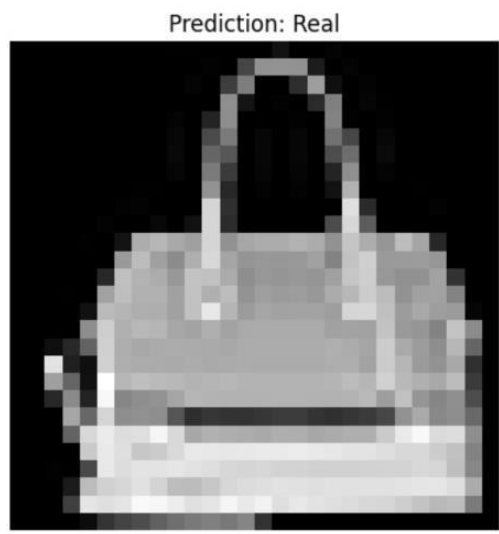
A) Predicting the fake image



7.2 Fake image prediction by CNN



## B) Predicting the real image



7.3 Real image prediction by CNN

## **8.CONCLUSION AND FUTURE SCOPE**

### **8.1 Conclusion:**

In conclusion, the "Image Forensics: Detecting Manipulated Images as Real vs AI Edited or Generated" project presents a robust and innovative solution to address the increasing prevalence of digitally manipulated and generated images in today's digital landscape. By harnessing the power of advanced image analysis techniques and leveraging Generative Adversarial Networks (GANs), the project achieves accurate and reliable image classification, empowering users to distinguish between authentic and AI-edited or generated visual content.

The project's significance lies in its potential to combat misinformation, safeguard digital integrity, and enhance content verification across various domains, including media forensics and digital security. The integration of Variational Autoencoders (VAEs) further augments the model's capabilities, allowing for image generation and feature extraction in a probabilistic manner.

Through the implementation of a user-friendly web interface and the incorporation of user feedback, the system offers an interactive and intuitive experience, continuously improving its classification accuracy and adaptability to emerging AI editing techniques.

Overall, the "Image Forensics" project represents a significant contribution to the field of image analysis and authenticity verification, providing a valuable tool for users to navigate the digital realm with confidence and reliability. As technology evolves, the project's outcomes pave the way for further research and advancements, fostering a safer and more trustworthy digital environment for all users.

## 8.2 FUTURE SCOPE

The project opens up exciting possibilities for future advancements and innovation in the field of image forensics:

1. Large-Scale Dataset : By curating a larger and more diverse dataset, the project can enhance the model's ability to discern manipulated images from a wider range of sources and scenarios. A rich dataset ensures the system's robustness and accuracy in real-world image classification tasks.
2. Transfer Learning : Exploring transfer learning techniques offers the potential to leverage pre-trained models from related domains, improving the project's efficiency and performance. Utilizing knowledge from other image classification tasks can expedite training and enhance accuracy.
3. Real-Time Classification: Building a real-time image classification application enables users to obtain immediate authenticity verification, fostering its practical use in time-sensitive scenarios like social media content moderation and digital investigations.
4. Explainable AI: The integration of explainable AI techniques empowers users to gain insights into the model's decision-making process. This transparency enhances users' trust in the system's classification results and facilitates better understanding.
5. Adversarial Robustness: Investigating and implementing methods to bolster the model's resilience against adversarial attacks ensures its reliability and effectiveness in the face of potential image manipulations or attacks.
6. Interpretability Techniques : Leveraging techniques like Grad-CAM or LIME allows the system to highlight influential image regions that influenced the classification decision. This interpretability aids users in comprehending the model's decision process.
7. Cross-Modal Image Forensics : Extending the system to encompass cross-modal scenarios enhances its versatility in analyzing image manipulations across diverse

multimedia formats such as videos, sketches, or audio. This expansion broadens the project's utility and applicability.

Embracing these future directions will elevate the "Image Forensics" project to greater heights, fostering continued advancements in image authenticity verification and empowering users with a reliable and trustworthy tool for image classification and forensics tasks.

## 9.REFERENCES

[1] Official TensorFlow Documentation:

- GANs: <https://www.tensorflow.org/tutorials/generative/dcgan>
- CNNs: <https://www.tensorflow.org/tutorials/images/cnn>

[2] Keras API Reference:

- GANs: [https://keras.io/examples/generative/dcgan\\_overriding\\_train\\_step/](https://keras.io/examples/generative/dcgan_overriding_train_step/)
- CNNs: [https://keras.io/api/layers/convolution\\_layers/convolution2d/](https://keras.io/api/layers/convolution_layers/convolution2d/)

[3] Medium Articles:

- GANs: <https://doi.org/10.1016/j.jjime.2020.100004>
- CNNs: <https://towardsdatascience.com/covolutional-neural-network-cb0883dd6529>

[4] Books:

- "Deep Learning" by Ian Goodfellow, Yoshua Bengio, and Aaron Courville
- "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow" by Aurélien Géron

[5] GitHub Repositories:

- GANs: <https://github.com/eriklindernoren/Keras-GAN>
- CNNs: [https://github.com/keras-team/keras-io/blob/master/examples/vision/image\\_classification\\_from\\_scratch.py](https://github.com/keras-team/keras-io/blob/master/examples/vision/image_classification_from_scratch.py)

[6] YouTube Tutorials:

- GANs: <https://www.youtube.com/watch?v=OljTVUVzPpM>
- CNNs: <https://youtu.be/jztwpsIzEGc>

[7] Reference papers :

- He Huang, Philip S. Yu and Changhu Wang , “An Introduction to Image Synthesis with Generative Adversarial Nets” , 2018.
- Ian J. Goodfellow, Jean Pouget-Abadie\* , Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair , Aaron Courville, YoshuaBengio , “Generative Adversarial

Nets” , 2014.

- Jie Gui, Zhenan Sun, Yonggang Wen, Dacheng Tao, Jieping Ye , “A Review on Generative Adversarial Networks: Algorithms, Theory, and Applications “ , 2020.
- Pourya Shamsolmoali , Masoumeh Zareapoor , Eric Granger , HuiyuZhou , Ruili Wang , M. Emre Celebi and Jie Yang , “Image Synthesis with Adversarial Networks: a Comprehensive Survey and CaseStudies” , 2020.