



"Unveiling the Truth: Exploring AI Solutions to Identify Real vs Synthetic Images"

¹Mr. Gorakshan Bapurao Ingole, ²Mr. Animesh Sanjay Timande, ³Mr. Rutvik Vilasrao Behare.

⁴Mr. Purvesh Purushottam Dabhade, ⁵Prof. Snehal V. Raut,

^{1,2,3,4} Student, ⁵ Prof

^{1,2,3,4} Student, Dr. Rajendra Gode Institute of Technology and Research, Amravati, IN,

⁵ Guide, Prof Dr. Rajendra Gode Institute of Technology and Research, Amravati, IN

Abstract : In an era where synthetic media and deepfakes are becoming increasingly sophisticated, ensuring the authenticity of digital content has emerged as a critical challenge. This project explores an innovative AI-driven solution to distinguish real images from synthetic ones by utilizing the capabilities of Large Language Models (LLMs), specifically through Gemini technology. Unlike traditional methods that rely solely on pixel-level analysis, our approach leverages the multimodal understanding of LLMs to interpret visual data with contextual depth and semantic reasoning. The system is deployed via a web-based platform, supported by cloud services like Supabase for storage and authentication. Through this work, we aim to enhance digital trust by providing a scalable, intelligent tool for detecting manipulated visual content, contributing to the broader field of AI-based media forensics.

Index Terms - Real vs Fake Image Classification, Image Authenticity Verification, Synthetic Media Analysis, AI-Based Media Forensics, Digital Image Integrity.

I. INTRODUCTION

1.1 Motivation

In an era where visual content dominates social media, journalism, entertainment, and even scientific publications, trust in the authenticity of images is paramount. However, advances in generative models such as Generative Adversarial Networks (GANs) and diffusion models have blurred the line between real and synthetic content. For instance, AI-generated images of fictional political events or non-existent individuals can easily circulate as genuine, causing misinformation to spread rapidly. Traditional detection methods, while valuable, are increasingly insufficient against sophisticated fakes. Thus, there is an urgent need for adaptive, intelligent solutions capable of maintaining the integrity of digital media.

1.2 Problem Statement

The detection of synthetic images is no longer a simple task of identifying obvious artifacts. Modern AI-generated images can mimic real-world textures, lighting, and semantic coherence at levels previously thought impossible. Complicating matters further, new AI models are often released without accompanying detection methods, creating an arms race between creators and detectors. Our focus is to explore how AI-based methods, especially those incorporating LLMs like Gemini, can bridge this gap and offer more reliable, context-aware detection solutions.

II. BACKGROUND: SYNTHETIC IMAGE GENERATION

2.1 Early Developments

The concept of synthetic image generation dates back to early computer graphics, where artists manually created digital images. However, true AI-based generation began with models that could learn distributions of real images, such as autoencoders and variational autoencoders (VAEs).

2.2 Rise of GANs and Diffusion Models

The invention of GANs by Ian Goodfellow in 2014 marked a turning point. GANs operate via a game-theoretic framework where a generator tries to create convincing images while a discriminator tries to detect fakes. This led to realistic image synthesis capabilities in domains ranging from human faces to art.

More recently, diffusion models such as DALL-E 2 and Stable Diffusion use iterative noise addition and removal processes to generate stunningly realistic visuals. These models have achieved unprecedented levels of photorealism, challenging detection mechanisms.

III. LITERATURE REVIEW

3.1 Classical Methods

Early detection strategies focused on identifying visual inconsistencies. Researchers looked for:

Boundary Artifacts: Poorly blended edges.

Lighting Mismatches: Inconsistent shadows or reflections.

Color Anomalies: Slight variations imperceptible to humans.

Metadata Analysis: Altered or missing EXIF data.

Despite being computationally inexpensive, these methods lacked robustness, especially when fakes improved.

3.2 Deep Learning Approaches

CNNs revolutionized synthetic image detection. Models like XceptionNet and MesoNet were trained to classify real vs synthetic images by learning fine-grained texture patterns. Transfer learning techniques, using models pre-trained on ImageNet, boosted performance. However, deep CNNs struggled with generalization: a model trained on one type of fake (e.g., GAN-generated) often failed against others (e.g., diffusion-generated).

Key Techniques:

Pixel-level feature extraction.

Frequency domain analysis.

Patch-based learning.

3.3 Transformer-Based and LLM Methods

Transformers introduced self-attention mechanisms that model long-range dependencies within images. Vision Transformers (ViT) and hybrid CNN-Transformer models improved detection by capturing global context.

Recently, multimodal LLMs such as Gemini and GPT-4V have been applied. Instead of only pixel-wise analysis, these models:

Analyze semantic coherence (e.g., does the background match the foreground context?),

Cross-reference image metadata,

Provide explanations along with decisions ("This image appears synthetic because...").

Their ability to process images and text together enables reasoned detection, offering a significant leap beyond CNNs.

IV. DATASETS AND BENCHMARKS

4.1 CIFAKE

The CIFAKE dataset is widely used for training and evaluating synthetic image detection models, especially in the domain of face forgery. This dataset includes over 500,000 synthetic images generated by GAN-based models and real images from publicly available datasets. It serves as a benchmark for models trained to detect fake facial images, with a focus on the subtle texture inconsistencies introduced by GANs.

4.2 Deep-Fake Detection Challenge (DFDC)

The DFDC dataset, developed for the DeepFake Detection Challenge organized by Facebook AI, contains over 100,000 video frames from real and deepfake videos. The dataset was designed to reflect real-world conditions and includes diverse actors, lighting conditions, and scenarios. It provides a valuable resource for developing models capable of detecting deepfake videos as well as still images. This dataset emphasizes the importance of generalizing models across varied generative methods, as it includes both traditional GAN-based fakes and newer diffusion-based fakes.

4.3 Fake-Chain

Fake Chain is a relatively new dataset focused on deepfake detection in the blockchain context, aiming to address the verification of image provenance. This dataset contains synthetic and real media along with traceable metadata, such as timestamps and transaction histories. The goal of Fake Chain is to improve the transparency of synthetic media in the digital space by ensuring that images can be linked back to their origin and verified for authenticity.

4.4 Other Notable Datasets

Other datasets include CelebA (focused on celebrity faces), FFHQ (Flickr-Faces-HQ for high-quality face generation), and ImageNet (which is not specifically focused on synthetic images but is useful for pretraining general-purpose deep learning models). These datasets are used in a variety of applications, including image classification, object recognition, and style transfer.

V. SYSTEM ARCHITECTURE

5.1 System Architecture Overview

The system developed for identifying real versus synthetic images is based on a lightweight, scalable architecture that leverages modern web technologies and powerful AI services. Unlike traditional machine learning systems that require extensive datasets and model training, our approach utilizes Google's Gemini Large Language Model (LLM) via API integration to perform synthetic image detection in real-time. The system is divided into three main components: Frontend Interface, Backend Processing, and Cloud Database Storage.

5.2 Backend Processing – Leveraging Gemini LLM

The core intelligence of our system resides in the backend, where the Gemini LLM is accessed through secure API requests. When a user uploads an image, the backend server processes this input and sends a properly formatted request to the Gemini model. The LLM evaluates the image, drawing upon its extensive pretrained multimodal knowledge to determine whether the input is real or synthetic. The response received from Gemini is parsed and formatted appropriately to deliver a meaningful output to the user.

Key aspects of the backend:

API-driven architecture.
No need for manual training on CIFAKE, DFDC, or other datasets.
Real-time predictions and immediate feedback.

Enhanced flexibility to adapt to updates in Gemini's capabilities.

This approach greatly reduces computational overhead, minimizes latency, and allows for easy scalability without the need for powerful on-premise GPUs.

5.3 Frontend Development – User Interaction Layer

The frontend of the application is built using simple yet effective technologies: HTML, CSS, and JavaScript. Its primary functions include:

Allowing users to upload images directly from their devices.

Providing a clean, responsive user interface to interact with the system.

Displaying prediction results (Real or Synthetic) in an understandable manner.

Handling simple error checking (such as unsupported file formats or failed uploads).

User Experience (UX) has been prioritized to ensure that even non-technical users can easily interact with the system.

5.4 Database Storage – Supabase Integration

All relevant data is securely stored using Supabase, an open-source backend-as-a-service platform. The database handles:

Secure storage of uploaded images (if required).

Logging prediction results for analytics and future improvement.

Managing user sessions and authentication (if implemented).

Supabase's real-time capabilities and PostgreSQL foundation make it ideal for handling the relatively lightweight data flows associated with synthetic image detection applications.

5.5 Workflow of the System

The overall process flow of the system is as follows:

1. The user accesses the web interface and uploads an image.
2. The frontend sends the image to the backend server.
3. The backend processes the image and sends it to the Gemini LLM through API.
4. Gemini analyzes the image and returns a prediction (real or synthetic).
5. The backend receives the result, stores it (if needed) in Supabase, and sends the result back to the frontend.
6. The frontend displays the result to the user in a user-friendly format.

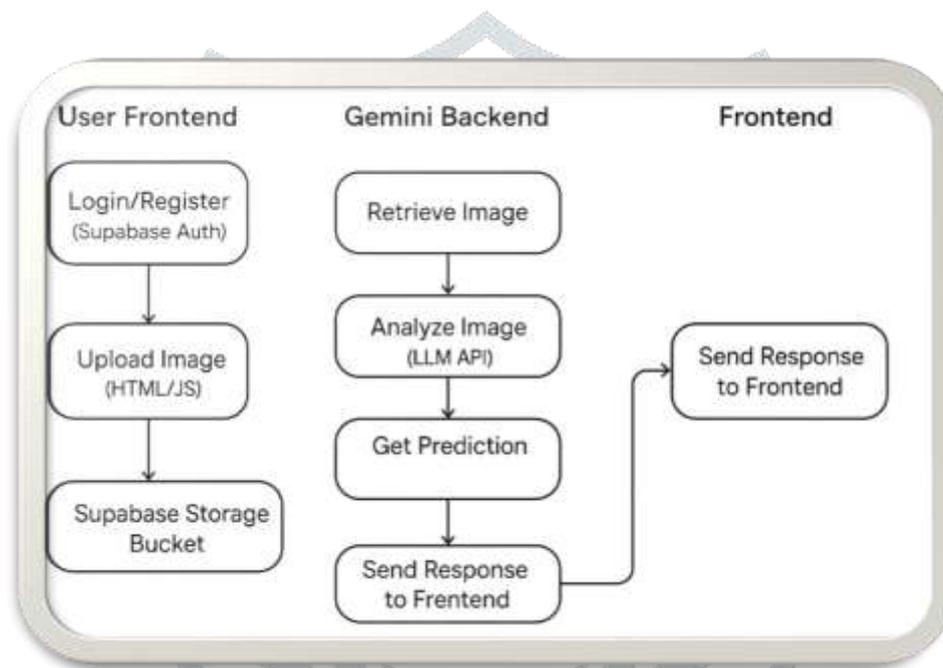


Fig 1. Flow Implementation

VI. TECHNICAL CHALLENGES

6.1 Transferability of Detection Models

One of the most significant technical challenges in synthetic image detection is transferability. Deep learning models trained on specific datasets often fail to generalize well to unseen synthetic content, especially when generated by different models. For instance, a model trained on detecting GAN-based images may not perform as well when detecting images generated by newer models like diffusion models. This creates a challenge in developing a universal detection system capable of identifying synthetic media across diverse platforms.

6.2 Data Imbalance

Many existing datasets suffer from data imbalance, where real images significantly outnumber synthetic images. This imbalance can lead to bias in model predictions, where the model may become overly conservative and favor predicting images as real. Techniques such as data augmentation and class balancing are often used to mitigate this issue, but it remains a persistent challenge in the field.

6.3 Adversarial Evasion

As detection models become more accurate, generative models evolve to evade detection. Adversarial attacks involve subtly modifying synthetic images to make them more challenging for detection models. These attacks can take the form of pixel-level perturbations or subtle adjustments to image features that are imperceptible to humans but confuse machine learning models. Developing detection systems that are robust to adversarial attacks is an ongoing challenge.

6.4 Model Interpretability

While deep learning models, especially CNNs and transformers, have achieved impressive performance in detecting synthetic images, they often operate as "black boxes," making it difficult to understand why a model makes a certain decision. Model interpretability is critical for building trust in detection systems, especially when they are deployed in sensitive areas like legal proceedings or public safety. Advances in explainable AI (XAI) are necessary to enhance the transparency of detection models and provide human-understandable explanations for their decisions.

6.5 Real-time Detection

As synthetic images and videos become more prevalent in real-time applications (e.g., live-streaming, video conferences), there is a growing demand for real-time detection systems. These systems must be capable of processing high-resolution images and videos with low latency while maintaining high detection accuracy. Meeting this requirement will require innovations in model efficiency, such as model pruning and quantization, to reduce the computational overhead of detection systems.

6.6 Generalization Across Diverse Generators

One of the major challenges in synthetic image detection is the generalization of detection models. Models trained on one type of generative model (e.g., GANs) may not perform well on another (e.g., diffusion models), as each model introduces different artifacts or lacks them entirely. For instance, GANs often exhibit pixel-level inconsistencies, while diffusion models generate smooth, photorealistic images that are harder to detect. Therefore, detectors need to handle a variety of generative methods.

6.7 High Computational Complexity

Deep learning-based methods for synthetic image detection, particularly CNNs and transformer models, often require significant computational power. Large datasets, especially those needed for training robust models, demand extensive computational resources. Training times can span days or even weeks, making it difficult to deploy these systems in real-time applications, such as social media platforms or digital forensics.

6.8 Lack of Comprehensive Datasets

While datasets like CIFAKE, DFDC, and FakeChain are commonly used for evaluating synthetic image detection methods, there is still a lack of diverse, real-world datasets. Most existing datasets focus on a specific type of synthetic content (e.g., faces, videos) and may not cover the broad spectrum of generated media (e.g., landscapes, objects). This limits the ability of detection models to generalize across different use cases.

6.9 Adversarial Attacks and Evasion

As detection methods improve, so do attempts to evade detection. Adversarial attacks, such as perturbations on synthetic images, can render deep learning models ineffective. These attacks involve subtly modifying synthetic images to confuse the detection system without changing the image's perceptual quality. A key challenge is developing detection techniques resilient to such attacks.

VII. EMERGING TREND AND FUTURE DIRECTION

7.1 Multimodal Detection Models

The integration of multimodal models is an exciting emerging trend in synthetic image detection. While traditional models focus only on image data, multimodal models, especially those incorporating LLMs like Gemini, can leverage both textual and visual information. For example, a model might analyze the image while cross-referencing its contextual metadata (e.g., timestamps, geolocation) to assess authenticity. This enables more reliable detections by combining evidence from multiple sources.

7.2 Watermarking and AI-Generated Metadata

As a response to increasing concerns over synthetic media, researchers are investigating techniques like watermarking and embedding traceable information into AI-generated images. This could involve embedding invisible digital watermarks or generating metadata that is inherently tied to the image's creation. These markers would allow easier verification of authenticity, as detectors could extract and verify these watermarks automatically.

7.3 Real-Time Detection Systems

With the increasing use of synthetic images in live media (e.g., livestreams, video calls), there is a growing need for real-time detection systems. This presents several challenges, including the need for extremely low-latency models that can process high-resolution images or videos in real-time without compromising detection accuracy. Innovations in model efficiency, such as pruning and quantization, may help achieve these goals.

7.4 Cross-Domain Application

The detection techniques developed for images can be extended to other domains such as audio and video. Fake audio and deepfake videos are becoming more prevalent, necessitating multi-modal detection systems that can analyze video, audio, and visual cues simultaneously. Future research could look into creating cross-domain detection models that can perform multimodal analysis to identify synthetic media across formats.

VIII. ETHICAL, LIGAL AND SOCIETAL IMPLICATION

8.1 Bias and Fairness in Detection Models

An important ethical consideration when developing synthetic image detection methods is bias. If the training data for these models are not sufficiently diverse, detection systems may fail to identify synthetic content in certain demographics or cultural contexts. For example, models may struggle with detecting fake images of people from underrepresented racial groups or regions. Ensuring fairness in synthetic image detection is crucial, and addressing this bias will require diverse training datasets and transparent evaluation metrics.

8.2 Privacy and Surveillance

With the ability to detect synthetic media comes the potential for surveillance. AI-based detection models can be deployed to monitor digital content at scale, raising concerns about privacy. Governments, corporations, and malicious actors could exploit these technologies for intrusive monitoring. Ethical considerations around consent, data ownership, and the use of detection systems must be addressed, especially when dealing with personal content or biometric data.

8.3 Legal Implications: Copyright and Misinformation

Legally, the widespread dissemination of synthetic images raises important questions about copyright infringement. Who owns a generated image, especially when it is indistinguishable from real content? Additionally, the proliferation of AI-generated images used in deepfakes and misinformation campaigns challenges existing legal frameworks. Governments and regulatory bodies are grappling with how to legislate these issues, such as who should be held accountable for the creation and distribution of synthetic media.

8.4 Public Trust and Transparency

As AI detection systems evolve, it is critical to maintain public trust. For AI models to be adopted and trusted by society, they must be transparent in their decision-making. The “black box” nature of many AI systems, particularly deep learning models, often makes it difficult for users to understand how a model arrived at a decision. Efforts to increase transparency in these systems, such as the development of explainable AI (XAI), will be essential in securing trust and ensuring that these systems are used responsibly.

IX. FUTURE DIRECTION

Advancements in AI Models:

Future research will likely focus on developing more advanced detection models, such as transformer-based architectures, which have shown promising results in other areas of computer vision.

Cross-Domain Detection:

Researchers are exploring cross-domain approaches that allow models trained on one type of synthetic media (e.g., deepfake images) to be applied to other domains, such as synthetic videos or voices.

Real-Time Detection:

Real-time synthetic media detection remains a key goal, particularly in streaming platforms and live broadcasts. Developing lightweight models capable of detecting deepfakes in real-time without sacrificing accuracy is a significant area of ongoing research.

Collaboration with Other Fields:

Collaboration between AI researchers, policymakers, and law enforcement will be crucial in addressing the societal impacts of synthetic media. Developing regulations and standards for synthetic media detection could help create a safer digital environment.

X. CONCLUSION

Synthetic image detection is a critical area of research with significant societal implications. As generative AI technologies continue to improve, the ability to accurately distinguish between real and fake images will be paramount in preventing misinformation, fraud, and digital manipulation. The review of current methods highlights the importance of multimodal detection systems, particularly those incorporating LLMs, as promising solutions for the future. While technical challenges remain, the growing attention to ethical and legal issues underscores the need for responsible development and deployment of these technologies. Looking forward, a combination of advanced AI models, robust datasets, and international collaboration will be essential to navigating the complexities of synthetic media detection.

XI. REFERENCES

- [1]. Zhang, X., & Dong, X. (2020). Deep learning for detecting image forgery: A survey. *Pattern Recognition*, 105, 107294. <https://doi.org/10.1016/j.patcog.2020.107294>
- [2]. Dolhansky, B., et al. (2020). The deepfake detection challenge. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 1-9. <https://doi.org/10.1109/CVPR42600.2020.00010>
- [3]. Yli-Huumo, J., Ko, D., Choi, S., & Park, S. (2016). A survey of blockchain: A comprehensive review of applications, challenges, and opportunities. *Proceedings of the 2016 IEEE 13th International Conference on Embedded Software and Systems (ICCESS)*, 18-24. <https://doi.org/10.1109/ICCESS.2016.17>
- [4]. Ali, M. M., & Vasilenko, D. (2018). Blockchain technology and its applications in image provenance. *International Journal of Computer Science and Network Security*, 18(10), 37-45. <https://doi.org/10.22937/IJCSNS.2018.18.10.37>
- [5]. Goodfellow, I. J., Shlens, J., & Szegedy, C. (2015). Explaining and harnessing adversarial examples. *Proceedings of the International Conference on Machine Learning (ICML)*, 1-9. <https://arxiv.org/abs/1412.6572>
- [6]. Carlini, N., & Wagner, D. (2017). Towards evaluating the robustness of neural networks. *Proceedings of the IEEE Symposium on Security and Privacy (SP)*, 39-57. <https://doi.org/10.1109/SP.2017.49>

