

DETECTION OF AI-GENERATED IMAGES VS REAL IMAGES USING EFFICIENTNET-B0

Seerapu Yamini¹, Vantakula Lakshmi²

¹ MCA Student, Gayatri Vidya Parishad College of Engineering(A), Visakhapatnam – 530048, Andhra Pradesh, India

² Assistant Professor, Department of Computer Applications, Gayatri Vidya Parishad College of Engineering(A), Visakhapatnam – 530048, Andhra Pradesh, India

Abstract - Artificial intelligence-generated images are rapidly increasing due to advancements in deep learning models and image synthesis technologies. Large volumes of synthetic images are produced every year using techniques such as Generative Adversarial Networks (GANs) and diffusion models. However, most of these images are highly realistic and difficult for users to distinguish from real photographs without advanced analytical methods. This research presents the development of an AI-generated image detection system using EfficientNet-B0 and BLIP to identify whether an image is authentic or synthetically created by artificial intelligence. The proposed system integrates EfficientNet-B0 with transfer learning to analyze hidden visual patterns, structural inconsistencies, and image artifacts for accurate classification. The system also incorporates the BLIP (Bootstrapping Language-Image Pretraining) model to generate descriptive captions for uploaded images, improving interpretability and user understanding. The system converts image data into classification results and textual descriptions through image preprocessing, feature extraction, and caption generation. By transforming complex deep learning processes into a simple and accessible solution, the proposed system simplifies image authenticity verification and helps reduce the spread of misleading AI-generated content across digital media platforms.

Key Words: AI-Generated Images, EfficientNet-B0, Deep Learning, Image Classification, BLIP, Image Authenticity Detection

I. INTRODUCTION

Artificial intelligence plays a vital role in transforming digital content creation, image processing, and media generation across various domains. Modern AI technologies such as Generative Adversarial Networks (GANs), diffusion models, and deep learning architectures are capable of producing highly realistic synthetic images that closely resemble real-world photographs. These advanced technologies are widely used in fields including entertainment, social media, advertising, virtual reality, and digital design.

Organizations and online platforms continuously generate and share large volumes of digital images through social media applications, websites, and AI-based content creation systems. These images include both authentic photographs and AI-generated synthetic content created using advanced image generation techniques. Although such technologies provide significant advantages in creativity and automation, they also create serious challenges related to misinformation, deepfakes, image manipulation, and digital authenticity verification.

Traditional image verification methods mainly rely on manual inspection and basic image analysis techniques. This makes it difficult for users to accurately distinguish between real and AI-generated images, especially when synthetic images contain highly realistic textures, lighting, and visual patterns. With the increasing availability of AI-generated content, there is a growing need for intelligent

systems that can automatically analyze images and provide accurate authenticity verification in a reliable and efficient manner.

Deep learning techniques provide an effective solution for this challenge. Convolutional Neural Networks (CNNs) convert complex image data into meaningful feature representations that allow systems to identify hidden artifacts, visual inconsistencies, and structural patterns more effectively. Transfer learning techniques further improve performance by utilizing pre-trained models, reducing computational complexity and training time while maintaining high classification accuracy.

This research proposes an AI-generated image detection system using EfficientNet-B0 and BLIP that integrates deep learning-based classification and image caption generation. The system allows users to analyze uploaded images, identify whether they are real or AI-generated, and generate descriptive captions through interactive processing. By transforming complex image analysis into meaningful classification and textual outputs, the system improves image authenticity verification and supports secure and reliable digital media analysis.

II. RELATED WORK

A. AI-Generated Image Detection and Deep Learning Research

Deep learning techniques have become essential for analyzing and identifying AI-generated images.

Goodfellow et al. introduced Generative Adversarial Networks (GANs), a framework capable of generating highly realistic synthetic images through adversarial training between generator and discriminator networks [1]. Their work demonstrated the rapid advancement of image generation technologies and highlighted the growing challenge of distinguishing synthetic images from real photographs. Rössler et al. developed FaceForensics++, a large-scale dataset for detecting manipulated and AI-generated images [2]. Their research showed that deep learning models can identify hidden artifacts and inconsistencies present in synthetic images. The study

emphasized the importance of benchmark datasets in training reliable image authenticity detection systems. Li et al. further explored AI-generated image analysis using pixel-level artifact detection methods, demonstrating that deep learning techniques can effectively identify subtle structural inconsistencies in generated content [4].

B. Deep Learning Models for Image Classification Several studies have focused on improving image classification performance using advanced deep learning architectures. Tan and Le proposed EfficientNet, a family of convolutional neural networks that achieves high accuracy while reducing computational complexity through compound scaling methods [3]. EfficientNet-B0 provides an optimal balance between efficiency and performance, making it suitable for image authenticity detection systems. He et al. introduced Residual Networks (ResNet), which improved deep learning performance by solving the vanishing gradient problem in deep neural networks [5]. Their work significantly enhanced feature extraction and classification capabilities in computer vision tasks.

Dosovitskiy et al. proposed Vision Transformers (ViT), which apply transformer-based architectures to image recognition problems and achieve competitive performance in image classification tasks [10].

C. Vision-Language Models and Image Caption Generation

Vision-language models combine image understanding with natural language processing to improve interpretability and content understanding. Radford et al. introduced CLIP, a model that learns visual representations using natural language supervision [6]. Their research demonstrated the effectiveness of combining visual and textual information for image understanding tasks.

Li et al. proposed BLIP (Bootstrapping Language Image Pretraining), a vision-language model capable of generating descriptive captions for images [7]. The model improves image interpretation by connecting visual content with

natural language descriptions. This approach enhances user understanding and supports explainable AI systems. These studies demonstrate that integrating image classification with caption generation significantly improves interpretability and user interaction compared to traditional image analysis systems.

D. Limitations of Existing Systems and Research Gap
Despite significant advancements in AI-generated image detection, existing systems still exhibit several limitations. Many detection models primarily focus only on classification and do not provide descriptive explanations of image content [2], [4]. Additionally, several systems require large computational resources and struggle to maintain efficiency in real-time applications.

Furthermore, many traditional image verification approaches rely heavily on manual analysis or static detection techniques, making them less effective against modern AI-generated images [1], [6]. This limits the ability of users to accurately verify image authenticity and understand image content simultaneously.

To address these limitations, this research proposes an AI-generated image detection system using EfficientNet-B0 and BLIP, which integrates efficient image classification and descriptive caption generation. The proposed system improves accessibility, interpretability, and reliability by providing both authenticity verification and image description through an interactive and intelligent framework.

III. DATASET DESCRIPTION

The dataset used in this research is collected from publicly available image datasets and AI-generated image repositories used for computer vision and deep learning research. The dataset contains both real images and synthetic images generated using advanced artificial intelligence models such as Generative Adversarial

Networks (GANs) and diffusion-based image generation techniques.

The dataset consists of multiple types of image information, which are described below:

- **Real Image Data:** Contains authentic images collected from publicly available datasets. These images are used as reference data for training and evaluating the classification model.
 - **AI-Generated Image Data:** Includes synthetic images generated using AI-based image generation models. This data is used to identify hidden artifacts and inconsistencies present in generated images.
 - **Caption Data:** Contains descriptive textual information related to images. This dataset is used by the BLIP model to generate meaningful captions for uploaded images.
 - **Classification Labels:** Includes labels such as "Real" and "AI-Generated" assigned to images for supervised learning and accurate classification.
- The image datasets are available in formats such as JPG, JPEG, and PNG, which are directly suitable for image processing and deep learning tasks. Therefore, no additional file conversion is required before training and analysis.

Before classification and caption generation, the datasets undergo preprocessing steps such as:

- Removing corrupted and duplicate images
- Resizing images into fixed dimensions
- Normalizing pixel values

- Organizing images into structured format These steps ensure that the data is clean, balanced, and ready for accurate image classification and caption generation.

IV. METHODOLOGY

The proposed AI-generated image detection system is developed using a structured methodology consisting of three major phases: Dataset Collection and Preprocessing, Model Development and Training, and Image Classification with Caption Generation.

Each phase of the methodology is explained in detail below.

The overall workflow of the proposed system is illustrated in Fig. 1.

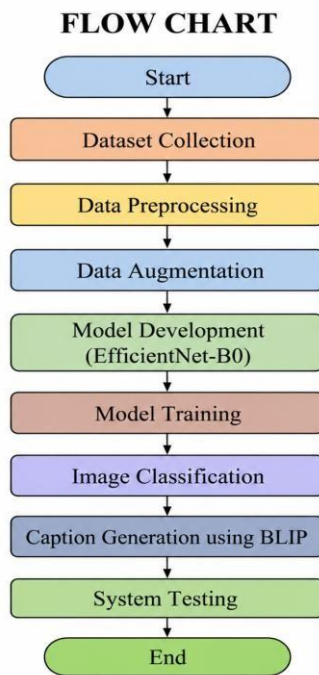


Fig. 1: Flowchart of AI-Generated Image Detection and Caption Generation System

[1] PHASE 1: DATA COLLECTION AND PREPROCESSING

Dataset Collection:

The datasets used in this research were collected from publicly available image datasets and AI-generated image repositories used for deep learning and computer vision research. These datasets contain both authentic images and synthetic images generated using advanced artificial intelligence techniques such as GANs and diffusion models.

The collected datasets include:

- Real image datasets
- AI-generated image datasets
- Image caption datasets
- Classification label information

The datasets contain thousands of images belonging to both real and AI-generated categories for accurate model training and evaluation.

Data Preprocessing:

Raw image datasets often contain corrupted files, duplicate images, inconsistent image sizes, and unnecessary noise. Therefore, preprocessing is performed to clean and prepare the datasets for classification and caption generation.

The preprocessing steps include:

- Resizing images into fixed dimensions
- Removing duplicate and corrupted images
- Normalizing pixel intensity values
- Organizing datasets into training and testing categories
- Applying image augmentation techniques such as rotation and flipping

Python libraries such as Pandas, NumPy, TensorFlow, and OpenCV are used to perform image preprocessing and transformation. After preprocessing, the datasets are organized into structured directories for efficient model training and image classification.

[2] PHASE 2: MODEL DEVELOPMENT AND TRAINING

In this phase, the processed image datasets are integrated into the deep learning framework for model development and training. The EfficientNet-B0 model and BLIP model are configured to perform image classification and caption generation tasks.

Model Development:

The pre-processed datasets containing real and AI-generated images are loaded into the training environment. Transfer learning is applied using the EfficientNet-B0 architecture with pre-trained weights to improve classification performance and reduce training time.

Image datasets are divided into training and testing sets to ensure proper evaluation and validation of the model.

Model Configuration:

Several model parameters are configured to support accurate image classification. These include:

- Image input dimensions
- Batch size and learning rate
- Number of training epochs
- Optimizer and loss function
- Transfer learning configuration

The BLIP (Bootstrapping Language-Image Pretraining) model is also integrated to generate descriptive captions for uploaded images.

Training Techniques:

Various deep learning techniques are used to improve model performance effectively. These include:

- Feature extraction using EfficientNet-B0
- Transfer learning for efficient training
- Image classification for authenticity detection
- Caption generation using BLIP
- Performance evaluation using accuracy, precision, and recall metrics

[3] PHASE 3: IMAGE CLASSIFICATION AND CAPTION GENERATION

The final phase involves developing an interactive system that allows users to upload images and obtain classification results along with descriptive captions dynamically.

System Implementation:

Several modules are developed to perform different image analysis tasks, including:

- Image Upload Module
- Image Preprocessing Module
- EfficientNet-B0 Classification Module
- BLIP Caption Generation Module
- Prediction Result Display Module
- Performance Evaluation Module
- User Interaction Interface

Interactive Processing:

Interactive processing allows users to upload and analyze images easily. Users can provide input images through the system interface and receive instant prediction results.

The system performs operations such as:

- Image preprocessing and normalization

- Feature extraction using EfficientNet-B0
- Real vs AI-generated image classification
- Caption generation using BLIP
- Displaying prediction results and image descriptions

When an image is uploaded, the system dynamically processes the image and updates the output results automatically.

Caption Generation and Output:

Image caption generation is implemented using the BLIP vision-language model. The system analyzes uploaded images and generates descriptive textual captions based on image content. The generated captions help users understand image context along with authenticity verification. This process improves interpretability and allows users to observe meaningful descriptions for both real and AI-generated images.

V. RESULTS AND SYSTEM ANALYSIS

The developed AI-generated image detection system successfully transforms complex deep learning operations into accurate image classification and caption generation results.

The EfficientNet-B0 classification module accurately identifies whether uploaded images are real or AI-generated by analyzing hidden visual patterns and structural inconsistencies. The system provides fast and reliable prediction results through an interactive user interface. These results allow users to quickly verify image authenticity and understand whether the uploaded image is synthetically generated.

The BLIP caption generation module produces meaningful textual descriptions for uploaded images through vision-language processing techniques. These captions help users understand the content and context of images along with authenticity verification.

The system also includes several important functionalities, such as:

- Real image detection
- AI-generated image detection
- Image preprocessing and normalization
- Feature extraction using EfficientNet-B0
- Caption generation using BLIP
- Performance evaluation using accuracy, precision, and recall metrics

One of the most significant features of the system is the integration of image classification and caption generation within a single framework. The system not only identifies whether an image is real or AI-generated but also generates descriptive captions that improve interpretability and user understanding.

Overall, the proposed system provides an intelligent and interactive platform for detecting AI-generated images and generating image descriptions efficiently. The developed system simplifies image authenticity verification and supports secure analysis of digital media content.

VI. CONCLUSION

This research presented the development of an AI-generated image detection system using EfficientNet-B0 and BLIP for identifying real and AI-generated images. The system integrates deep learning-based image classification and caption generation techniques to provide a comprehensive image authenticity verification platform.

By transforming complex image data into classification results and descriptive captions, the system simplifies image analysis and improves accessibility for users. The system enables accurate authenticity detection, hidden artifact analysis, feature extraction, and automatic image caption generation through an interactive interface.

The results demonstrate that modern deep learning models such as EfficientNet-B0 and vision-language models like BLIP can significantly enhance the detection and interpretation of AI-generated images. The proposed system can support digital media verification, cybersecurity, image forensics, and educational research by providing reliable and intelligent image analysis capabilities.

Future work may include integrating real-time image verification, video deepfake detection, advanced transformer-based architectures, and multimodal AI techniques for improved image authenticity analysis and content understanding.

VII. REFERENCES

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative Adversarial Networks," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 27, pp. 2672–2680, 2014.
- [2] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner, "FaceForensics++: Learning to Detect Manipulated Facial Images," *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 1–11, 2019.
- [3] M. Tan and Q. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 6105–6114, 2019.
- [4] A. Radford et al., "Learning Transferable Visual Models From Natural Language Supervision," *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [5] Y. Wang et al., "CNN-generated images are surprisingly easy to spot... for now," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [6] A. Nichol and P. Dhariwal, "Improved Denoising Diffusion Probabilistic Models," *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [7] J. Li, D. Li, C. Xiong, and S. Hoi, "BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation," 2022.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- [9] A. Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," *International Conference on Learning Representations (ICLR)*, 2021.
- [10] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. Chen, "MobileNetV2: Inverted Residuals and Linear Bottlenecks," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4510–4520, 2018.