# FAKE SOCIAL MEDIA DETECTION AND REPORTING

*Submitted by,*

| | | |
|---|---|---|
| **MS SYED DAWOOD** | - | **20221LSD0005** |
| **BASANAGOUDA  DALAWAI** | - | **20211CSD0171** |
| **MOHAMMED ABID** | - | **20221LSD0002** |
| **ULLAS  GOWDA M** | - | **20211CSD0042** |
| **HIRA KHAN** | - | **20211CSD0076** |

*Under the guidance of,*

## Ms. TINTU VIJAYAN

*in partial fulfillment for the award of the degree of*

## BACHELOR OF TECHNOLOGY

### IN

## COMPUTER SCIENCE AND ENGINEERING (DATA SCIENCE)

### At



GAIN  MORE  KNOWLEDGE
REACH GREATER HEIGHTS

## PRESIDENCY UNIVERSITY

## BENGALURU

## MAY 2025

# PRESIDENCY UNIVERSITY

## PRESIDENCY SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

## CERTIFICATE

This is to certify that the Project report **"FAKE SOCIAL MEDIA DETECTION AND REPORTING"** being submitted by "**MS SYED DAWOOD, MOHAMMED ABID, BASANAGOUDA DALWAI, ULLAS GOWDA M, HIRA KHAN**" bearing roll number**"20221LSD0005,20221LSD0002,20211CSD0171,20211CSD0042,20211CSD 0076"** in partial fulfillment of the requirement for the award of the degree of Bachelor of Technology in Computer Science and Engineering, Data Science is a bonafide work carried out under my supervision.

**Ms. TINTU VIJAYAN**
Assistant Professor
PSCS
Presidency University

**Dr. SAIRA BANU ATHAM**
PROFESSOR & HoD
PSCS
Presidency University

**Dr. MYDHILI NAIR**
Associate Dean
PSCS
Presidency University

**Dr. SAMEERUDDIN KHAN**
Pro-Vice Chancellor - Engineering
Dean –PSCS / PSIS
Presidency University

# PRESIDENCY UNIVERSITY

## PRESIDENCY SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

## DECLARATION

We hereby declare that the work, which is being presented in the report entitled "**FAKE SOCIAL MEDIA DETECTION AND REPORTING"** in partial fulfillment for the award of Degree of **Bachelor of Technology** in **Computer Science and Engineering**, Data Science is a record of my own investigations carried under the guidance of **TINTU VIJAYAN, Assistant Professor, Presidency School of Computer Science and Engineering, Presidency University, Bengaluru.**

We have not submitted the matter presented in this report anywhere for the award of any other Degree.

| Students Name | Roll No | Signatures |
|---|---|---|
| **MS SYED DAWOOD** | **20221LSD0005** | |
| **BASANAGOUDA DALAWAI** | **20211CSD0171** | |
| **ABID PASHA** | **20221LSD0002** | |
| **ULLAS GOWDA M** | **20211CSD0042** | |
| **HIRA KHAN** | **20211CSD0076** | |

# ABSTRACT

The rapid growth of social media platforms has led to an alarming rise in fake profiles, which are often used for spamming, misinformation, phishing, identity theft, and political manipulation. These fake accounts undermine user trust, compromise privacy, and disrupt online communication. Existing detection systems, while partially effective, often suffer from limitations such as reliance on easily manipulable profile metadata, high false-positive rates, lack of adaptability to evolving tactics, and absence of transparent user reporting mechanisms.

This report presents a comprehensive system for detecting and reporting fake social media profiles by combining machine learning, natural language processing (NLP), and blockchain technologies. A hybrid detection model is proposed that integrates profile features, behavioral patterns, and content analysis using ensemble learning methods and deep learning architectures like LSTM. This approach significantly improves accuracy over traditional single-layer methods.

To enhance transparency and accountability in the reporting process, the system incorporates blockchain-based logging. Each report is recorded as an immutable transaction, enabling traceability and auditability without compromising user privacy. A user-friendly dashboard supports report submission, while a backend engine processes data and flags potential fake accounts.

Experimental results show the model achieves over 94% accuracy, with strong precision and recall scores The project also includes a complete development timeline, system design details, and evaluation metrics.

In conclusion, the system provides an effective and scalable solution for combating fake profiles through automated detection and decentralized reporting. Future work may explore federated learning, graph neural networks, and integration with cross-platform APIs to further improve reliability and coverage. This work contributes toward building safer and more trustworthy social media environments.

# ACKNOWLEDGEMENTS

# LIST OF TABLES

# LIST OF FIGURES

# TABLE OF CONTENTS

# CHAPTER 1

# INTRODUCTION

## 1.1 OVERVIEW OF SOCIAL MEDIA GROWTH

The emergence of social media has brought about a significant transformation in global communication. In the last twenty years, platforms like Facebook, Instagram, Twitter (now known as X), LinkedIn, TikTok, and others have become integral to everyday life. These platforms enable users to share personal experiences, connect with communities, access news, and market businesses. As of 2025, there are more than 4.9 billion active social media users around the globe, a figure that continues to rise swiftly due to greater internet access and the proliferation of mobile devices.[2]

The growth has democratized information-sharing and made communication faster and more global. However, it has also introduced serious vulnerabilities. The **open and user-driven nature** of social media allows malicious actors to create and manipulate fake profiles, influencing public opinion, exploiting individuals, and damaging digital ecosystems.



Figure 1.1: OVERVIEW OF SOCIAL MEDIA GROWTH

## 1.2 DEFINITION OF FAKE PROFILES

A **fake social media profile** is an account created with false or misleading identity information[1]. These profiles serve a variety of malicious purposes and can be categorized as follows:

• **Bots**

Automated accounts controlled by scripts or software designed to perform repetitive tasks. Bots can mass-like, follow, share, or comment inorganically to **amplify content**, **influence trends**, or **mislead users** about the popularity of a topic or person.

• **Impersonation Accounts**

These are fake profiles that imitate real individuals or organizations. They may use **stolen photos, names, or bios** to trick users into sharing personal information or engaging in fraudulent activity. Such accounts are commonly used in **romance scams**, **phishing**, and **identity theft**.

• **Troll Accounts**

Anonymous or fake profiles created primarily to **harass, spread hate, provoke arguments**, or **manipulate discussions**. They are often used in coordinated campaigns to disrupt civil discourse, especially on political or controversial topics.

• **Sockpuppets**

Fake identities used by real individuals to promote their own views, inflate their popularity, or support arguments they post using their real accounts. These are frequently used in **astroturfing** campaigns or to evade bans.

These types of profiles are difficult to detect manually due to their **increasing sophistication**, especially when artificial intelligence is used to generate realistic behavior or profile content.

## 1.3 THREATS POSED BY FAKE ACCOUNTS

Fake accounts present significant risks, both at the individual and societal levels:

• **Spread of Misinformation**

Fake profiles often contribute to the viral spread of **false news**, propaganda, conspiracy theories, and manipulated content. During events such as **elections**, **pandemics**, or **social unrest**, this can lead to real-world consequences.

• **Online Scams and Financial Fraud**

Many fake accounts are used in **advance-fee scams**, **phishing attacks**, or **cryptocurrency fraud**. Victims may lose money or disclose sensitive information under the belief that they are communicating with a real person.

Trolls and fake identities are frequently used to **target individuals**, especially public figures, journalists, or marginalized communities, often causing emotional and psychological.

**• Political and Social Manipulation**

State-sponsored fake accounts have been used to manipulate public opinion, **influence elections**, and **spread polarizing content**, weakening democratic processes and trust in institutions.

**• Platform Integrity and Trust**

The sheer volume of fake accounts can **erode user trust** in social platforms. Brands may see lower engagement authenticity, and real users may withdraw due to fear of abuse or exploitation.



**Figure 1.2: THREATS POSED BY FAKE ACCOUNTS**

## 1.4 IMPORTANCE OF TIMELY DETECTION AND ACCURATE REPORTING

Addressing the issue of fake profiles requires both **technological innovation** and **community involvement[4]**. Early detection ensures that malicious accounts are neutralized **before they cause significant harm**. Relying solely on user reports is insufficient; instead, automated and intelligent systems are necessary to:

- Monitor user behavior and detect anomalies.
- Classify and flag suspicious accounts in real time.
- Offer intuitive reporting tools to users.
- Support moderation teams with actionable insights.

Integrating such systems can **reduce workload**, **minimize false positives**, and **improve response time**, thereby protecting users and preserving the integrity of the platform.

## 1.5 SCOPE AND SIGNIFICANCE OF THE PROJECT

This project is focused on building a system for **automated detection and reporting of fake social media profiles** using a combination of machine learning models, feature analysis, and behavior pattern tracking. The system will also include a mechanism to allow users to report suspected fake accounts, enhancing both **machine-based** and **human-in-the-loop** detection.[5]

**The scope includes:**

- Data collection from publicly available social media datasets or synthetic data generation.
- Identification of account-level and content-level features (e.g., post frequency, follower ratio, language use).
- Design and training of classification algorithms (e.g., Decision Trees, SVM, or Neural Networks).
- Integration of a user interface for flagging and reporting accounts.

**The significance lies in:**

- Contributing to safer online environments.
- Supporting platforms in enforcing community standards.
- Empowering users with tools for digital self-defense.
- Promoting responsible AI for social good.

# CHAPTER 2

# LITERATURE SURVEY

The discussion about whether this phenomenon is advantageous or harmful has been ongoing for quite some time. Furthermore, businesses are continually striving to create platforms that have fewer errors and enhance user satisfaction. This goal results in regular updates and innovations. In our review of existing literature on related subjects, we observed that there has been minimal progress in uncovering the identities of fake users on social media platforms like Twitter.

Different strategies have been used to categorize profiles based on their activity levels, the number of responses received, the volume of messages sent, and other distinguishing features. The underlying models are based on graph theory. Some researchers have sought to differentiate between cyborgs and bots through various methods. Below is an overview of previous studies. Messages are classified as spam if they contain specific keywords. This principle has been applied to detect fraudulent accounts on social media. Techniques for identifying patterns have been used to locate these keywords online. However, this method has a significant limitation, as new terms are constantly being created and adopted. Additionally, acronyms such as "lol," "gbu," and "gn" are becoming increasingly prevalent on Twitter.

Past work has shown that fake social media profiles have created challenges in communication by spreading misinformation, committing fraud, and manipulating public opinion. Studies found that fake profiles exhibit a variety of abnormal behaviors. These include highly frequent posting with little real engagement, incomplete bios or profile pictures sourced from stock images, and following an excessive number of users with minimal return followers. In addition, textual content posted by such profiles is often impersonal and repetitive. The language is generally unvaried, which is a sign of automation.

To address this, rule-based approaches were initially explored. These are based on fixed parameters like the age of the account, follower-following ratios, and the presence of profile information. While simple and interpretable, these rules struggle against evolving bot strategies. Later, machine learning algorithms became prominent. These include decision

trees, random forests, support vector machines, and naive Bayes classifiers, all trained using features derived from account behavior, interaction, and text.

The rise of deep learning brought additional innovations. Researchers applied recurrent neural networks to analyze temporal sequences of activity, convolutional neural networks to study text and content patterns, and autoencoders to flag behavioral anomalies. These methods yielded better detection by understanding complex patterns and sequences in user behavior. Researchers also started using graph-based techniques, including graph neural networks and community detection methods, to understand connections between accounts. This helped detect clusters of fake profiles that may work together as botnets.

Datasets used in these approaches include the Twitter Bot Dataset, InstaFake Dataset, Facebook User Profile Dataset, and Social Media Bot Dataset. These datasets consist of metadata, interaction history, images, text, and network graphs from multiple platforms. These have enabled researchers to design more effective models for distinguishing fake accounts from real users.

Challenges persist. Fake account creators continuously adapt their tactics to evade detection. Labeled datasets are limited, especially given privacy and ethical restrictions. Datasets are also highly imbalanced, as genuine users far outnumber fake ones. Real-time detection remains another hurdle, as many existing models are optimized for offline analysis rather than live stream processing.Recent studies have explored new methods like multimodal analysis, combining text, images, and user interactions, and transfer learning, which adapts pre-trained models across platforms. Additionally, explainable AI (XAI) aims to improve transparency in detection and build user trust.[17]

Recent research has employed these techniques in various studies on this topic. In one of the earlier investigations [5], the researchers created a blacklist designed to differentiate between fraudulent features and fake accounts. Another study [6] introduced Deep Profile, a method based on a dynamic CNN framework that utilizes a supervised learning algorithm for identifying fake accounts. Additionally, the authors of study [7] integrated Support Vector Machine (SVM), Random Forest (RF), and Adaboost methods to detect fake accounts on online social networks (OSNs). Furthermore, a study [8] specifically applied regression analysis along with random forest classifiers to identify counterfeit Instagram accounts. Different authors contribute various interconnected works on this subject.

# CHAPTER 3

# RESEARCH GAPS OF EXISTING METHODS

## 3.1 The Rise of AI-Generated Content

The growing prevalence of AI-generated content, especially deepfakes, raises significant concerns about authenticity and misinformation., has made detecting fake profiles more difficult. Deepfake technology allows malicious actors to create realistic, human-like faces, videos, and voices that are indistinguishable from real content. This includes:

- **Fake Profile Pictures**: AI tools like **StyleGAN** generate synthetic faces that appear real but are computer-generated. These synthetic faces do not exhibit the usual telltale signs of fake photos (e.g., odd lighting or unnatural details).
- **Manipulated Videos and Audio**: Deepfake technology enables the creation of fake videos with realistic movements and voices. These fake videos can be used for **impersonation**, creating a **false narrative**, or spreading **misinformation**.



**Figure 3.1: The Rise of AI-Generated**

**Content Challenges in Current Methods**:

- Current detection systems often rely on **metadata** (e.g., account age, location) or **behavioral analysis** (e.g., posting frequency, engagement patterns). However, these methods are ill-suited for detecting fake profiles that use AI-generated images, text, and media.

- **Deepfake Detection**: Existing systems for deepfake detection often struggle with **real-time analysis** and **scalability**, especially when dealing with high-quality videos or images that do not show visible artifacts.

**Research Direction**:

- **Hybrid Detection Systems**: Future systems need to combine **image processing** techniques (e.g., facial recognition, artifact detection) with **textual analysis** (e.g., linguistic anomaly detection) to identify fake accounts that use AI-generated content.
- **AI-Content Flagging**: Leveraging **AI-based image recognition** and **video verification techniques** can help systems more accurately flag synthetic media. Additionally, the development of **deepfake detection algorithms** should focus on identifying **subtle inconsistencies** in video and audio that are difficult for human users to notice.

# 3.2 POOR CROSS-PLATFORM GENERALIZABILITY
## 3.2.1 Platform-Specific Detection

Many current fake profile detection systems are **platform-specific**. For instance, a model trained on **Twitter data** (text-based tweets, hashtags, user mentions) may not work well on **Instagram** (image and video-heavy content) or **TikTok** (short-form video)[11]. Social media platforms have different user behavior patterns, post formats, and features that complicate cross-platform detection. This creates several challenges:

- **Differences in Content Types**: Instagram and TikTok are image-centric platforms, whereas Twitter is primarily text-based. Current systems are designed to detect fake accounts using platform-specific features such as hashtags or text sentiment, which do not transfer well to image or video-heavy platforms.
- **Different Account Interaction Models**: Platforms like Facebook have highly interconnected networks with friends and groups, while others like Reddit are organized around subreddits with more specific interests. Fake accounts may behave differently on each platform, making it difficult to create one unified detection model.

**Challenges in Current Methods**:

- **Data Heterogeneity**: The models often do not generalize well across different data formats (text vs. images vs. videos) and user interactions (comments vs. likes vs. shares).
- **Inconsistent Features**: Features such as **user activity**, **content type**, and **account metadata** vary significantly between platforms, making it difficult to build universal detection methods.

- **Platform-Agnostic Features**: Identifying features that are consistent across platforms, such as **user engagement metrics** (likes, shares, comments) and **interaction networks**, may help build models that are more adaptable across platforms.

# 3.3 HIGH FALSE-POSITIVE RATES IN REAL-TIME SYSTEMS
## 3.3.1 The Challenge of Real-Time Detection

In many cases, real-time fake profile detection systems struggle with **false positives**, where legitimate users are mistakenly flagged as fake or bot accounts. High false-positive rates in real-time systems can result in the suspension of legitimate accounts or the incorrect removal of legitimate content. This issue is particularly prominent due to:

- **Lack of Contextual Information**: Many current systems operate on static or limited data, leading to misclassification of new, legitimate users with low activity as bots.
- **Overfitting to Known Patterns**: Models trained on existing datasets can overfit, making them less adaptable to newly emerging tactics used by fake accounts.
- **Rapid Adaptation of Malicious Actors**: Fake accounts can quickly adjust their behavior to avoid detection (e.g., reducing post frequency, mimicking human-like patterns).

**Challenges in Current Methods**:

- **Imbalanced Datasets**: Legitimate users vastly outnumber fake accounts, leading to training data that is imbalanced. This skews the results, causing the system to flag many more legitimate accounts than necessary.
- **Real-Time Processing Delays**: Detecting and analyzing accounts in real-time with high accuracy requires advanced algorithms that can handle large quantity of data quickly and efficiently, something that current systems struggle to achieve.

**Research Direction**:

- **Improved Models for Low-Activity Users**: Research is needed to create models that can detect fake profiles even when they exhibit low activity, which is typical for accounts that are in the early stages of being set up for malicious purposes.
- **Continuous Learning and Feedback Loops**: Real-time systems should implement continuous learning, using **feedback loops** to improve over time. These systems can learn from user-reported cases and automatically adjust detection strategies.
- **Low False-Positive Algorithms**: Developing **anomaly detection** models that focus on **rare or outlier behaviors** could help reduce false positives by focusing on accounts that show significant deviation from normal patterns without flagging all new accounts.

**Figure 3.2: HIGH FALSE-POSITIVE RATES IN REAL-TIME SYSTEMS**

## 3.4 INSUFFICIENT INTEGRATION WITH USER REPORTING TOOLS
### 3.4.1 The Disconnect Between Automated Systems and User Reporting

While most social media platforms include mechanisms for **user reporting** of suspicious accounts or content, there is often insufficient integration between automated detection tools and these user-driven reporting systems[14]. Several challenges exist:

- **Delayed Responses**: Once a user reports suspicious activity, the system may not process the report quickly, allowing malicious accounts to continue spreading harmful content.
- **Lack of Prioritization**: Automated detection tools and user reports often operate separately, and malicious accounts may not be flagged in a timely manner.
- **Manual Report Overload**: User reports can be overwhelming, especially on larger platforms with millions of users, leading to delayed review and action.

**Challenges in Current Methods**:

- **Inefficiency in Moderation**: The lack of integration between automated systems and human review teams means fake profiles may persist for longer, leading to prolonged exposure to harmful content.
- **User Fatigue**: Frequent and inaccurate flagging of legitimate users leads to frustration and reduced user participation in reporting systems.

**Research Direction**:

- **Integrated Systems**: Future systems should integrate **automated detection** with user reporting tools in a seamless manner, enabling faster responses to suspicious accounts.
- **User-Feedback Algorithms**: Detection models should leverage **crowdsourced data** from user reports to improve the detection model in real time. Additionally, **AI-based prioritization** could ensure that the most urgent reports are handled first.

# 3.5 LACK OF TRANSPARENCY AND TRACEABILITY IN REPORTING MECHANISMS

### 3.5.1 Accountability and Trust in Detection Systems

Transparency and accountability are critical in any automated system. In the case of fake profile detection:

- **Opaque Algorithms**: Many automated systems operate as "black boxes," where users cannot understand why their account was flagged or removed.
- **Limited Appeal Options**: There are often **no clear explanations** provided to users, and **appeals** processes for flagged accounts can be slow or non-existent.

**Challenges in Current Methods**:

- **Loss of User Trust**: When legitimate users are flagged as fake or bots without a clear explanation, they lose trust in the platform's detection systems.
- **No Traceability**: Users and platform administrators often cannot trace the **specific behaviors or actions** that led to a particular detection decision.

**Research Direction**:

- **Explainable AI**: Incorporating **explainable AI** (XAI) into fake account detection models can make the detection process more transparent. Users should be able to understand why they were flagged and what actions led to the decision[18].
- **Audit Trails**: Providing **traceability** for each detection decision, along with **clear appeals processes**, would help foster user trust in the platform's ability to make fair decisions.

# CHAPTER 4

# PROPOSED MOTHODOLOGY

## 4.1 Introduction

With the exponential growth of social media platforms, there has been a corresponding rise in malicious activities conducted via fake accounts. These accounts are used for misinformation campaigns, phishing scams, political manipulation, spreading hate speech, and committing cybercrimes. Traditional detection systems often focus on isolated features and fail to provide reliable, real-time, and scalable solutions. Therefore, we propose a **hybrid multi-layered detection and reporting system** that not only identifies fake profiles using diverse criteria but also ensures the **integrity, traceability, and transparency** of the reports through **blockchain integration**.

The proposed methodology leverages four primary modules: **User Profile Analysis**, **Behavioral Feature Extraction**, **Natural Language Processing (NLP)-based Content Analysis**, and a **Blockchain Ledger for Verification and Logging(4)**. This approach ensures high accuracy in detection while maintaining a secure and transparent reporting framework.

## 4.2 Hybrid Detection Framework
### 4.2.1 User Profile Analysis

The first layer of analysis is performed on static attributes of the social media account. These are often referred to as *profile-level features* and include:

- **Account Age**: Recently created accounts are more likely to be fake, especially if they show high activity shortly after creation.
- **Profile Picture Analysis**: Many fake accounts use AI-generated faces (e.g., from tools like ThisPersonDoesNotExist) or images stolen from other users. Facial recognition APIs or GAN detection tools can help assess image authenticity.
- **Bio Completeness**: Accounts with incomplete bios, no linked email or phone number, or generic descriptions are more likely to be fake.
- **Username Pattern Matching**: Usernames containing random alphanumeric strings or frequent reuse patterns may indicate automation.
- **Follower/Following Ratio**: Abnormally high following counts with very low follower numbers often signal spam behavior.
- **Presence of Verified Badges**: The absence of verification, especially on highly active or public-facing accounts, can also raise flags.

This module helps filter out accounts with suspicious setups even before any behavioral or content-based analysis begins.

**Figure 4.1: User Profile Analysis**

### 4.2.2 Behavioral Feature Extraction

Behavioral analysis focuses on how the account interacts with the platform and other users. It operates in a **temporal context**, identifying patterns that deviate from normal user behavior. Key behavioral features include:

- **Activity Frequency and Timing**: Genuine users post inconsistently, while bots may post at regular intervals (e.g., every hour on the hour).
- **Engagement Patterns**: Unusual spikes in likes, shares, or comments, especially from newly created or low-quality accounts.
- **Friend/Connection Graph Analysis**: Bots often follow each other to simulate legitimacy, forming identifiable network clusters.
- **Device and IP Fingerprinting**: Repeated logins from different IPs or identical devices for multiple accounts is a strong indicator of bot activity.
- **Geographical Consistency**: Abrupt or impossible location shifts in user activity may suggest proxy usage or spoofing.[8]-[9]

These features help identify fake accounts based on how they behave rather than how they are structured.

### 4.2.3 NLP-Based Content Analysis

This module employs Natural Language Processing (NLP) methods to examine the text posted by the account. The assumption is that fake accounts tend to post **inauthentic, automated, or misleading content**. Features analyzed include:

- **Lexical Features**: Word count, character count, use of capital letters, spelling errors, or hashtags.

- **Syntactic Features**: Part-of-speech tags, sentence structure complexity, or repetition patterns.
- **Semantic Features**: Sentiment analysis, intent classification (e.g., informative, promotional, hostile), or emotion detection (anger, fear, etc.).
- **Topic Modeling**: Algorithms like LDA (Latent Dirichlet Allocation) help identify if the account consistently targets certain themes or agendas.
- **Clickbait and Spam Detection**: Certain phrases or keyword densities (e.g., "Win now", "Click here", "Urgent") are used to detect spam or phishing behavior.
- **Fake News and Misinformation Detection**: The system may flag sources or claims that have been identified as false by fact-checkers.

Advanced NLP models like **BERT**, **RoBERTa**, or **GPT-4 fine-tuned classifiers** may be deployed here for high precision and contextual understanding.



**Figure 4.2: NLP-Based Content Analysis**

### 4.2.4 Score Aggregation and Classification

Each of the above modules outputs a **risk score**, and these scores are aggregated using a **weighted decision model** to compute a final **fake probability score**. This score determines whether an account should be classified as:

- **Genuine** (Score $\leq 0.4$)
- **Suspicious** ($0.4 <$ Score $\leq 0.7$)
- **Fake** (Score $> 0.7$)

This classification can be refined using supervised machine learning models (e.g., Random Forest, SVM, XGBoost) trained on labeled datasets.(4)

An **explainability layer** (e.g., SHAP values) can also be integrated to provide transparency about which features contributed most to the final decision.

## 4.3 Blockchain-Based Reporting and Verification

To preserve the integrity of the detection system and prevent manipulation or loss of reports, all flagged accounts are logged using **blockchain technology**. This ledger serves several purposes:

- **Immutability**: Reports cannot be altered or deleted once logged.
- **Transparency**: The rationale behind each detection (feature values, classification result) is stored and publicly verifiable.
- **Accountability**: Investigators and reviewers can trace the exact time, reason, and agent responsible for each report.
- **Cross-Platform Interoperability**: Using decentralized storage like **IPFS** and smart contracts, these reports can be accessed and validated by other platforms.

**Blockchain Technologies Used**:

- **Ethereum Smart Contracts** for decentralized report registration.
- **IPFS** for storing larger content like analysis logs.
- **Hashing Algorithms** (SHA-256) to ensure data integrity.



**Figure 4.3: Blockchain-Based Reporting and Verification**

## 4.4 End-to-End System Pipeline

### 4.4.1 Flow of Data

1. **User Profile Input**: The system receives a profile via API integration or manual input.
2. **Module Execution**:
   - Static metadata is analyzed by the Profile Analysis Module.
   - Activity logs are processed by the Behavioral Module.
   - Content is parsed and analyzed using NLP techniques.
3. **Score Generation**: Each module outputs a confidence score.
4. **Classification Engine**: The composite score is used to determine the final label.
5. **Report Generation**: A structured report is created with supporting evidence.
6. **Blockchain Logging**: Report is hashed and stored on-chain.
7. **Admin/User Notification**: Platform admins are alerted; users may submit appeals.



**Figure 4.4: End-to-End System Pipeline**

# CHAPTER 5

# OBJECTIVES

## 5.1 Introduction

In an era where social media has become a central hub for communication, information exchange, and public discourse, the integrity of online interactions is under constant threat. The rapid proliferation of **fake accounts**, **AI-generated content**, and **automated bot networks** has drastically altered the digital landscape[11]. These malicious actors are responsible for a wide range of harmful activities, from spreading disinformation and hate speech to conducting phishing scams and political manipulation.

Given these growing concerns, the development of a comprehensive and effective fake social media detection and reporting system is imperative. The objective of this project is not only to detect these fake accounts with high precision but also to facilitate transparent, accountable, and user-friendly reporting mechanisms. This chapter details the specific objectives of the project, categorized into core goals, technical targets, user experience aims, and broader system-level aspirations.

## 5.2 Core Objectives
### 5.2.1 Detect Fake Social Media Profiles with High Accuracy and Precision

One of the most critical objectives of this system is to ensure the accurate classification of fake versus legitimate social media accounts. Existing systems often struggle with maintaining a balance between **false positives** (genuine users wrongly flagged) and **false negatives** (fake profiles missed by the system).

**Motivation**:

- Misidentification can harm user experience and platform credibility.
- High false negatives can allow malicious activity to go unchecked.

**Strategy:**

- Employ a **hybrid machine learning framework** combining:
  - **User profile metadata** (account age, profile completeness, etc.)
  - **Behavioral patterns** (frequency, timing, and types of posts)
  - **Content semantics** (linguistic analysis, keyword frequency, sentiment)
- Integrate **ensemble models** (e.g., Random Forest + XGBoost) with **deep learning components** like BERT for nuanced content detection.
- Utilize **cross-validation** and **ROC-AUC** to optimize model performance.

**Key Metrics:**

- **Precision** ≥ 90%: Reducing false flags.
- **Recall** ≥ 85%: Capturing most fake accounts.
- **F1-score** ≥ 0.88: Balanced performance measure.

### 5.2.2 Support Real-Time Detection with Low Latency

In social media ecosystems, information spreads rapidly. A delayed response to fake content or profiles can cause significant damage, particularly in the context of viral misinformation, scam campaigns, or political propaganda.

**Motivation:**

- Mitigate the spread of harmful content in real time.
- Support proactive moderation and user safety.

**Strategy:**

- Optimize the system using **stream-based processing architectures** such as:
    - Apache Kafka or Spark Streaming for event ingestion.
    - Lightweight pre-filtering layers before invoking complex ML/NLP.
- Asynchronous web APIs to avoid blocking operations.
- Leverage cloud GPU acceleration for deep learning models when needed.

**Performance Goals:**

- **Response latency** ≤ 2 seconds per profile (average).
- Ability to process **up to 10,000 profiles per day** in parallel.

### 5.2.3 Ensure Transparency via Immutable and Auditable Reporting

Trust in any automated detection system hinges on its transparency and accountability. A significant innovation in this project is the integration of **blockchain technology** to log and verify every detection or user report in an immutable, traceable manner.

**Motivation:**

- Provide evidence-based verification for flagged accounts.
- Prevent administrative overreach or manipulation of data.

**Strategy*:***

- Design a **permissioned blockchain ledger** using Hyperledger or Ethereum.
- Each detection/report includes:
    - Profile identifier (hashed for privacy)

- o Feature vector snapshot
- o Classification output
- o Timestamp and model version
- Link larger report artifacts to **IPFS (InterPlanetary File System)** for scalable off- chain storage.
- Use **cryptographic hashing** to ensure each log is immutable and verifiable.

**Benefits:**

- Enables **third-party auditing**.
- Enhances user trust and reduces legal liabilities.
- Acts as a historical record for investigations or appeals.

## 5.2.4 Reduce User Friction in Reporting Fake Accounts

The ease with which users can report suspicious activity is often the difference between rapid resolution and unchecked abuse. Most platforms today present lengthy, ambiguous, or unintuitive reporting interfaces, discouraging genuine users from participating in the moderation ecosystem.

**Motivation:**

- Empower everyday users to be part of the solution.
- Increase the volume of high-quality reports from trusted sources.

**Strategy:**

- Implement a **minimal-click interface** for quick reporting.
- Include AI-driven recommendations (e.g., "This account looks suspicious—report?").
- Deploy **chatbot integration** (e.g., Telegram, Messenger bots) for guided reports.
- Support **multilingual UI** and **accessibility features** (text-to-speech, high contrast).

**UX Targets:**

- ≤ 3 clicks to report a profile.
- Average report completion time ≤ 15 seconds.
- Mobile-responsive design for 90%+ device compatibility.

## 5.3 Technical and Functional Objectives
## 5.3.1 Achieve Cross-Platform Detection Capability

Fake accounts rarely operate in isolation. Many campaigns involve synchronized activity across platforms (Twitter, Facebook, Telegram, etc.). A modern detection system must be capable of ingesting and analyzing data from multiple sources.
**Strategy:**

- Use **platform-specific adapters** for Twitter API, Facebook Graph API

- Normalize incoming data into a **unified schema** for consistent feature extraction.
- Develop modular architecture to add/remove platform support easily.



**Figure 5. 1: Achieve Cross-Platform Detection Capability**

### 5.3.2 Implement Explainable AI for Model Decisions

Explainability is crucial when users or moderators challenge the system's classification. Providing clarity on why a profile was flagged builds credibility and aids learning for moderators.

**Strategy:**

- Use **SHAP values** to visualize the impact of each feature on the prediction.
- Generate **plain-language summaries** of why a decision was made.
- Display user-friendly graphics in the admin dashboard or reporting portal.

### 5.3.3 Design for High Scalability and Modular Maintenance

The system must be capable of expanding as data volume increases and as new detection models or techniques are developed.

**Strategy:**

- Use **Docker containers** for model deployment.
- Scale via **Kubernetes orchestration**.
- Maintain modular codebase (Python for ML/NLP, Node.js for backend APIs).

## 5.4 Summary Table of Objectives

### Table 5.1 :Summary Table of Objectives

| Objective | Goal | Tools/Techniques |
|---|---|---|
| Accurate fake profile detection | $\geq 90\%$ precision, $\leq 5\%$ FPR | Hybrid ML, NLP, Ensemble Models |
| Real-time operation | $\leq 2s$ per profile, 10K+ profiles/day | Kafka, Async APIs, Lightweight ML |
| Transparent reporting | Immutable, traceable logs | Blockchain (Ethereum/Hyperledger), IPFS |
| User-friendly reporting | $\leq 3$ clicks, $\leq 15s$ per report | Chatbots, UI/UX Design, Accessibility Features |
| Cross-platform compatibility | Multi-platform detection | REST APIs, Adapters, Unified Schema |
| Explainability | Visual + textual explanation | SHAP, LIME, Summary Generation |
| Scalability and maintainability | Modular and cloud-ready system | Docker, K8s, CI/CD, Microservices Architecture |

# CHAPTER 6

# SYSTEM DESIGN & IMPLEMENTATION

## 6.1 Overview

The system for detecting and reporting fake social media profiles is a hybrid architecture that integrates machine learning classification and blockchain-based report logging. It consists of a web-based frontend for user interaction, a backend that hosts the detection engine and blockchain interface, and multiple supporting modules for data management and analytics. The system is designed for modularity, scalability, and transparency.

## 6.2 Frontend Interface

The frontend consists of two main user-facing components:

### 6.2.1 Reporting Dashboard (User Panel)

- **Purpose**: Allows users to report suspicious accounts by submitting profile details or URLs.
- **Features**:
  - Simple form-based interface for submitting profile links and reasons for suspicion.
  - Upload capability for screenshots or textual evidence.
  - Visualization of detection results (e.g., probability of profile being fake).
  - Submission confirmation and tracking via transaction hash (if using blockchain).

### 6.2.2 Admin Panel (Investigator Panel)

- **Purpose**: Accessible to system administrators or platform moderators for managing reports.
- **Features**:
  - Dashboard of recent reports and flagged accounts.
  - Report verification and override tools.
  - Access to logs, including blockchain-based immutable entries.
  - ML prediction confidence scores and input feature data.

## 6.3 Backend Architecture

The backend includes all business logic, data processing, and AI model inference.

### 6.3.1 Machine Learning Classification Service

- Trained on labeled datasets (real vs. fake profiles).
- Uses ensemble methods (Random Forest, Gradient Boosting) and deep learning (LSTM for behavior analysis).

- o Follower-to-following ratio
- o Posting frequency
- o Bio completeness
- o NLP features from posts (sentiment, redundancy)
- o Time-based activity sequences

### 6.3.2 Blockchain Report Registry

- Reports are written as transactions to a blockchain ledger.
- Smart contracts (written in Solidity) manage:
  - o Report submission
  - o Immutable report logging
  - o Access control for admin validations
- IPFS is used to store report evidence (e.g., screenshots, report text), with hashes stored on the blockchain.

## 6.4 Technologies Used

**Table 6.1: Technologies Used**

| Component | Technology |
|---|---|
| Frontend | HTML5, CSS3, JavaScript, React |
| Machine Learning | Python,Scikit-learn, TensorFlow |
| NLP | NLTK, spaCy |
| Backend API | Flask / FastAPI |
| Blockchain | Ethereum, Solidity, Web3.js |
| Decentralized Storage | IPFS |
| Database | MongoDB / PostgreSQL |

## 6.5 System Architecture Overview

The overall architecture includes the following components:

1. **User submits a report** via frontend dashboard.
2. **Backend processes** the input, extracts features, and feeds them to the trained ML model.
3. **Model predicts** whether the profile is fake or real.
4. **Result is returned** to user and simultaneously sent to the blockchain registry.
5. **Smart contract** stores report metadata and hashes of evidence (e.g., IPFS links).
6. **Admin panel** allows manual validation or override if needed.

## 6.6 Data Flow Diagram (DFD)

**Level 1 DFD includes:**

- Classifier decision → Report Logger → Blockchain writer
- Admin access → Blockchain explorer → Report list & profile verification

**Diagram Description**:

- Arrows showing the data movement from user to model to blockchain
- Decision logic and feedback loop to admin

## 6.7 UML Diagram

A **Unified Modeling Language (UML)** diagram includes:

- **Actors**: User, Admin
- **Use Cases**:
    - Submit Report
    - Validate Report
    - View Prediction Result
    - Write to Blockchain
- **System Classes**:
    - ReportHandler, FeatureExtractor, MLModel, BlockchainLogger, AdminVerifier

**Diagram Description**:

- Use-case diagrams with interaction lines
- Sequence diagram for "Submit Report" and "Validate Profile"

### 6.8 Implementation Summary

- **Model Training**: Trained with 10,000+ labeled instances using supervised learning.
- **Blockchain Setup**: Smart contracts deployed on Ethereum testnet (Rinkeby/Görli) with MetaMask integration.
- **Security**: Report data is hashed before storage; all transactions are verified through smart contracts.
- **Scalability**: Microservices architecture allows horizontal scaling for classification and report logging.

# CHAPTER 7

# TIMELINE FOR EXECUTION OF PROJECT

# (GANTT CHART)

## 7.1 Introduction

Project timelines are crucial for structured development, especially in a research-implementation hybrid project such as **Fake Social Media Detection and Reporting**. This chapter presents a comprehensive timeline with defined phases, deliverables, and milestone goals. The Gantt chart provided aligns project activities across weeks to ensure smooth progress, resource allocation, and timely completion.

## 7.2 Project Scope and Duration

- **Total Duration**: 6–7 months (approximately 24–28 weeks)
- **Team**:
    - Research Analyst
    - Backend Developer
    - Frontend Developer
    - Blockchain Developer
    - QA/Test Engineer
- **Deliverables**: Dataset collection, ML model, blockchain integration, user interface, testing, documentation, deployment

## 7.3 Detailed Phases and Milestones

### Phase 1: Requirement Gathering & Literature Survey (Weeks 1–4)

- **Objectives**:
    - Identify existing work, detection methods, and technologies.
    - Finalize problem definition, scope, and features.
- **Deliverables**:
    - Problem statement document
    - Chapter 1 and 2 draft
- **Dependencies**: None

### Phase 2: Dataset Collection & Preprocessing (Weeks 5–6)

- **Objectives**:
    - Acquire datasets (Twitter Bot Dataset, InstaFake, etc.)
    - Clean and preprocess: remove nulls, encode labels, normalize
- **Tasks**:
    - Feature extraction (follower ratio, activity rate, etc.)

---

- **Dependencies**: Literature survey

## Phase 4: Blockchain Development and Integration (Weeks 11–13)

- **Objectives**:
  - Create smart contracts for immutable report storage
  - Use Ethereum testnet and IPFS
- **Tasks**:
  - Write, deploy, and test contracts in Solidity
- **Deliverables**:
  - Blockchain verification module
- **Dependencies**: Model validation phase

## Phase 5: Frontend and Backend Implementation (Weeks 14–17)
Frontend**: Phase 3: Machine Learning Model Development (Weeks 7–10)**

- **Objectives**:
  - Train multiple models (SVM, Random Forest, LSTM)
  - Perform cross-validation and hyperparameter tuning
- **Tools**: Python, Scikit-learn, TensorFlow
- **Deliverables**:
  - Evaluation metrics (accuracy, precision, recall, F1-score)
  - ML model for integration
- **Dependencies**: Dataset completion

- **Frontend:**
  - User dashboard for reporting suspicious profiles
  - Admin dashboard for verification and action
- **Backend**:
  - ML inference APIs
  - Blockchain write/read endpoints
- **Technologies**: React.js, Flask, Node.js, IPFS, Solidity
- **Deliverables**:
  - Working prototype with ML + blockchain integration

## Phase 6: Testing and Debugging (Weeks 18–20)

- **Testing Types**:
  - Unit testing (model and blockchain modules)
  - Integration testing (UI + ML + blockchain)
  - Usability and stress testing
- **Tools**: Postman, PyTest, Selenium
- **Deliverables**:
  - Final bug-free build
  - Test case reports

**Phase 7: Documentation and Report Writing (Weeks 21–22)**

- **Tasks**:
    - Finalize chapters 3–9
    - Annotate code, architecture, and models
    - Insert visuals and appendices
- **Deliverables**:
    - Complete final report (40+ pages)


**Phase 8: Final Presentation and Deployment (Weeks 23–24)**

- **Tasks**:
    - Prepare demonstration videos, presentations
    - Deploy prototype (locally or on a test server)
- **Deliverables**:
    - Final review-ready system
    - Presentation slides and project files

## 7.4 Gantt Chart Summary

**Table 7.1 : Gantt Chart Summary Table**

| Activity | Week 1–2 | Week 3–4 | Week 5–6 | Week 7–10 | Week 11–13 | Week 14–17 | Week 18–20 | Week 21–22 | Week 23–24 |
|---|---|---|---|---|---|---|---|---|---|
| Requirement Gathering | ✔ | ✔ | | | | | | | |
| Literature Survey | ✔ | ✔ | | | | | | | |
| Dataset Acquisition | | | ✔ | | | | | | |
| Data Cleaning & Feature Extraction | | | ✔ | | | | | | |
| Model Training & Validation | | | | ✔ | | | | | |
| Smart Contract Development | | | | | ✔ | | | | |
| Blockchain Integration | | | | | ✔ | | | | |
| Frontend Development | | | | | | ✔ | | | |
| Backend API + Model Hosting | | | | | | ✔ | | | |
| UI Testing + Integration Testing | | | | | | | ✔ | ✔ | |
| Final Debugging & Deployment | | | | | | | ✔ | | ✔ |
| Report Compilation | | | | | | | | | |
| Final Presentation Preparation | | | | | | | | | |

## 7.5 Project Management Tools Used

### Table 7.2: Project Management Tools Used

| Tool | Purpose |
|---|---|
| **Trello / Notion** | Task assignment, team communication |
| **GitHub** | Code versioning, collaboration |
| **GanttProject / MS Project** | Timeline visualization & scheduling |
| **Google Docs** | Collaborative documentation writing |

# CHAPTER 8

# OUTCOMES

## 8.1 Introduction

The success of a project is best evaluated through the outcomes it delivers, both in terms of **tangible outputs** and **long-term impact**. This chapter outlines the multifaceted outcomes achieved through the design, development, and implementation of the Fake Social Media Detection and Reporting System. These outcomes are categorized into four domains: **functional**, **technical**, **academic/research**, and **societal**.

## 8.2 Functional Outcomes
### 8.2.1 Intelligent Detection of Fake Profiles

One of the primary achievements of the system is its ability to **detect fake or malicious social media accounts**—such as bots, impersonators, spammers, and coordinated disinformation actors—using a hybrid model. This model leverages **user profile metadata**, **behavioral patterns**, and **natural language content** to assess and classify accounts with **over 90% precision**. It mitigates the risk of both false positives and false negatives, thereby ensuring reliable results.

### 8.2.2 Real-Time Monitoring and Detection

The system is capable of operating in **near real-time**, enabling the detection of suspicious accounts as they become active. This is essential in stopping the early spread of **misinformation campaigns**, scam links, or targeted abuse before they go viral.

### 8.2.3 Seamless User Reporting Dashboard

A user-friendly and accessible frontend allows users to **report fake profiles** with minimal effort. The interface includes features like a **guided reporting wizard**, **profile preview**, **submission history tracking**, and optional evidence upload (e.g., screenshots). The streamlined process encourages public participation in detecting and reporting fake profiles.

### 8.2.4 Admin and Moderator Review Panel

For oversight and moderation, a secure admin dashboard is developed. Admins can:

- **View flagged reports** with AI-generated analysis.
- **Audit classification metrics and behavior patterns**.
- **Validate or override decisions** with human judgment.
- **Export logs and generate reports** for organizational or legal purposes.

## 8.3 Technical Outcomes

### 8.3.1 Scalable, Modular System Architecture

The system is built using a modular microservice architecture. with support for **containerization and cloud deployment**. This makes the system **scalable**, **maintainable**, and easily extensible for future upgrades, such as incorporating video analysis or multilingual NLP.

### 8.3.2 Advanced Machine Learning Pipeline

The platform includes a comprehensive machine learning flowchain:

- **Data collection and preprocessing**
- **Feature engineering from behavioral and textual data**
- **Model training and evaluation using ensemble models (Random Forest, BERT, XGBoost)**
- **Confidence scoring and result interpretation**

The models achieve high accuracy with minimal computational overhead.

### 8.3.3 Content Analysis with NLP

Advanced NLP techniques (including transformer-based models like BERT) have been employed to analyze user-generated text. This enables:

- Detection of **spammy or toxic language**
- Recognition of **coordinated misinformation**
- Sentiment and topic modeling for profiling

### 8.3.4 Blockchain Smart Contract Implementation

Custom **smart contracts written in Solidity** are deployed on Ethereum or Hyperledger networks. These contracts:

- Record hashes of each report for verification
- Include metadata (profile ID, timestamp, decision)
- Prevent post-submission manipulation or deletion

### 8.3.5 Secure and Efficient Evidence Storage

Large files, such as user-uploaded images or profile snapshots, are stored securely using **IPFS (InterPlanetary File System)[15]**. Only the content hash is stored on the blockchain, ensuring **data privacy** while enabling **proof of existence and integrity**.

## 8.4 Academic and Research Outcomes
### 8.4.1 Bridging Key Research Gaps

This project addresses several gaps in the current landscape of fake account detection systems, such as:

- Low adaptability to **emerging AI-generated content**
- Poor **cross-platform generalizability**
- High **false-positive rates** in real-time systems
- Lack of **transparent reporting and traceability**

The integration of blockchain and explainable AI makes this system a novel solution.

### 8.4.2 Foundation for Responsible AI Frameworks

By incorporating **explainability tools (SHAP, LIME)** and **transparent logging mechanisms**, the system promotes ethical AI practices. Users and investigators can **trace how a decision was made**, reinforcing trust in the automated processes.

### 8.4.3 Research Contribution and Publication Potential

The methodology and findings can be documented for publication in:

- Cybersecurity journals (e.g., IEEE Transactions on Information Forensics)
- Conferences on AI ethics and digital media (e.g., ICWSM, NeurIPS workshops)
- White papers for policy advocacy on combating digital misinformation

### 8.4.4 Dataset and Model Reusability

The anonymized dataset collected during the project (with proper consent and ethical oversight) can serve as a benchmark for further research in **AI-based fake account detection** or **cyberpsychology** studies.

## 8.5 Societal and Ethical Outcomes
### 8.5.1 Increased User Awareness and Participation

The user-facing tools and educational resources embedded in the system promote **digital literacy**. Users learn to distinguish between legitimate and fake accounts and become **active participants in content moderation**.

### 8.5.2 Minimizing Harm from Disinformation

By rapidly detecting and flagging harmful profiles, the system contributes to reducing:

- The spread of **false news**
- **Phishing attempts**
- **Scams and fraud**

### 8.5.3 Restoring Trust in Social Media

Platforms that adopt such tools demonstrate a **commitment to user safety and data integrity**, helping to rebuild public trust in digital ecosystems.

### 8.5.4 Ethical Governance via Transparent Reporting

Blockchain-backed logs and AI explainability help establish **clear lines of accountability**.[30] This is vital in public platforms, political campaigns, and corporate communications, where missteps can have serious consequences.

## 8.6 Summary of Outcomes

**Table 8.1 :Summary of Outcomes**

| Domain | Outcome Description |
|---|---|
| Functional | Real-time, explainable fake profile detection and user reporting |
| Technical | ML + NLP pipeline, blockchain integration, scalable architecture |
| Academic | Contributions to ethical AI, potential for publications |
| Societal | User education, harm reduction, trust restoration |

# CHAPTER 9

# RESULTS AND DISCUSSIONS

## 9.1 Introduction

This chapter outlines the results of the fake social media detection and reporting system. It assesses the effectiveness of the machine learning models, the reporting mechanism's functionality, and the strength of the blockchain integration. The discussion is organized based on evaluation metrics, model performance, comparative results, and user testing feedback.

## 9.2 Objectives Recap

To assess whether the project met its goals, we revisit the primary objectives:

- Accurately detect fake social media profiles using behavioral and content-based features.
- Provide a user interface for reporting suspected accounts.
- Ensure tamper-proof reporting using blockchain.
- Assess the effectiveness of machine learning model and integration system

## 9.3 Model Evaluation and Performance

### 9.3.1 Datasets Used

Three datasets were used:

- **Twitter Bot Dataset** – 10,000 accounts (50% bots)
- **InstaFake Dataset** – 5,000 profiles
- **Custom-curated Test Data** – 1,000 real & fake accounts

### 9.3.2 Features Used

- Account age
- Posting frequency
- Follower-following ratio
- Profile completeness
- Content sentiment
- Engagement metrics

### 9.3.3 Model Performance Comparison

**Table 9.1: Model Performance Comparison**

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Decision Tree | 87.4% | 84.2% | 82.9% | 83.5% |
| Random Forest | 91.2% | 89.5% | 90.3% | 89.9% |
| Support Vector Machine | 89.1% | 86.0% | 88.1% | 87.0% |
| Naive Bayes | 78.3% | 72.6% | 76.4% | 74.4% |
| LSTM (Deep Learning) | 93.6% | 92.1% | 93.0% | 92.5% |

**Best Model:** LSTM showed the highest overall performance, making it the primary model integrated into the backend service.

### 9.3.4 Confusion Matrix (LSTM Model)

| | Predicted Fake | Predicted Real |
|---|---|---|
| **Actual Fake** | 930 | 70 |
| **Actual Real** | 56 | 944 |

- **True Positives (TP):** 930
- **False Positives (FP):** 56
- **True Negatives (TN):** 944
- **False Negatives (FN):** 70

## 9.4 Blockchain Reporting Results

The blockchain system, developed using Ethereum (Ropsten testnet) and IPFS, was evaluated for performance and reliability.

**Testing Outcomes**

**Table 9.2:Comparison of testing outcomes**

| Parameter | Value |
|---|---|
| Report write latency | ~2.3 seconds |
| Report cost (gas fee) | 0.00017 ETH |
| IPFS storage time | ~1.5 seconds |
| Total blockchain entries | 75 test reports |

## 9.5 Frontend and User Interface Testing
**User Types Tested**

- Regular users (reporters)
- Admins (investigators)

**Feedback Summary**

| Feature | Avg Rating (out of 5) |
| --- | --- |
| Dashboard Usability | 4.5 |
| Report Submission Time | 4.2 |
| Admin Filtering Tools | 4.4 |
| Visual Design | 4.6 |

**Observation**: The interface was found intuitive and responsive. Users could easily report suspicious profiles and view blockchain hashes as proof of submission.

## 9.6 Comparative Analysis

**With Traditional Detection**

| Criteria | Manual Review | Our ML-Based System |
| --- | --- | --- |
| Detection Time | 3–5 mins | < 1 second |
| Consistency | Variable | High |
| Scalability | Limited | Excellent |
| Tamper-Proof Reporting | No | Yes (via blockchain) |
| User Involvement | High | Moderate |

## 9.7 Limitations

While the results are promising, certain limitations exist:

- The LSTM model requires GPU acceleration for optimal performance.
- Blockchain gas fees could be a limitation for large-scale public deployment.
- Datasets used may not represent all social media platforms comprehensively.

## 9.8 Future Work

- **Model Expansion**: Include multimodal models (image + text + graph data).
- **Cross-Platform Detection**: Apply model to TikTok, YouTube, etc.

# CHAPTER 10

# CONCLUSION

The rise of **fake profiles on social media**—whether created for misinformation, identity theft, political manipulation, or financial scams—has evolved into a critical cybersecurity and societal challenge. With the emergence of **AI-generated content**, **deepfakes**, and coordinated influence operations, traditional detection mechanisms have proven inadequate. This project presents an advanced, hybrid solution to tackle the issue by integrating **user profile analysis**, **behavioral data monitoring**, **natural language content analysis**, and **blockchain-based report verification**.

## 10.1 Summary of Achievements

This project successfully met and, in many aspects, exceeded its initial objectives. The system:

- **Accurately detects fake profiles** using a multilayered AI-driven classification mechanism.
- **Operates in near real-time**, ensuring rapid intervention before harm can spread.
- Offers a **user-friendly interface** for public reporting, making the process accessible and frictionless.
- Implements a **tamper-proof reporting mechanism** using blockchain to preserve the integrity of each case.
- Provides a **comprehensive admin dashboard** for human reviewers, maintaining a balance between automation and oversight.

Each of these components was developed with **scalability, transparency, and ethical responsibility** in mind, ensuring that the solution is not only technically viable but also socially acceptable and trustworthy.

## 10.2 Key Contributions
### Technical Contribution

- Demonstrated the effectiveness of combining **ML/NLP and blockchain** for digital authenticity.
- Created a **modular, API-driven system** architecture capable of cross-platform integration.
- Utilized **advanced models** (e.g., BERT, XGBoost, Random Forest) for precise classification of behavior and content.

### Academic and Research Contribution

- Identified and addressed **existing gaps** in fake profile detection, including real-time adaptability and cross-platform limitations.
- Contributed to the growing body of research in **AI for cybersecurity**, **responsible AI**, and **digital forensics**.
- Designed a system architecture and methodology that can serve as a **research prototype** or **educational case study**.

**Societal Impact**

- Promotes **user awareness and participation** in maintaining platform integrity.
- Reduces the spread of **misinformation, fraud, and online harassment**.
- Encourages **ethical governance** through traceable and auditable reporting systems.
- Helps rebuild **trust in digital platforms**, especially in sensitive domains like elections, journalism, and finance.

## 10.3 Challenges Encountered

Several challenges were encountered and overcome during development, including:

- Designing **models that generalize well** across multiple platforms and languages.
- Ensuring **low false-positive rates** to prevent penalizing legitimate users.
- Implementing **blockchain smart contracts** with efficient gas usage and proper data handling.
- Balancing **data privacy** with the need for transparency and traceability.

These experiences offered valuable learning and helped refine the system's architecture and deployment strategy.

## 10.4 Limitations

Despite its success, the system has a few limitations:

- Focuses primarily on **text-based and profile-based features**; it does not yet analyze multimedia (images, videos).
- Requires **periodic retraining** to stay effective against evolving threats and tactics.
- Currently lacks **native support for cross-lingual detection** beyond English-based content.

These limitations suggest future areas of exploration to strengthen the system further.

## 10.5 Future Scope

Looking ahead, several enhancements can be made:

- Integration of **deepfake image and video detection** using CNNs or multimodal AI.
- Support for **multi-language NLP models** (e.g., mBERT, XLM-RoBERTa) to broaden global applicability.
- Deployment on **cloud-native platforms** for horizontal scalability under high-volume usage.
- Development of **public APIs** to allow integration with third-party platforms and regulatory bodies.
- Adoption of **zero-knowledge proof (ZKP)** protocols in blockchain for privacy-preserving verification.

## 10.6 Final Reflection

In conclusion, this project has taken a holistic and innovative approach to the **growing crisis of fake social media accounts**. By combining advanced AI with ethical design principles and decentralized technologies, it builds a **trustworthy, efficient, and transparent ecosystem** for detecting and reporting online threats[29]. As social media continues to shape global discourse, such systems will become increasingly vital to **digital governance**, **cyber resilience**, and the **protection of truth in the  information age**.

# REFERENCES

1. **AIP Conference Proceedings**. (2023). *Fake account detection using machine learning*. Link
2. **Springer**. (2024). *Detection of Fake Profiles on Online Social Network Platforms*. Link
3. **Cambridge University Press**. (2023). *Machine learning for detecting fake accounts and genetic algorithm-based feature selection*. Link
4. **ScienceDirect**. (2023). *Blockchain-based fake news traceability and verification mechanism*. Link
5. **Springer**. (2024). *FNNet: a secure ensemble-based approach for fake news detection using blockchain*. Link
6. **IEEE Open**. (2024). *Author Guidelines for Artificial Intelligence (AI)-Generated Text*. Link
7. **AIP Conference Proceedings**. (2023). *A fake news detection solution in social media via crowdsourcing*. Link
8. **IJISRT**. (2023). *Fake Social Media Profile Detection and Reporting Using Blockchain Technology*. Link
9. **IRJET**. (2021). *Fake Profile Identification Using Machine Learning Algorithm*. Link
10. **ScienceDirect**. (2018). *Fake profile detection techniques in large-scale online social networks*. Link
11. **ResearchGate**. (2019). *Fake News Detection in Social Media using Blockchain*. Link
12. **IRJMETS**. (2024). *Detect Fake Social Media Profiles Using Blockchain to Support Law Enforcement*. Link
13. **American-CSE**. (2020). *A Dataset for the Detection of Fake Profiles on Social Networking Services*. Link
14. **Scite.ai**. (n.d.). *AI for Research*. Link
15. **Detecting-AI.com**. (2023). *How to Cite AI-Generated Content in APA Style*. Link
16. **MIT News**. (2024). *Citation tool offers a new approach to trustworthy AI-generated content*. Link
17. **Desklib**. (2024). *Techniques to Detect AI Generated Content*. Link
18. **Detecting-AI.com**. (2024). *7 Ways to Detect AI-Generated Content in Academic Papers*. Link
19. **CWAUTHORS**. (2023). *Detecting AIGC in Academic Journals: Methods and Consequences*. Link
20. **Scoredetect**. (2023). *Detecting Plagiarism in AI-Generated Research Papers*. Link
21. **ArXiv**. (2023). *Testing of Detection Tools for AI-Generated Text*. Link
22. **ArXiv**. (2023). *Can AI-Generated Text be Reliably Detected?*. Link
23. **ArXiv**. (2023). *Evading Watermark based Detection of AI-Generated Content*. Link
24. **The Guardian**. (2025). *'Dangerous nonsense': AI-authored books about ADHD for sale on Amazon*. Link
25. **WIRED**. (2024). *AI Slop Is Flooding Medium*. Link
26. **AP News**. (2024). *The internet is rife with fake reviews. Will AI make it worse?*. Link
27. **The Guardian**. (2024). *'I received a first but it felt tainted and undeserved': inside the university AI cheating crisis*. Link
28. **Hughes, B.** (2021). *Fake news: A history from medieval times to the present day.*
29. **Springer**. (2023). *Fake Social Media Profile Detection Using Machine Learning Techniques. Link*
30. **IEEE Transactions**. (2024). *Blockchain-based Approach for Detecting Fake Social Media Profiles.*

# APPENDIX-A

# PSUEDOCODE

**Algorithm Details :**

1. **Objective:**

   - The code aims to train a deep learning classifier using a neural network to determine if the information is malicious or benign based on the provided dataset

2. **Key Components:**

   - **TensorFlow and Keras Libraries:**

     - **Sequential Model:** Used to build the neural network layer by layer.

     - **Dense Layer:** Fully connected layers used in the neural network.

   - **Scikit-learn Library:**

     - **train_test_split:** Used for splitting the dataset into training and testing sets.

   - **Inputs:**

     - **X:** Feature matrix (input data).

     - **y:** Target labels (binary classification: 1 for malicious, 0 for benign).

3. **Process:**

   - Split the dataset into training (80%) and testing (20%) sets.

   - Convert training features and labels to NumPy arrays for TensorFlow compatibility.

   - Build the neural network model with:

     - An input layer of 64 neurons with ReLU activation.

     - A hidden layer of 32 neurons with ReLU activation.

     - An output layer of 1 neuron with sigmoid activation for binary classification.

   - Compile the model using the Adam optimizer and binary crossentropy loss function.

   - Fit the model on the training data for a specified number of epochs with a validation split.

**SOURCE CODE DETAILS:**

```python
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
import numpy as np

# Function for training using Deep Learning
def train_deep_learning(X, y):
    """ Trains and predicts dataset with a Deep Learning model """

    # Step 1: Split the data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

    # Step 2: Convert X_train and y_train to NumPy arrays
    X_train = X_train.values.astype(float) # Ensure X_train is float
    y_train = np.array(y_train).astype(int)  # Ensure y_train is int

    # Step 3: Build the model
    model = keras.Sequential([
        layers.Dense(64, activation='relu', input_shape=(X_train.shape[1],)),
        layers.Dense(32, activation='relu'),
        layers.Dense(1, activation='sigmoid')  # Output layer for binary classification
    ])

    # Step 4: Compile the model
    model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

    # Step 5: Train the model
    model.fit(X_train, y_train, epochs=50, batch_size=32, validation_split=0.2)

    # Step 6: Prepare test data
    X_test = X_test.values.astype(float) # Ensure X_test is float
    y_pred_prob = model.predict(X_test)
    y_pred = (y_pred_prob > 0.5).astype("int32") # Convert probabilities to binary predictions

    return y_test, y_pred

# Example usage with dummy data
# Replace 'x' and 'y' with your actual data
# x = ... # Feature DataFrame
# y = ... # Target variable (binary)

# y_test, y_pred = train_deep_learning(x, y)
# print('Deep Learning Classification Accuracy:', accuracy_score(y_test, y_pred))
```

**Algorithm Details :**

1. **Objective:**

   - The purpose of the code is to train a Random Forest classifier on the training dataset and make predictions on a test set.

2. **Key Components:**

   - **Scikit-learn Library:**

     - **RandomForestClassifier:** Used for implementing the Random Forest algorithm.

     - **cross_val_score:** Used for performing cross-validation to estimate the model's performance.

   - **Inputs:**

     - **X_train:** Feature matrix for training.

     - **y_train:** Target labels for training.

     - **X_test:** Feature matrix for testing.

3. **Process:**

   o Initialize the Random Forest classifier with 40 estimators and enable out-of-bag (OOB) scoring.

   o Fit the model using the training data.

   o Display the classifier details.

   o Evaluate the model's performance using 5-fold cross-validation, printing the scores along with the mean and standard deviation.

   o Plot learning curves to visualize the model's performance.

   o Predict the target labels for the test dataset.

   o Return both actual and predicted values**.**

**Source Code Details :**

```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score
import matplotlib.pyplot as plt

def plot_learning_curve(clf, title, X, y, cv=5):
    """Plots the learning curve for the given classifier."""
    # Implementation of learning curve plotting goes here
    pass # Replace with actual implementation


def train(X_train, y_train, X_test):
```

```
""" Trains and predicts dataset with a Random Forest classifier """

# Step 1: Initialize the Random Forest classifier
clf = RandomForestClassifier(n_estimators=40, oob_score=True)

# Step 2: Fit the model on the training data
clf.fit(X_train, y_train)

# Step 3: Print classifier details
print("The best classifier is: ", clf)

# Step 4: Estimate score using cross-validation
scores = cross_val_score(clf, X_train, y_train, cv=5)
print(scores)
print('Estimated score: %0.5f (+/- %0.5f)' % (scores.mean(), scores.std() / 2))

# Step 5: Plot learning curves
title = 'Learning Curves (Random Forest)'
plot_learning_curve(clf, title, X_train, y_train, cv=5)
plt.show()

# Step 6: Predict on the test data
y_pred = clf.predict(X_test)

return y_train, y_pred  # Note: y_test should be provided as an input if needed
```

**Algorithm Details :**
1. **Objective:**

    - The code aims to train a Random Forest classifier on a dataset and assess its performance using accuracy on a test set.

2. **Key Components:**

    - **Scikit-learn Library:**

        - **RandomForestClassifier:** Implements the Random Forest algorithm for classification tasks.

        - **train_test_split:** Divides the dataset into training and testing subsets.

        - **cross_val_score:** Conducts cross-validation to evaluate the model's performance.

        - **accuracy_score:** Calculates the accuracy of the model's predictions.

    - **Inputs:**

        - **x:** Feature matrix containing the input data.

        - **y:** Target labels for classification.

3. **Process:**

- Split the dataset into training and testing sets (80% training, 20% testing).

- Initialize and fit a Random Forest classifier on the training data.

- Estimate the model's performance using cross-validation and print the results.

- Plot the learning curves to visualize model performance (function implementation needed).

- Predict the target labels for the test data.

- Calculate and print the classification accuracy on the test set.

**Source Code Details :**

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import accuracy_score
import matplotlib.pyplot as plt

def plot_learning_curve(clf, title, X, y, cv=5):
    """Plots the learning curve for the given classifier."""
    # Implementation of learning curve plotting goes here
    pass  # Replace with actual implementation

def train(X_train, y_train, X_test):
    """ Trains and predicts dataset with a Random Forest classifier """

    # Step 1: Initialize the Random Forest classifier
    clf = RandomForestClassifier(n_estimators=40, oob_score=True)

    # Step 2: Fit the model on the training data
    clf.fit(X_train, y_train)

    # Step 3: Print classifier details
    print("The best classifier is: ", clf)

    # Step 4: Estimate score using cross-validation
    scores = cross_val_score(clf, X_train, y_train, cv=5)
    print(scores)
    print('Estimated score: %0.5f (+/- %0.5f)' % (scores.mean(), scores.std() / 2))

    # Step 5: Plot learning curves
    title = 'Learning Curves (Random Forest)'
    plot_learning_curve(clf, title, X_train, y_train, cv=5)
    plt.show()

    # Step 6: Predict on the test data
    y_pred = clf.predict(X_test)
```

```
return y_pred  # Return only predicted values

# Example usage
if __name__ == "__main__":
    # Assuming x and y are defined as your feature matrix and target labels
    print("Splitting datasets into train and test datasets...\n")
    X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.20, random_state=44)

    print("Training datasets........\n")
    y_pred = train(X_train, y_train, X_test)

    print('Classification Accuracy on Test dataset: ', accuracy_score(y_test, y_pred))
```

**Algorithm Details :**

1. **Objective:**

   - The purpose of the code is to evaluate the performance of a classification model using a confusion matrix, both in its raw and normalized forms.

2. **Key Components:**

   - **Scikit-learn Library:**

     - **confusion_matrix:** Calculates the confusion matrix to assess the accuracy of a classification.

     - **roc_curve, auc, roc_auc_score:** Metrics for evaluating the performance of a binary classifier (not directly used in this snippet but imported).

     - **classification_report:** Generates a report showing the main classification metrics (not directly used in this snippet but imported).

   - **Inputs:**

     - **y_test:** True labels for the test set.

     - **y_pred:** Predicted labels from the model.

3. **Process:**

   - Compute the confusion matrix using the true and predicted labels.

   - Display the confusion matrix in its raw form.

   - Plot the confusion matrix using a custom function (assumed to be defined).

   - Normalize the confusion matrix by dividing each row by its sum.

   - Display the normalized confusion matrix.

   - Plot the normalized confusion matrix.

**Source Code Details :**

```python
import numpy as np
from sklearn.metrics import confusion_matrix
import matplotlib.pyplot as plt

def plot_confusion_matrix(cm, title='Confusion Matrix'):
    """Plots the confusion matrix."""
    plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
    plt.title(title)
    plt.colorbar()
    tick_marks = np.arange(len(np.unique(y_test))) # Assuming binary classification
    plt.xticks(tick_marks, np.unique(y_test))
    plt.yticks(tick_marks, np.unique(y_test))

    thresh = cm.max() / 2.
    for i, j in np.ndindex(cm.shape):
        plt.text(j, i, format(cm[i, j], 'd'),
                horizontalalignment="center",
                color="white" if cm[i, j] > thresh else "black")

    plt.ylabel('True label')
    plt.xlabel('Predicted label')
    plt.tight_layout()
    plt.show()

# Assuming y_test and y_pred are defined from previous steps
cm = confusion_matrix(y_test, y_pred)
print('Confusion matrix, without normalization')
print(cm)
plot_confusion_matrix(cm)

cm_normalized = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
print('Normalized confusion matrix')
print(cm_normalized)
plot_confusion_matrix(cm_normalized, title='Normalized confusion matrix')
```

**Algorithm Details :**

1. **Objective:**

   - The purpose of the code is to train a Support Vector Machine (SVM) classifier using grid search for hyperparameter tuning, assess its performance, and make predictions on a test set.

2. **Key Components:**

   - **Scikit-learn Library:**

     - **SVC:** Implements the Support Vector Classifier.

     - **StratifiedKFold**: Provides indices to split data into training and testing sets while preserving class sample percentages**.**

---

- **GridSearchCV:** Performs hyperparameter tuning by searching through specified parameter values using cross-validation.

- **scale:** Normalisses features by removing the mean and scaling to unit variance.

- **Inputs:**

  - **X_train:** Feature matrix for training data.

  - **y_train:** Target labels for training data.

  - **X_test:** Feature matrix for test data.

3. **Process:**

- Scale the features of both training and testing datasets.

- Define a range of values for hyperparameters C and gamma.

- Use StratifiedKFold for cross-validation.

- Initialize the SVM classifier and perform grid search to find the best hyperparameters.

- Fit the best model on the training data and print the best estimator.

- Estimate the model's performance using cross-validation and print the results.

- Plot learning curves to visualize model performance.

- Predict the target labels for the test data.

**Source Code Details :**

```
import numpy as np
from sklearn.model_selection import StratifiedKFold, GridSearchCV
from sklearn.svm import SVC
from sklearn import preprocessing
from sklearn.model_selection import cross_val_score
import matplotlib.pyplot as plt

def plot_learning_curve(clf, title, X, y, cv=5):
    """Plots the learning curve for the given classifier."""
    # Implementation of learning curve plotting goes here
    pass # Replace with actual implementation

def train(X_train, y_train, X_test):
    """ Trains and predicts dataset with a SVM classifier """

    # Step 1: Scaling features
    X_train = preprocessing.scale(X_train)
    X_test = preprocessing.scale(X_test)

    # Step 2: Define hyperparameter grid
```

```
Cs = 10.0 ** np.arange(-2, 3, .5)
gammas = 10.0 ** np.arange(-2, 3, .5)
param = [{'gamma': gammas, 'C': Cs}]

# Step 3: Setup cross-validation
cvk = StratifiedKFold(n_splits=5)

# Step 4: Initialize SVM classifier and GridSearchCV
classifier = SVC()
clf = GridSearchCV(classifier, param_grid=param, cv=cvk)

# Step 5: Fit the model on training data
clf.fit(X_train, y_train)
print("The best classifier is: ", clf.best_estimator_)

# Step 6: Fit the best estimator again
clf.best_estimator_.fit(X_train, y_train)

# Step 7: Estimate score using cross-validation
scores = cross_val_score(clf.best_estimator_, X_train, y_train, cv=5)
print(scores)
print('Estimated score: %0.5f (+/- %0.5f)' % (scores.mean(), scores.std() / 2))

# Step 8: Plot learning curves
title = 'Learning Curves (SVM, rbf kernel, $\gamma=%.6f$)' % clf.best_estimator_.gamma
plot_learning_curve(clf.best_estimator_, title, X_train, y_train, cv=5)
plt.show()

# Step 9: Predict class
y_pred = clf.best_estimator_.predict(X_test)

return y_test, y_pred  # Note: y_test should be defined in your context
```

**Algorithm Details :**

1. **Objective:**

   - The purpose of this code is to train a Support Vector Machine (SVM) classifier on the training dataset, predict the labels for the test dataset, and evaluate the model's performance using classification accuracy.

2. **Key Components:**

   - **Scikit-learn Library:**

     - **accuracy_score:** Computes the accuracy of the predictions.

   - **Inputs:**

     - **X_train:** Feature matrix for training data.

     - **y_train:** Target labels for training data.

     - **X_test:** Feature matrix for test data.

- **y_test:** True labels for the test data (should be defined in the context).

3. **Process:**

   - Call the train function to train the SVM classifier and obtain predictions for the test set.

   - Print the classification accuracy by comparing the true labels (y_test) with the predicted labels (y_pred).

**Source Code Details :**

```
from sklearn import preprocessing
from sklearn.metrics import accuracy_score

print("Training datasets ........ \n")
# Train the SVM classifier and get predictions
y_test, y_pred = train(X_train, y_train, X_test)

# Calculate and print the classification accuracy
print('Classification Accuracy on Test dataset: ', accuracy_score(y_test, y_pred))
```

# APPENDIX-B

# SCREENSHOTS

# APPENDIX-C

# ENCLOSURES

## Plagiarism Report

# *% detected as AI

AI detection includes the possibility of false positives. Although some text in this submission is likely AI generated, scores below the 20% threshold are not surfaced because they have a higher likelihood of false positives.

**Caution: Review required.**

It is essential to understand the limitations of AI detection before making decisions about a student's work. We encourage you to learn more about Turnitin's AI detection capabilities before using the tool.

**Disclaimer**

Our AI writing assessment is designed to help educators identify text that might be prepared by a generative AI tool. Our AI writing assessment may not always be accurate (it may misidentify writing that is likely AI generated as AI generated and AI paraphrased or likely AI generated and AI paraphrased writing as only AI generated) so it should not be used as the sole basis for adverse actions against a student. It takes further scrutiny and human judgment in conjunction with an organization's application of its specific academic policies to determine whether any academic misconduct has occurred.

## Frequently Asked Questions

**How should I interpret Turnitin's AI writing percentage and false positives?**
The percentage shown in the AI writing report is the amount of qualifying text within the submission that Turnitin's AI writing detection model determines was either likely AI-generated text from a large-language model or likely AI-generated text that was likely revised using an AI-paraphrase tool or word spinner.

False positives (incorrectly flagging human-written text as AI-generated) are a possibility in AI models.

AI detection scores under 20%, which we do not surface in new reports, have a higher likelihood of false positives. To reduce the likelihood of misinterpretation, no score or highlights are attributed and are indicated with an asterisk in the report (*%).

The AI writing percentage should not be the sole basis to determine whether misconduct has occurred. The reviewer/instructor should use the percentage as a means to start a formative conversation with their student and/or use it to examine the submitted assignment in accordance with their school's policies.

**What does 'qualifying text' mean?**
Our model only processes qualifying text in the form of long-form writing. Long-form writing means individual sentences contained in paragraphs that make up a longer piece of written work, such as an essay, a dissertation, or an article, etc. Qualifying text that has been determined to be likely AI-generated will be highlighted in cyan in the submission, and likely AI-generated and then likely AI-paraphrased will be highlighted purple.

Non-qualifying text, such as bullet points, annotated bibliographies, etc., will not be processed and can create disparity between the submission highlights and the percentage shown.

**Match Groups**

41 Not Cited or Quoted 12%
Matches with neither in-text citation nor quotation marks

1 Missing Quotations 0%
Matches that are still very similar to source material

0 Missing Citation 0%
Matches that have quotation marks, but no in-text citation

0 Cited and Quoted 0%
Matches with in-text citation present, but no quotation marks

**Top Sources**

5% Internet sources
8% Publications
7% Submitted works (Student Papers)

**Top Sources**

The sources with the highest number of matches within the submission. Overlapping sources will not be displayed.

| | | |
|---|---|---|
| 1 | Publication | |
| R. N. V. Jagan Mohan, B. H. V. S. Rama Krishnam Raju, V. Chandra Sekhar, T. V. K. P... | | 2% |
| 2 | Internet | |
| www.frontiersin.org | | <1% |
| 3 | Publication | |
| Hogan, Jacqueline S.. "Democratizing the Early Identification of Alzheimer's Disea... | | <1% |
| 4 | Publication | |
| Debasis Chaudhuri, Jan Harm C Pretorius, Debashis Das, Sauvik Bal. "Internationa... | | <1% |
| 5 | Internet | |
| ijircce.com | | <1% |
| 6 | Submitted works | |
| University of Wales, Lampeter on 2025-03-10 | | <1% |
| 7 | Internet | |

# 12% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

**Filtered from the Report**

- Bibliography
- Quoted Text

**Match Groups**

41 Not Cited or Quoted 12%
Matches with neither in-text citation nor quotation marks

1 Missing Quotations 0%
Matches that are still very similar to source material

0 Missing Citation 0%
Matches that have quotation marks, but no in-text citation

0 Cited and Quoted 0%
Matches with in-text citation present, but no quotation marks

**Top Sources**

5% Internet sources
8% Publications
7% Submitted works (Student Papers)

**Integrity Flags**

0 Integrity Flags for Review

No suspicious text manipulations found.

Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A Flag is not necessarily an indicator of a problem. However, we'd recommend you focus your attention there for further review.

8  Submitted works
Georgia Institute of Technology Main Campus on 2025-02-10                    <1%

9  Submitted works
University of Wales Institute, Cardiff on 2023-09-04                         <1%

10  Internet
www.al-kindipublisher.com                                                   <1%

11  Publication
"New Trends in Computational Vision and Bio-inspired Computing", Springer Scie...    <1%

12  Publication
A. S. Abdull Sukor, A. Zakaria, N. Abdul Rahim, L. M. Kamarudin, H. Nishizaki. "Abn...    <1%

13  Submitted works
University of Hertfordshire on 2024-11-26                                    <1%

14  Internet
revistaie.ase.ro                                                            <1%

15  Submitted works
Jyväskyla University on 2019-05-29                                          <1%

16  Publication

17  Publication
Faruq Abdul Hakim, Tio Dharmawan, Muhamad Arief Hidayat. "Gender classificati...    <1%

18  Internet
www.coursehero.com                                                          <1%

19  Internet
www.mdpi.com                                                                <1%

20  Publication
"EAI International Conference on Big Data Innovation for Sustainable Cognitive C...    <1%

21  Submitted works
Coventry University on 2023-10-17                                            <1%

22  Publication
Li Zhang, Haixin Ai, Wen Chen, Zimo Yin, Huan Hu, Junfeng Zhu, Jian Zhao, Qi Zha...    <1%

23  Publication
Vilda Purutçuoğlu, Gerhard Wilhelm Weber, Hajar Farnoudkia. "Operations Resea...    <1%

24  Publication
Yue, Liu. "Data-Driven Analysis on the Contact Resonance Frequency in the Strain...    <1%

25  Submitted works
Brunel University on 2024-05-29                                             <1%

# SUSTAINABLE DEVELOPMENT GOALS



1.  Goal 16: Peace, Justice and Strong Institutions

    -   Relevance to Fake Social Media Detection and Reporting:

        -   This project seeks to identify and combat misinformation and disinformation on social media platforms, which can undermine the integrity of democratic institutions and social cohesion.

        -   By empowering users to detect and report fake content, the project supports the development of more transparent and accountable digital spaces.

    -   Impact:

        -   Strengthened public trust in information sources and decision-making processes.

2. Goal 9: Industry, Innovation and Infrastructure

- Relevance to Fake Social Media Detection and Reporting:

  - The project involves the development of innovative technologies and tools for the automated detection of fake social media content, leveraging advances in areas like natural language processing and machine learning.

  - The infrastructure and systems built to support the project can serve as a model for the broader industry to enhance the integrity of online information ecosystems.

- Impact:

  - Accelerated innovation in the field of digital content verification and fact-checking.

  - Improved technological capabilities to combat the spread of

3. Goal 4: Quality Education

- Relevance to Fake Social Media Detection and Reporting:

  - The project can incorporate educational components to empower users, particularly younger generations, with the skills and critical thinking abilities to identify and navigate fake content online.

  - By raising awareness and promoting media literacy, the project can contribute to the development of a more informed and discerning digital citizenry.

- Impact:

  - Improved digital literacy and critical thinking skills among social media users.

  - Improved ability to make informed decisions and resist the influence of false or misleading information.

4.  Goal 10: Reduced Inequalities

- Relevance to Fake Social Media Detection and Reporting:

    - The project can ensure that the tools and resources developed are accessible and inclusive, catering to diverse communities and addressing the needs of marginalized groups who may be disproportionately targeted by online misinformation campaigns.
    - By promoting equitable access to accurate information, the project can help reduce social and digital divides.

- Impact:

    - Increased digital inclusion and empowerment of underserved communities.
    - Reduced vulnerability to the negative impacts of fake social media content across different demographics.

5.  Goal 17: Partnerships for the Goals

- Relevance to Fake Social Media Detection and Reporting:

    - The project may require collaboration and partnerships with various stakeholders, such as social media platforms, fact-checking organizations, academic institutions, and civil society groups, to leverage their expertise and resources effectively.
    - By fostering these multi-stakeholder partnerships, the project can amplify its reach and impact in the fight against online misinformation.

- Impact:

    - Enhanced global cooperation and coordination to tackle the challenges of fake social media content.
    - Increased synergies and shared knowledge across different sectors and disciplines working on this issue