
Report on AI Image Detector

LAM, Tsz Kit
SID: 1155194085
The Chinese University of Hong Kong
1155194085@link.cuhk.edu.hk

CHOI Ho Yan
SID: 1155194468
The Chinese University of Hong Kong
1155194468@link.cuhk.edu.hk

Abstract

This paper targets the problem of high similarity between AI-generated images and human-crafted images, proposes an AI Image Detector using three image classification models, namely custom Convolutional Neural Network, ResNet and Vision Transformer, and analyzes their performance and effectiveness.

1 Introduction

Recently, AI-generated images adapted to the style of Ghibli, a Japanese animation studio well-known for its unique anime style behind *My Neighbor Totoro*, *Spirited Away* and other classical films, have sparked heated discussions about the role of AI in artwork. The latest GPT-4o model from OpenAI has incorporated this artistic style into its image generator, resulting in an 11% increase in global ChatGPT app downloads and a 5% rise in weekly active users [1].

The underlying problem behind this discussion is the decrease in disparity between human-crafted images and AI-generated images, leading to many potential threats of misleading and misuses of information. People are more sensitive to images when receiving information, and there is an average of 2 million images generated daily [2], highlighting the severe consequences of AI images. For example, forgery AI photos and information in celebrity reputations or political campaigns could frame public opinion, with more than 80% of Americans expressing concerns about AI misuse in the U.S. presidential election [3].

To mitigate the potential concerns about AI images, this project presents a binary classification model that classifies between AI-generated and human-designed images. We hope that this AI image detector can effectively recognize AI photos so that people can be more cautious about the image sources when receiving new information.

2 Problem Statements

Given a dataset divided into "AI-generated image" and "human-generated image" two classes, for each image in the dataset, the goal of our classifier is to predict which class it belongs to. We implement the classifier by designing a custom CNN architecture and fine-tuning an existing ViT architecture. In each case, the model is trained with labelled images and tested with unseen images. Each model performance is improved by means of the accuracy and loss for each epoch, and analyzed by its confusion matrix. We further discuss the comparison and improvements of these two models.

3 Data

3.1 Data Source

We would use a dataset named "AI vs. Human-Generated Images" published in Kaggle [4]. The image data are sampled from the Shutterstock platform across different categories and labelled either "human-created images" or "AI-generated images", in which the images are pairings of an authentic image and its equivalent generated image using AI generative models. We perform the model training and testing based on the provided dataset, assuming this data source is reliable and general.



Figure 1: An example of an authentic image paired with its AI-generated equivalent [4]

3.2 Data Extraction

The original dataset contains approximately 80,000 training data and 5,500 testing data, with CSV files indicating the class labels for each image filename. However, we found that the labeling for the test set is not clear and the folder structures are not well organized.

Therefore, we randomly extracted 25,000 training data and 5,000 testing data from the training dataset for each "AI-generated image" class and "human-generated image" class.

3.3 Data Preprocessing

After extracting the necessary dataset, the images are resized to 224x224 pixels to reduce training computations and unify the input size [5]. We apply ImageNet normalization to the pixel values. The ImageNet normalization involves subtracting the mean and dividing by the standard deviation of the ImageNet dataset RGB channels (mean = [0.485, 0.456, 0.406], std = [0.229, 0.224, 0.225]), which is a common practice and provide a better model training convergence [5].

We also perform data argumentation on training data. The random transformations, such as horizontal flips, rotations and color jitter, are applied to increase the dataset diversity and prevent overfitting. A sample of code snippet can be found in Appendix A.

4 Model Architectures

Image classification is a fundamental task in computer vision which assigns a specific category to an image based on its visual content. Early image classification methods flatten images into pixels and use traditional feedforward neural networks or heuristic algorithms like SIFT and SURF [6]. With a successful breakthrough in image classification accuracy by AlexNet using the deep layers, ReLU activation and dropout regularization, convolutional neural networks (CNNs) begin to be widely adopted in image classification [7]. Later with an introduction of the attention mechanism and Vision Transformer (ViT), the performance of classification is further improved for large-scale tasks [8].

Following the history and literature reviews, we build our AI image classifier by implementing a custom CNN and a transfer learning of ResNet and ViT architecture. We present the architecture

overviews, training and testing results in the following three subsections, and analyze the performance of three models in the next sections.

4.1 Custom Convolutional Neural Network

A Convolutional Neural Network is a special type of feedforward neural network designed to learn spatial features adaptively through convolutional filters or kernels [9]. Similar to many neural networks, neurons in CNN layers, however, are organized into three dimensions (height, width and depth) representing the spatial information. Therefore, CNN is suitable for performing image-related machine learning (color channels for depth).

A CNN is mainly composed of convolutional layers, pooling layers and fully-connected layers [9]. The convolutional layer is used to extract spatial features, such as edges and textures of an image, by applying learnable filters to the input; the pooling layer is responsible for reducing spatial dimensions to prevent overfitting while preserving important features; the fully-connected layer is used to combine features from previous layers for classification and output the final category class information of the image.

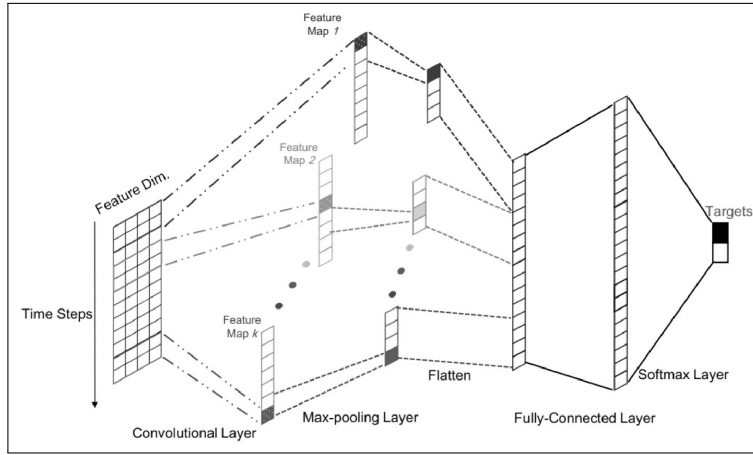


Figure 2: An example of CNN layer structure [10]

4.1.1 Architecture Design

Following the introduction, we design our CNN architecture for AI image classifier. See Figure 3.

In convolutional layers, the spatial features are extract from low level (edges, shapes) to high level by layers. Each output size could be calculated by $(\lfloor \frac{W-K+2P}{S} \rfloor + 1) \times (\lfloor \frac{H-K+2P}{S} \rfloor + 1)$, where W is input width, H is input width, K is kernel size, P is padding, S is stride.

After each convolutional layer, we apply Batch Normalization and ReLU activation function. Batch Normalization helps maintain a stable input distribution and activation function helps the network learn more effectively by avoiding dead neurons.

In pooling layers, the dimensions are reduced by $(\lfloor \frac{W-K}{S} \rfloor + 1) \times (\lfloor \frac{H-K}{S} \rfloor + 1)$, where W is input width, H is input width, K is kernel size, S is stride.

After three blocks of convolutional layer and the pooling layer, the outputs are flatten from $28 \times 28 \times 64$ to 50176 neutrons in the input layer of fully-connected layer, and dropouted by 40% probability to prevent overfitting. The output layer consists of 2 neutrons for binary classification.

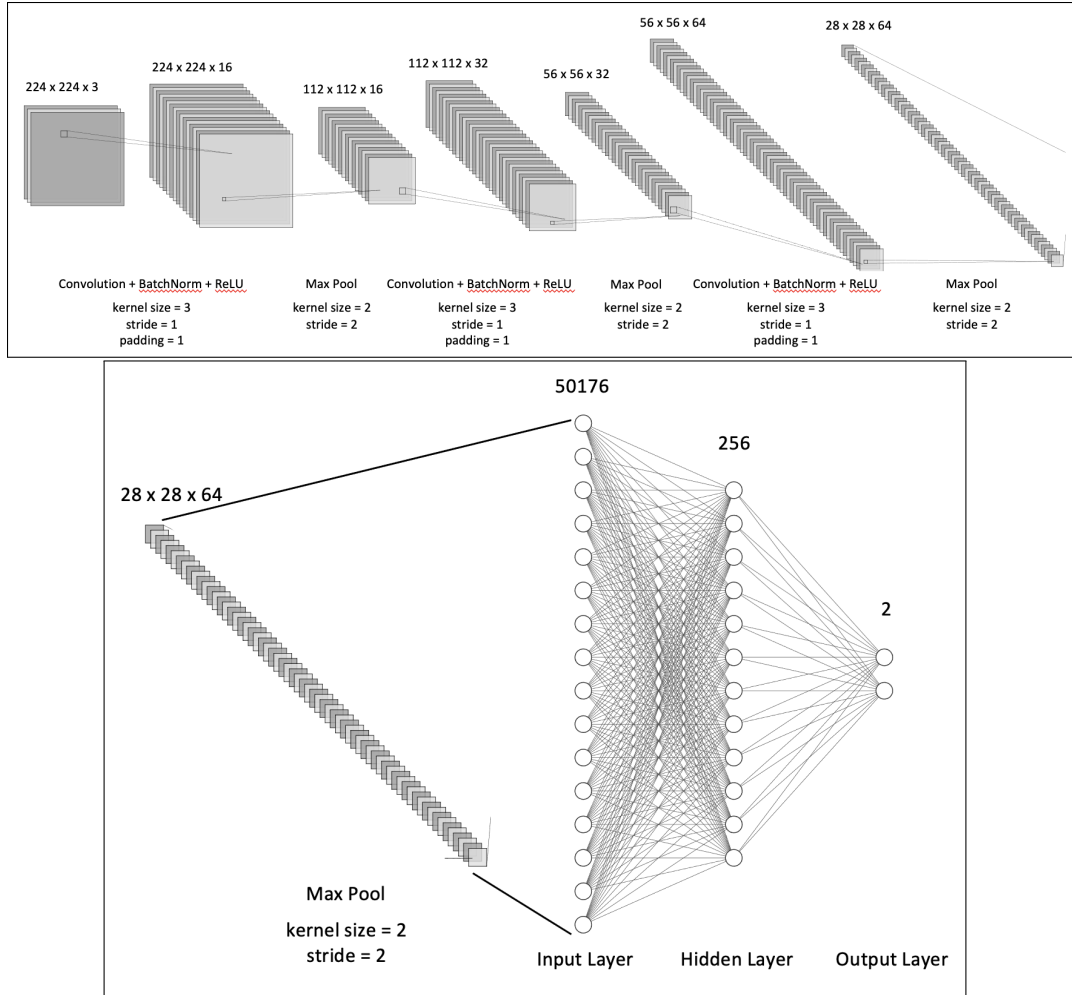


Figure 3: CNN architecture of AI Image Classifier

4.1.2 Training Process

The dataset is first loaded into two lists, the train and test datasets, where each image is labelled 0 for the AI-generated image, 1 for the human-crafted image. The dataset is fixed throughout the training.

During the model training, each image in the train dataset is passed through the model architecture and we obtain the two neuron logit values in the output layer. The logits are directly converted to class predictions by choosing the label with the highest logit. For example, if the output logits are [2.2, -1.3], the predicted class is 0 (AI Image) as $2.2 > -1.3$. We do not apply softmax function for prediction probabilities as the loss function we use is Cross Entropy Loss expecting logits as inputs. We also use Adaptive Moment Estimation (Adam) as an optimizer for adaptive learning rate and faster convergence.

After the classification of one image, the loss is calculated by Cross Entropy Loss and propagated back to the architecture layers. The accuracy is calculated by comparing the predicted class with the ground-truth class. All the loss and accuracy are cumulated for data visualization and analysis.

The training is repeated for all the images in train dataset in one epoch. The model is trained for 10 epochs and the following shows a training log. A complete Python program is attached in the Appendix B.

```

1 ===== Start training =====
2 Epoch [1/10], Loss: 0.4317, Train Acc: 82.40%
3 Test Acc: 83.46%
4 Epoch [2/10], Loss: 0.3268, Train Acc: 86.66%
5 Test Acc: 88.50%
6 Epoch [3/10], Loss: 0.2971, Train Acc: 87.93%
7 Test Acc: 87.73%
8 Epoch [4/10], Loss: 0.2719, Train Acc: 88.91%
9 Test Acc: 91.11%
10 Epoch [5/10], Loss: 0.2498, Train Acc: 89.97%
11 Test Acc: 89.16%
12 Epoch [6/10], Loss: 0.2338, Train Acc: 90.69%
13 Test Acc: 90.25%
14 Epoch [7/10], Loss: 0.2177, Train Acc: 91.40%
15 Test Acc: 90.51%
16 Epoch [8/10], Loss: 0.2025, Train Acc: 92.09%
17 Test Acc: 93.12%
18 Epoch [9/10], Loss: 0.1887, Train Acc: 92.56%
19 Test Acc: 91.15%
20 Epoch [10/10], Loss: 0.1749, Train Acc: 93.12%
21 Test Acc: 93.46%
22 ===== Finish training =====

```

Listing 1: CNN Model Training Log

4.1.3 Results and Analysis

Figure 4 shows two data visualizations of training loss and test accuracy over epochs.

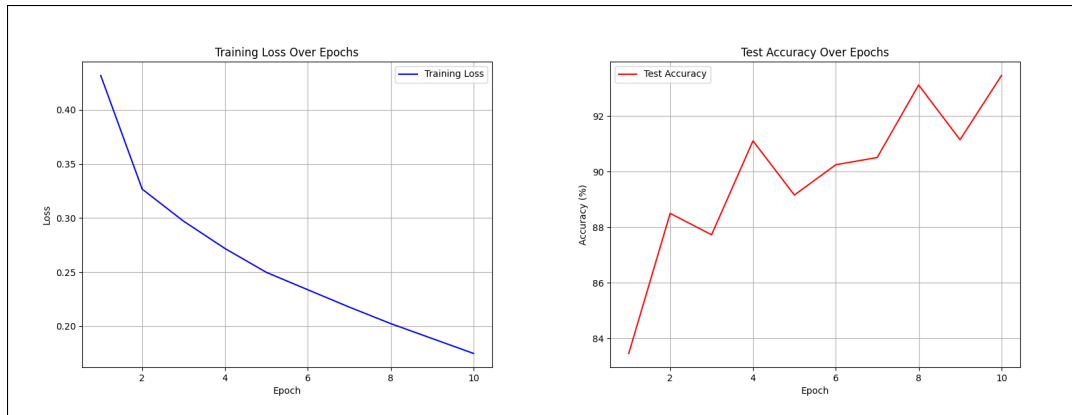


Figure 4: CNN Training Loss (left) and Test Accuracy (right)

In the training loss, the loss starts at 0.43 and constantly decreases to about 0.17 after 10 epochs. The steady downward trend indicates that the model learns and optimizes well on the training data. The loss has not fully flattened out by epoch 10, suggesting that the model can be improved with more training epochs.

In the test accuracy, the accuracy starts at 0.82 and increases to approximately 0.93 after 10 epochs, with some fluctuations along the trend. The overall upward trend is positive, showing that the model generalizes well to the test set as training progresses. The fluctuations suggest some instability in the learning process. This could be due to the reason of the model overshooting optimal weights with a high learning rate, or the model struggles with some variability in test dataset.

Figure 5 shows a confusion matrix of the classification model after 10 epochs. We calculate the accuracy, precision, recall, and F1 score based on the confusion matrix.

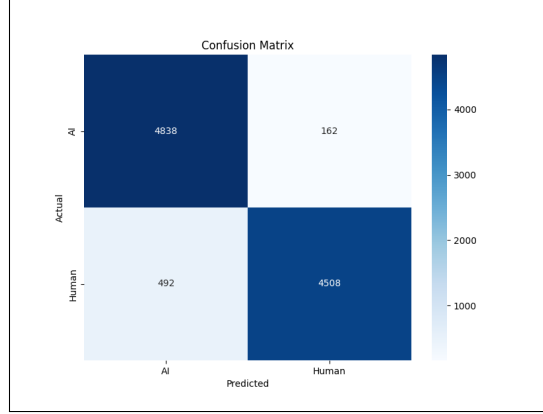


Figure 5: CNN Confusion Matrix

Accuracy	Precision	Recall	F1-Score
93.46%	90.77%	96.76%	93.67%

Table 1: Accuracy, Precision, Recall, and F1 Score

The model performs better at identifying AI images (higher recall) but has more false positives (492 human images misclassified as AI images) than false negatives (162 AI images misclassified as human images). This suggests the model might be slightly biased toward predicting "AI images". The model in general achieves a satisfactory accuracy of 93%.

4.2 ResNet

ResNet, or Residual Neural Network, is a pioneering deep learning architecture introduced in 2015 by researchers at Microsoft Research [11]. It tackles the challenge of training very deep neural networks by incorporating residual connections, which enable layers to learn residual functions relative to their inputs. This approach facilitates the training of networks with hundreds of layers, significantly enhancing performance in image recognition tasks. ResNet achieved remarkable success by winning the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2015, demonstrating its robustness and effectiveness in computer vision applications. Therefore, we applied ResNet for the classification task.

4.2.1 Architecture Design

ResNet50 is a 50-layer deep convolutional neural network designed to address vanishing gradients in deep networks through residual connections. The architecture begins with an input layer processing $224 \times 224 \times 3$ images, followed by an initial 7×7 convolution (stride 2) and 3×3 max pooling (stride 2), reducing spatial dimensions to $56 \times 56 \times 64$. These layers extract low-level features while halving resolution for computational efficiency[12].

The core of ResNet50 consists of four stages of bottleneck residual blocks. Each block compresses channels with a 1×1 convolution, processes features via a 3×3 convolution, and restores dimensions with another 1×1 convolution. A residual connection skips these operations, adding the input directly to the output to preserve gradient flow. **Stage 2** (3 blocks) maintains 56×56 resolution but increases depth to 256 channels, balancing detail retention and feature complexity[12].

Stage 3 (4 blocks) downsamples to 28×28 resolution using strided convolutions in the first block, expanding to 512 channels. This stage captures mid-level features like textures and patterns. **Stage 4** (6 blocks) further reduces resolution to 14×14 while increasing channels to 1024, enabling the network to learn hierarchical representations of objects[12].

Layer Name	Output Size	Configuration	Repetitions
Input	$224 \times 224 \times 3$	-	-
Conv1	$112 \times 112 \times 64$	7×7 , 64, stride 2	1
Pool1	$56 \times 56 \times 64$	3×3 max pool, stride 2	1
Stage 2 (Conv2_x)	$56 \times 56 \times 256$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix}$ Bottleneck block, Residual connection	3
Stage 3 (Conv3_x)	$28 \times 28 \times 512$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix}$ Bottleneck block, Residual connection	4
Stage 4 (Conv4_x)	$14 \times 14 \times 1024$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix}$ Bottleneck block, Residual connection	6
Stage 5 (Conv5_x)	$7 \times 7 \times 2048$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix}$ Bottleneck block, Residual connection	3
Global Avg Pool	$1 \times 1 \times 2048$	7×7 global average pool	1
FC	1000	Fully connected, softmax	1

Table 2: ResNet50 Architecture

The final residual stage (**Stage 5**, 3 blocks) operates at 7×7 resolution with 2048 channels, focusing on high-level semantic features. The network concludes with global average pooling, collapsing spatial dimensions to $1 \times 1 \times 2048$, and a fully connected layer for classification. This design ensures efficient training while achieving state-of-the-art accuracy on tackling tasks like image classification[12].

4.2.2 Training Process

The implementation uses a transfer learning approach with a pre-trained ResNet50 model to classify between human and AI-generated images. The base ResNet architecture, pre-trained on ImageNet, serves as a feature extractor with its weights frozen to preserve learned representations. Only the final classifier has been customized with a two-layer neural network (512 hidden units with ReLU activation and dropout regularization) that outputs two classes. This design allows the model to leverage general visual features while specializing in detecting subtle differences between human and AI-generated content.

During training, the model processes images in batches of 16 through 10 epochs, using Adam optimizer with a learning rate of 0.001 and Cross Entropy Loss. Performance is tracked through training losses and test accuracies, with the best-performing epoch's model weights saved for inference. The validation process produces a confusion matrix showing the model's discrimination capability between human and AI-generated content. This efficient approach enables accurate classification while requiring minimal computational resources by training only the task-specific layers rather than the entire network, complete Python program is attached in the Appendix C.

```

1 ===== Start training =====
2 Epoch [1/10], Loss: 0.3948, Train Acc: 81.87%
3 Test Acc: 88.98%
4 Epoch [2/10], Loss: 0.3383, Train Acc: 85.35%
5 Test Acc: 89.36%
6 Epoch [3/10], Loss: 0.3226, Train Acc: 86.22%
7 Test Acc: 90.44%
8 Epoch [4/10], Loss: 0.3118, Train Acc: 86.65%
9 Test Acc: 91.37%
10 Epoch [5/10], Loss: 0.3061, Train Acc: 87.02%
11 Test Acc: 90.33%

```

```

12 Epoch [6/10], Loss: 0.2957, Train Acc: 87.58%
13 Test Acc: 92.03%
14 Epoch [7/10], Loss: 0.2915, Train Acc: 87.65%
15 Test Acc: 91.61%
16 Epoch [8/10], Loss: 0.2889, Train Acc: 87.94%
17 Test Acc: 92.23%
18 Epoch [9/10], Loss: 0.2841, Train Acc: 88.46%
19 Test Acc: 92.52%
20 Epoch [10/10], Loss: 0.2802, Train Acc: 89.93%
21 Test Acc: 92.91%
22 ===== Finish training =====

```

Listing 2: ResNet Model Training Log

4.2.3 Results and Analysis

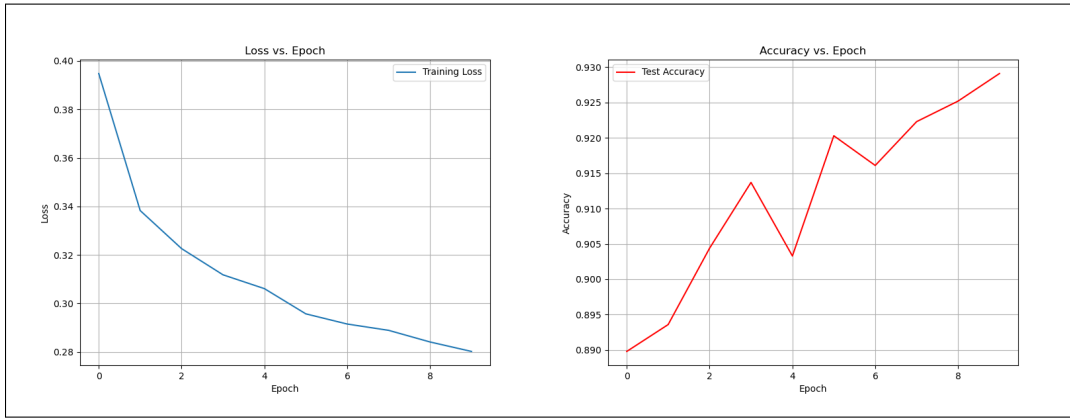


Figure 6: ResNet50 Training Loss (left) and Test Accuracy (right)

In the training loss, the loss starts at 0.39 and constantly decreases to about 0.12 after 10 epochs. The steady downward trend indicates that the model learns and optimizes moderately on the training data. The loss tends to be flattened out by epoch 10, indicating that the model can be improved slightly with more training epochs.

In the test accuracy, the accuracy starts at 0.82 and increases to approximately 0.93 after 10 epochs, with some fluctuations along the trend. The overall upward trend is positive, showing that the model performs good generalization to the test set as training progresses. The fluctuations suggest some instability in the learning process. This could be due to the reason of the batch training process, causing variation among the data samples in the batch.

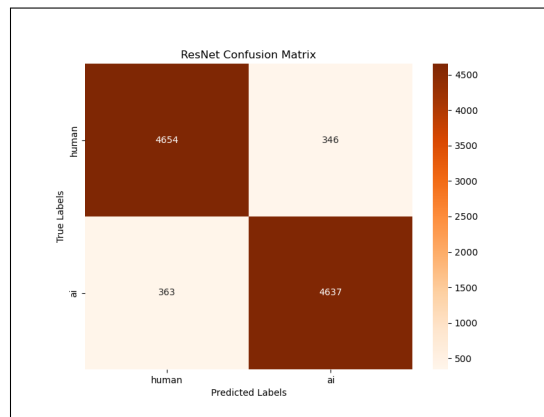


Figure 7: ResNet50 Confusion Matrix of Best Epoch

Accuracy	Precision	Recall	F1-Score
92.91%	93.05%	92.74%	92.89%

Table 3: Accuracy, Precision, Recall, and F1 Score

Figure 7 shows a confusion matrix of the classification model after 10 epochs. We can observe that the model is well balanced on the Accuracy, Precision, Recall, F1-Score as the True Negative equally high as True Positive, and False Negative equally low as False positive indicates the maturity of ResNet50 model on performing classification task among images.

4.3 Vision Transformer

The Vision Transformer (ViT) is a deep learning model introduced in 2020 by Dosovitskiy et al. for image classification, adapting the transformer architecture from natural language processing to computer vision [1]. Unlike convolutional neural networks (CNNs), ViT processes images by dividing them into fixed-size patches (e.g., 16×16 pixels), flattening each patch into a vector, and embedding it into a sequence of tokens. These tokens, along with a special classification (CLS) token, are fed into a series of transformer encoder layers that use self-attention to capture global relationships across the image. Positional embeddings are added to retain spatial information. ViT’s output, typically the CLS token’s embedding, is passed through a linear layer for classification. ViT excels in tasks like ImageNet classification when pre-trained on large datasets, offering advantages in capturing long-range dependencies compared to CNNs, though it requires substantial computational resources [13].

4.3.1 Architecture Design

Layer Name	Output Size	Configuration	Repetitions
Input	$224 \times 224 \times 3$	-	-
Patch Embedding	197×768	Patch size 16×16 , 768 dim, flatten and linear projection	1
CLS Token	1×768	Prepend learnable CLS token	1
Positional Embedding	197×768	Add learnable 1D positional embeddings	1
Transformer Encoder	197×768	Multi-Head Self-Attention (12 heads, 768 dim) LayerNorm, MLP (3072 dim, GELU) Residual connections, LayerNorm Dropout (0.1)	12
LayerNorm	197×768	Final LayerNorm	1
CLS Token Output	1×768	Extract CLS token embedding	1
MLP Head	1000	Linear layer, softmax	1

Table 4: ViT-B/16 Architecture

ViT-B/16 processes images as sequences of 16×16 patches, linearly projecting them to 768D embeddings. A learnable [CLS] token and positional embeddings are added to the 197 resulting tokens (14×14 patches + [CLS]). The core architecture consists of 12 identical transformer layers, each featuring:

- Multi-head self-attention (12 heads)
- MLP (3072 hidden units)
- Pre-layer normalization
- Residual connections

The model maintains a 197×768 representation throughout, applying dropout ($p = 0.1$) for regularization. After the final layer norm, the [CLS] token's 768D embedding is extracted for classification via a linear head. This pure transformer approach demonstrates that global self-attention can effectively replace convolutions for vision tasks.

Key features: Patch embeddings, learned positional encoding, [CLS] token aggregation, and standard transformer components adapted for visual data.

4.3.2 Training Process

During training, the model processes images in batches of 16 through 10 epochs, using Adam optimizer with a slightly lower learning rate of 0.0005 to avoid destabilizing the sensitive transformer weights, and Cross Entropy Loss for classification. Performance metrics including training loss and test accuracy are tracked throughout the process, with the best-performing model weights saved for inference. The validation phase generates a confusion matrix that visualizes the model's effectiveness in discriminating between human and AI-generated content, complete Python program is attached in the Appendix D.

```

1 ===== Start training =====
2 Epoch [1/10], Loss: 0.1763, Train Acc: 93.42%
3 Test Acc: 95.24%
4 Epoch [2/10], Loss: 0.1683, Train Acc: 94.54%
5 Test Acc: 96.36%
6 Epoch [3/10], Loss: 0.1526, Train Acc: 94.63%
7 Test Acc: 96.44%
8 Epoch [4/10], Loss: 0.1398, Train Acc: 93.56%
9 Test Acc: 95.37%
10 Epoch [5/10], Loss: 0.1312, Train Acc: 94.52%
11 Test Acc: 96.33%
12 Epoch [6/10], Loss: 0.1301, Train Acc: 95.10%
13 Test Acc: 96.91%
14 Epoch [7/10], Loss: 0.1247, Train Acc: 95.01%
15 Test Acc: 96.82%
16 Epoch [8/10], Loss: 0.1178, Train Acc: 95.23%
17 Test Acc: 97.04%
18 Epoch [9/10], Loss: 0.0941, Train Acc: 95.08%
19 Test Acc: 96.89%
20 Epoch [10/10], Loss: 0.0902, Train Acc: 95.30%
21 Test Acc: 97.11%
22 ===== Finish training =====

```

Listing 3: ViT Model Training Log

4.3.3 Results and Analysis

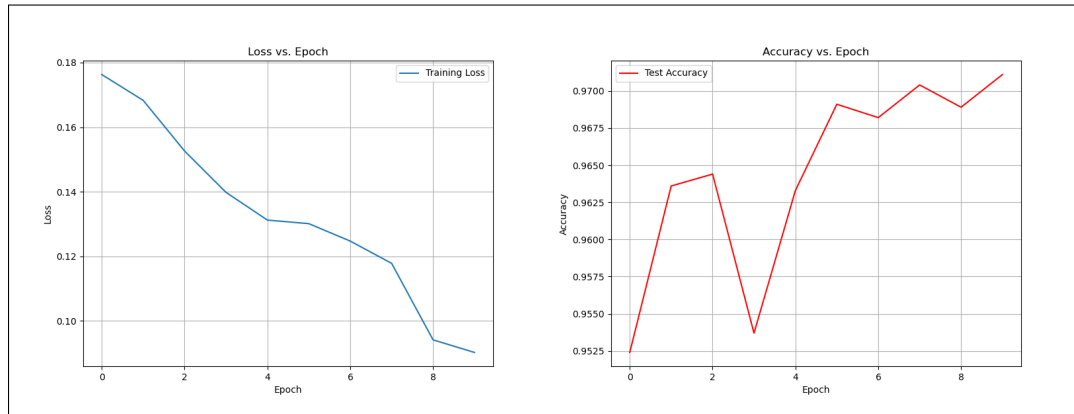


Figure 8: ViT Training Loss (left) and Test Accuracy (right)

In the training loss, the loss starts at 0.18 and constantly decreases to less than 0.1 after 10 epochs. The training loss is satisfactory that indicates the incredibility of ViT on handling image classification tasks. The loss is still having a linear trend after 10 epochs, indicating that the model can be improved with more epochs.

In the test accuracy, the accuracy starts at 0.95 and increases to over 0.97 after 10 epochs, with some fluctuations along the trend. The overall upward trend is positive, showing that the model performs good generalization to the test set as training progresses. The fluctuations suggest some instability in the learning process. This could be due to the reason of the batch training process, causing variation among the data samples in the batch.

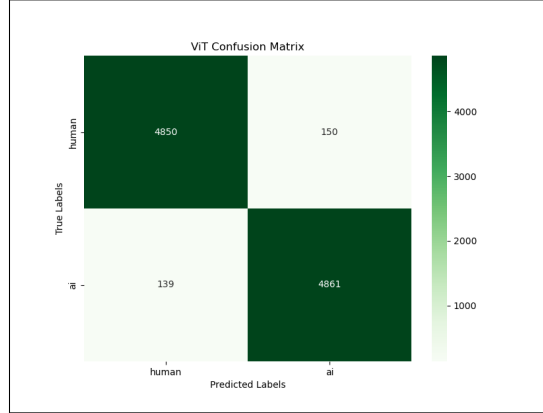


Figure 9: ViT Confusion Matrix of Best Epoch

Accuracy	Precision	Recall	F1-Score
97.11%	97.01%	97.22%	97.11%

Table 5: Accuracy, Precision, Recall, and F1 Score

Figure 9 shows a confusion matrix of the classification model after 10 epochs. We can observe that the model is well balanced on the Accuracy, Precision, Recall, F1-Score of roughly 97% as the True Negative equally high as True Positive, and False Negative equally low as False positive indicates the maturity of ViT model on performing classification task among images.

5 Discussion

Among all 3 models implemented, ViT successfully achieved high accuracy with 97%, which is higher than our expectation of 95% , which provide a good evidence on its ability to perform classification tasks. While the traditional hierarchical processing like out custom CNN model and ResNet50 performs slightly lower than our expectaion, yet still achive a considerable accuracy of around 93%.

We have discussed regarding the reason, and we believe that Vision Transformer (ViT) outperforms ResNet50 primarily because its self-attention mechanism captures global relationships between image patches simultaneously, unlike ResNet’s hierarchical local processing through convolutional layers. ViT’s ability to model long-range dependencies across the entire image is particularly advantageous for detecting subtle artifacts and unnatural patterns characteristic of AI-generated content. Additionally, ViT’s patch-based representation and fewer inductive biases about image structure allow it to better identify the statistical irregularities and perceptual inconsistencies that distinguish AI-generated images from human-created ones [14].

Also, the accuracy of our custom CNN and ResNet50 performs similarly, and we believe the CNN's tailored architecture is optimized specifically for this binary classification task, while we have used the pretrained ResNet50 model that may not specifically suitable for such task, with a too small learning rate that hindered the learning process.

Due to the limited project time and report space, we could not do further studies on our AI image detector. Therefore, we suggest one possible further action is to deploy our models to applications such as websites, so that our detector is practical to use, and we can collect feedback and more test datasets from users to improve our model. Another possible further study is to investigate the datasets and the potential reasons for misclassified images. There could be some false data or findings of AI-generated images perfectly imitating human images, which enhances our understanding towards this topic.

6 Conclusion

This paper targets the problem of high similarity between AI-generated images and human-crafted images, proposes three classification models, namely custom Convolutional Neural Network, ResNet and Vision Transformer, and analyzes their performance and effectiveness.

As the generative models progress, AI-generated images, which forge the human-crafted images, are widely spread on the Internet, posing artwork copyright infringement and information security challenges. An effective solution would be to develop a highly accurate AI image detector that classifies the AI-generated images and human-crafted images. The architecture behind the detector is a binary classification model, and three different classification models, namely custom Convolutional Neural Network, ResNet and Vision Transformer are presented.

Data has been first collected, extracted, and pre-processed by normalization and augmentation methods. Custom Convolutional Neural Network, mainly composed of 3 blocks of convolutional layers and max-pooling layers, and a 3-fully-connected layers, provide a good accuracy of 93%. For ResNet, it achieves a considerable accuracy of 93%, which is in line with the model performance. For Visual Transformer, it achieves a satisfactory accuracy of 97%, which we believe is capable of providing a reliable suggestion on distinguishing AI-generated images and human-crafted images.

The result shows that all three approaches are robust models for AI-human image classification. It is suggested that deployment to the application and in-depth studies on the dataset could be further explored.

References

- [1] "South China Morning Post," *South China Morning Post*, Apr. 02, 2025. https://www.scmp.com/tech/big-tech/article/3304828/chatgpt-usage-hits-record-studio-ghibli-style-ai-images-go-viral?module=perpetual_scroll_0&pgtype=article (accessed Apr. 04, 2025).
- [2] N. Tiku, "AI can now create any image in seconds, bringing wonder and danger," *Washington Post*, Sep. 28, 2022. <https://www.washingtonpost.com/technology/interactive/2022/artificial-intelligence-images-dall-e/>
- [3] H. Y. Yan, G. Morrow, K.-C. Yang, and J. Wihbey, "The origin of public concerns over AI supercharging misinformation in the 2024 U.S. presidential election," *The origin of public concerns over AI supercharging misinformation in the 2024 U.S. presidential election*, Jan. 2025, doi: <https://doi.org/10.37016/mr-2020-171>.
- [4] A. Sala, "AI vs. Human-Generated Images," *Kaggle.com*, 2025. <https://www.kaggle.com/datasets/alessandrasala79/ai-vs-human-generated-dataset> (accessed May 05, 2025).
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, May 2012.
- [6] "Understanding CNN for Image Processing | Svitla Systems," *Svitla Systems*, Jun. 21, 2024. <https://svitla.com/blog/cnn-for-image-processing/>
- [7] "The evolution of image classification explained," *Stanford.edu*, 2012. <https://stanford.edu/shervine/blog/evolution-image-classification-explained>
- [8] L. Chen, S. Li, Q. Bai, J. Yang, S. Jiang, and Y. Miao, "Review of Image Classification Algorithms Based on Convolutional Neural Networks," *Remote Sensing*, vol. 13, no. 22, p. 4712, Nov. 2021, doi: <https://doi.org/10.3390/rs13224712>.
- [9] Keiron O'Shea and R. R. Nash, "An Introduction to Convolutional Neural Networks," *arXiv* (Cornell University), Nov. 2015, doi: <https://doi.org/10.48550/arxiv.1511.08458>.
- [10] "Convolutional Neural Network - an overview | ScienceDirect Topics," *www.sciencedirect.com*. <https://www.sciencedirect.com/topics/engineering/convolutional-neural-network>
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *www.cv-foundation.org* Dec. 2015. Available: https://www.cv-foundation.org/openaccess/content_cvpr_2016/papers/He_Deep_Residual_Learning_CVPR_2016_paper.pdf
- [12] B. Koonce, "ResNet 50," Apress eBooks, pp. 63–72, Jan. 2021, doi:https://doi.org/10.1007/978-1-4842-6168-2_6
- [13] S.-H. Tsang, "Review: Vision Transformer (ViT)," *Medium*, Feb. 05, 2022. <https://sh-tsang.medium.com/review-vision-transformer-vit-406568603de0>
- [14] Nouar AlDahoul and Y. Zaki, "Detecting AI-Generated Images Using Vision Transformers: A Robust Approach for Safeguarding Visual Media Integrity," Jan. 2025, doi: <https://doi.org/10.2139/ssrn.5029893>

Appendices

A Data Transformation Python Code Snippet

```
1 import torchvision.transforms as transforms
2
3 train_transform = transforms.Compose([
4     transforms.Resize((224, 224)),
5     transforms.RandomHorizontalFlip(p=0.5),
6     transforms.RandomRotation(degrees=10),
7     transforms.ColorJitter(brightness=0.2, contrast=0.2),
8     transforms.ToTensor(),
9     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
10     0.224, 0.225])
11 ])
12
13 test_transform = transforms.Compose([
14     transforms.Resize((224, 224)),
15     transforms.ToTensor(),
16     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
17     0.224, 0.225])
18 ])
```

B CNN Image Classifier Python Program

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torchvision
5 import torchvision.transforms as transforms
6 from torch.utils.data import DataLoader
7 from torchvision.datasets import ImageFolder
8 import os
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11 import numpy as np
12 from sklearn.metrics import confusion_matrix
13
14 train_transform = transforms.Compose([
15     transforms.Resize((224, 224)),
16     transforms.RandomHorizontalFlip(p=0.5),
17     transforms.RandomRotation(degrees=10),
18     transforms.ColorJitter(brightness=0.2, contrast=0.2),
19     transforms.ToTensor(),
20     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
21     0.224, 0.225])
22 ])
23
24 test_transform = transforms.Compose([
25     transforms.Resize((224, 224)),
26     transforms.ToTensor(),
27     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
28     0.224, 0.225])
29 ])
30
31 def load_dataset(data_dir):
32     train_dataset = ImageFolder(root=os.path.join(data_dir, 'train'),
33     transform=train_transform)
34     test_dataset = ImageFolder(root=os.path.join(data_dir, 'test'),
35     transform=test_transform)
```

```

34     train_loader = DataLoader(train_dataset, batch_size=16, shuffle=
True)
35     test_loader = DataLoader(test_dataset, batch_size=16, shuffle=
False)
36
37     return train_loader, test_loader
38
39 class CNN(nn.Module):
40     def __init__(self):
41         super(CNN, self).__init__()
42
43         self.relu = nn.ReLU(inplace=True)
44         self.pool = nn.MaxPool2d(kernel_size=2, stride=2)
45
46         # Convolutional Layer 1
47         self.conv1 = nn.Conv2d(in_channels=3, out_channels=16,
kernel_size=3, stride=1, padding=1, bias=False)
48         self.bn1 = nn.BatchNorm2d(16)
49
50         # Convolutional Layer 2
51         self.conv2 = nn.Conv2d(in_channels=16, out_channels=32,
kernel_size=3, stride=1, padding=1, bias=False)
52         self.bn2 = nn.BatchNorm2d(32)
53
54         # Convolutional Layer 3
55         self.conv3 = nn.Conv2d(in_channels=32, out_channels=64,
kernel_size=3, stride=1, padding=1, bias=False)
56         self.bn3 = nn.BatchNorm2d(64)
57
58         # Fully Connected Layers
59         self.dropout = nn.Dropout(0.4)
60         self.fc1 = nn.Linear(28 * 28 * 64, 256)
61         self.fc2 = nn.Linear(256, 2)
62
63     def forward(self, x):
64
65         x = self.pool(self.relu(self.bn1(self.conv1(x))))
66         x = self.pool(self.relu(self.bn2(self.conv2(x))))
67         x = self.pool(self.relu(self.bn3(self.conv3(x))))
68
69         x = x.view(x.size(0), -1)
70
71         x = self.dropout(x)
72         x = self.relu(self.fc1(x))
73         x = self.fc2(x)
74         return x
75
76 def train_model(model, train_loader, test_loader, criterion, optimizer
, num_epochs, device):
77
78     best_val_acc = 0.0
79     train_losses = []
80     test accuracies = []
81
82     for epoch in range(num_epochs):
83         model.train()
84         running_loss = 0.0
85         correct = 0
86         total = 0
87
88         for images, labels in train_loader:
89             images, labels = images.to(device), labels.to(device)
90             optimizer.zero_grad()
91             outputs = model(images)
92             loss = criterion(outputs, labels)

```

```

93         loss.backward()
94         optimizer.step()
95
96         running_loss += loss.item()
97         _, predicted = torch.max(outputs, 1)
98         total += labels.size(0)
99         correct += (predicted == labels).sum().item()
100
101     train_loss = running_loss / len(train_loader)
102     train_losses.append(train_loss)
103     train_acc = 100 * correct / total
104     print(f"Epoch [{epoch+1}/{num_epochs}], Loss: {running_loss/len(train_loader):.4f}, Train Acc: {train_acc:.2f}%")
105
106     # Testing
107     model.eval()
108     test_correct = 0
109     test_total = 0
110     with torch.no_grad():
111         for images, labels in test_loader:
112             images, labels = images.to(device), labels.to(device)
113             outputs = model(images)
114             _, predicted = torch.max(outputs, 1)
115             test_total += labels.size(0)
116             test_correct += (predicted == labels).sum().item()
117
118     val_acc = 100 * test_correct / test_total
119     test accuracies.append(val_acc)
120     print(f"Test Acc: {val_acc:.2f}%")
121
122     # Save best model
123     if val_acc > best_val_acc:
124         best_val_acc = val_acc
125         torch.save(model.state_dict(), 'cnn_model.pth')
126
127     return train_losses, test accuracies
128
129 # Plot Graphs
130 def plot_training_loss(losses):
131     save_path='training_loss.png'
132     plt.figure(figsize=(8, 6))
133     plt.plot(range(1, len(losses) + 1), losses, 'b-', label='Training Loss')
134     plt.title('Training Loss Over Epochs')
135     plt.xlabel('Epoch')
136     plt.ylabel('Loss')
137     plt.grid(True)
138     plt.legend()
139     plt.savefig(save_path)
140     plt.close()
141
142 def plot_test_accuracy(accuracies):
143     save_path='test_accuracy.png'
144     plt.figure(figsize=(8, 6))
145     plt.plot(range(1, len(accuracies) + 1), accuracies, 'r-', label='Test Accuracy')
146     plt.title('Test Accuracy Over Epochs')
147     plt.xlabel('Epoch')
148     plt.ylabel('Accuracy (%)')
149     plt.grid(True)
150     plt.legend()
151     plt.savefig(save_path)
152     plt.close()
153
154 def plot_confusion_matrix(true_labels, pred_labels):

```



```

155     save_path='confusion_matrix.png'
156     classes = ['AI', 'Human']
157     cm = confusion_matrix(true_labels, pred_labels)
158     plt.figure(figsize=(8, 6))
159     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=
160     classes, yticklabels=classes)
161     plt.title('Confusion Matrix')
162     plt.xlabel('Predicted')
163     plt.ylabel('Actual')
164     plt.savefig(save_path)
165     plt.close()
166
167 def compute_confusion_matrix(model, val_loader, device):
168     model.eval()
169     true_labels = []
170     pred_labels = []
171
172     with torch.no_grad():
173         for images, labels in val_loader:
174             images, labels = images.to(device), labels.to(device)
175             outputs = model(images)
176             _, predicted = torch.max(outputs, 1)
177             true_labels.extend(labels.cpu().numpy())
178             pred_labels.extend(predicted.cpu().numpy())
179
180     return true_labels, pred_labels
181
182 if __name__ == "__main__":
183     data_dir = 'dataset'
184     device = torch.device("cuda" if torch.cuda.is_available() else "
185     cpu")
186     num_epochs = 10
187
188     try:
189         train_loader, test_loader = load_dataset(data_dir)
190
191         # Verify data loading
192         for images, labels in train_loader:
193             print(f"Training batch shape: {images.shape}, Labels: {
194             labels}")
195             break
196
197         for images, labels in test_loader:
198             print(f"Testing batch shape: {images.shape}, Labels: {
199             labels}")
200             break
201
202         print(f"Class mapping: {train_loader.dataset.class_to_idx}")
203
204         # Initialize model, loss, and optimizer
205         model = CNN().to(device)
206         criterion = nn.CrossEntropyLoss()
207         optimizer = optim.Adam(model.parameters(), lr=0.001)
208
209         # Train the model
210         print("==== Start training =====")
211         train_losses, test accuracies = train_model(model,
212         train_loader, test_loader, criterion, optimizer, num_epochs,
213         device)
214         print("==== Finish training =====")
215
216         # Plot graphs
217         plot_training_loss(train_losses)
218         plot_test_accuracy(test accuracies)

```

```

214         true_labels, pred_labels = compute_confusion_matrix(model,
215             test_loader, device)
216         plot_confusion_matrix(true_labels, pred_labels)
217
218     except Exception as e:
219         print(f"Error: {str(e)}")

```

C ResNet Image Classifier Python Program

```

1  import torch
2  import torch.nn as nn
3  import torch.optim as optim
4  from torch.utils.data import Dataset
5  from torchvision import models
6  from torchvision.models import ResNet50_Weights # Import weights
7  from PIL import Image
8  import os
9  import matplotlib.pyplot as plt
10 from tqdm import tqdm
11 from sklearn.metrics import confusion_matrix
12 import seaborn as sns
13
14 # Import transforms and loader from transform.py
15 from transform import test_transform, load_dataset
16
17 # Configuration
18 class Config:
19     IMAGE_SIZE = 224
20     BATCH_SIZE = 16
21     NUM_EPOCHS = 10
22     LEARNING_RATE = 0.001
23     NUM_CLASSES = 2
24     DATA_DIR = r"\dataset" # Directory to load data
25     DEVICE = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
26     FIGURE_DIR = r"\ResNet_Result" # Directory to save figures
27
28 os.makedirs(Config.FIGURE_DIR, exist_ok=True)
29
30 # Create datasets and dataloaders using transforms from transform.py
31 def create_dataloaders():
32     print("Loading datasets...")
33     train_loader, test_loader = load_dataset(Config.DATA_DIR, Config.BATCH_SIZE)
34     print(f"Using device: {Config.DEVICE}")
35     return train_loader, test_loader
36
37 # ResNet-based Model
38 class ResNet(nn.Module):
39     def __init__(self, num_classes=Config.NUM_CLASSES):
40         super(ResNet, self).__init__()
41         # Use a pretrained ResNet model
42         self.resnet = models.resnet50(weights=ResNet50_Weights.IMAGENET1K_V1)
43
44         for param in self.resnet.parameters():
45             param.requires_grad = False
46
47         num_fts = self.resnet.fc.in_features
48         self.resnet.fc = nn.Sequential(
49             nn.Linear(num_fts, 512),
50             nn.ReLU(),
51             nn.Dropout(0.5),

```

```

52         nn.Linear(512, num_classes)
53     )
54
55     def forward(self, x):
56         return self.resnet(x)
57
58 # Plot training and validation metrics
59 def plot_metrics(train_losses, train_accuracies, test_losses,
60                 test_accuracies):
61     plt.figure(figsize=(10, 5))
62     plt.plot(train_losses, label='Training Loss')
63     plt.xlabel('Epoch')
64     plt.ylabel('Loss')
65     plt.title('Loss vs. Epoch')
66     plt.legend()
67     plt.grid(True)
68     plt.savefig(os.path.join(Config.FIGURE_DIR, 'train_loss_plot.png'),
69                 dpi=300)
70
71     plt.figure(figsize=(10, 5))
72     plt.plot(test_accuracies, label='Test Accuracy')
73     plt.xlabel('Epoch')
74     plt.ylabel('Accuracy')
75     plt.title('Accuracy vs. Epoch')
76     plt.legend()
77     plt.grid(True)
78     plt.savefig(os.path.join(Config.FIGURE_DIR, 'test_accuracy_plot.
79                 png'), dpi=300)
80
81     plt.show()
82
83 # Plot confusion matrix
84 def plot_confusion_matrix(y_true, y_pred, class_names):
85     cm = confusion_matrix(y_true, y_pred)
86     plt.figure(figsize=(8, 6))
87     sns.heatmap(cm, annot=True, fmt='d', cmap='Oranges', xticklabels=
88                 class_names, yticklabels=class_names)
89     plt.xlabel('Predicted Labels')
90     plt.ylabel('True Labels')
91     plt.title('ResNet Confusion Matrix')
92     plt.savefig(os.path.join(Config.FIGURE_DIR, 'confusion_matrix.png'),
93                 dpi=300)
94     plt.show()
95
96 def train_model(model, dataloaders, criterion, optimizer, num_epochs=
97                 Config.NUM_EPOCHS):
98     best_acc = 0.0
99     best_epoch = 0
100
101     # Lists to store metrics for plotting
102     train_losses = []
103     train_accuracies = []
104     test_losses = []
105     test_accuracies = []
106     best_epoch_labels = []
107     best_epoch_predictions = []
108
109     for epoch in range(num_epochs):
110         print(f'Epoch {epoch+1}/{num_epochs}')
111         print('-' * 20)
112
113         current_epoch_labels = []
114         current_epoch_predictions = []
115
116         for phase in ['train', 'test']:

```

```

111         if phase == 'train':
112             model.train()
113         else:
114             model.eval()
115
116         running_loss = 0.0
117         running_corrects = 0
118
119         total = len(dataloaders[phase])
120
121         #progress bar
122         with tqdm(total=total, desc=f'{phase}') as pbar:
123             for inputs, labels in dataloaders[phase]:
124                 inputs = inputs.to(Config.DEVICE)
125                 labels = labels.to(Config.DEVICE)
126
127                 optimizer.zero_grad()
128
129                 with torch.set_grad_enabled(phase == 'train'):
130                     outputs = model(inputs)
131                     _, preds = torch.max(outputs, 1)
132                     loss = criterion(outputs, labels)
133
134                     if phase == 'train':
135                         loss.backward()
136                         optimizer.step()
137
138                     running_loss += loss.item() * inputs.size(0)
139                     running_corrects += torch.sum(preds == labels.data
140 )
141
142                     # Store test predictions for current epoch
143                     if phase == 'test':
144                         current_epoch_labels.extend(labels.cpu().numpy
145 )
146                         current_epoch_predictions.extend(preds.cpu().
147 numpy())
148
149                     pbar.update(1)
150                     pbar.set_postfix({'loss': f'{loss.item():.4f}'})
151
152                     # Calculate epoch statistics
153                     epoch_loss = running_loss / len(dataloaders[phase].dataset
154 )
155                     epoch_acc = running_corrects.double() / len(dataloaders[
156 phase].dataset)
157
158                     print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f
159 }')
160
161                     # Store metrics for plotting
162                     if phase == 'train':
163                         train_losses.append(epoch_loss)
164                         train_accuracies.append(epoch_acc.cpu().numpy())
165                     else:
166                         test_losses.append(epoch_loss)
167                         test_accuracies.append(epoch_acc.cpu().numpy())
168
169                     if epoch_acc > best_acc:
170                         best_acc = epoch_acc
171                         best_epoch = epoch
172                         torch.save(model.state_dict(), 'best_ResNet_model.
173 pth')
174
175                     best_epoch_labels = current_epoch_labels.copy()

```

```

169         best_epoch_predictions = current_epoch_predictions
170     .copy()
171
172     print()
173     print(f'Best test accuracy: {best_acc:.4f} at epoch {best_epoch+1}
174     ')
175
176     plot_metrics(train_losses, train_accuracies, test_losses,
177     test_accuracies)
178
179     print(f"Creating confusion matrix from epoch {best_epoch+1} (best
180     accuracy)")
181     plot_confusion_matrix(best_epoch_labels, best_epoch_predictions,
182     class_names=['human', 'ai'])
183
184     return model
185
186 def main():
187     print(f"PyTorch version: {torch.__version__}")
188
189     model = ResNet().to(Config.DEVICE)
190
191     criterion = nn.CrossEntropyLoss()
192
193     optimizer = optim.Adam(filter(lambda p: p.requires_grad, model.
194     parameters()),
195                             lr=Config.LEARNING_RATE)
196
197     train_loader, test_loader = create_dataloaders()
198     dataloaders = {'train': train_loader, 'test': test_loader}
199
200     print(model)
201
202     model = train_model(model, dataloaders, criterion, optimizer,
203     Config.NUM_EPOCHS)
204
205     print('Training complete. Models saved.')
206
207 if __name__ == '__main__':
208     main()
209
210 # Interface for classification
211 def predict_image(image_path, model_path='best_Resnet_model.pth'):
212     # Load the model
213     model = ResNet().to(Config.DEVICE)
214     model.load_state_dict(torch.load(model_path, map_location=Config.
215     DEVICE))
216     model.eval()
217
218     image = Image.open(image_path).convert('RGB')
219     image_tensor = test_transform(image).unsqueeze(0).to(Config.DEVICE
220     )
221
222     with torch.no_grad():
223         outputs = model(image_tensor)
224         _, preds = torch.max(outputs, 1)
225         probabilities = torch.nn.functional.softmax(outputs, dim=1)
226
227     class_names = ['human', 'ai']
228     return {
229         'class': class_names[preds.item()],
230         'confidence': probabilities[0][preds.item()].item(),
231         'human_prob': probabilities[0][0].item(),
232         'ai_prob': probabilities[0][1].item(),

```

```

225         'is_human': preds.item() == 0
226     }

```

D Visual Transformer Image Classifier Python Program

```

1  import torch
2  import torch.nn as nn
3  import torch.optim as optim
4  from torch.utils.data import Dataset
5  from torchvision import models
6  from torchvision.models import ViT_B_16_Weights # Import weights for
   ViT
7  from PIL import Image
8  import os
9  import matplotlib.pyplot as plt
10 from tqdm import tqdm
11 from sklearn.metrics import confusion_matrix
12 import seaborn as sns
13 from transform import test_transform, load_dataset
14
15 # Configuration
16 class Config:
17     IMAGE_SIZE = 224
18     BATCH_SIZE = 16
19     NUM_EPOCHS = 10
20     LEARNING_RATE = 0.0005
21     NUM_CLASSES = 2
22     DATA_DIR = r"\dataset"
23     DEVICE = torch.device('cuda' if torch.cuda.is_available() else '
   cpu')
24     FIGURE_DIR = r"\ViT_Result"
25
26 os.makedirs(Config.FIGURE_DIR, exist_ok=True)
27
28 # Create datasets and dataloaders
29 def create_dataloaders():
30     print("Loading datasets...")
31     train_loader, test_loader = load_dataset(Config.DATA_DIR, Config.
   BATCH_SIZE)
32     print(f"Using device: {Config.DEVICE}")
33     return train_loader, test_loader
34
35 class ViTClassifier(nn.Module):
36     def __init__(self, num_classes=Config.NUM_CLASSES):
37         super(ViTClassifier, self).__init__()
38         self.vit = models.vit_b_16(weights=ViT_B_16_Weights.
   IMAGENET1K_V1)
39
40         for param in self.vit.parameters():
41             param.requires_grad = False
42
43         # Replace the head (classification layer)
44         self.vit.heads = nn.Sequential(
45             nn.Linear(768, 512),
46             nn.LayerNorm(512),
47             nn.GELU(),
48             nn.Dropout(0.1),
49             nn.Linear(512, num_classes)
50         )
51
52         for param in self.vit.heads.parameters():
53             param.requires_grad = True
54

```

```

55         self.vit.encoder.pos_embedding.requires_grad = True
56
57     def forward(self, x):
58         return self.vit(x)
59
60 # Plot training and validation metrics
61 def plot_metrics(train_losses, train_accuracies, test_losses,
62                 test_accuracies):
63     plt.figure(figsize=(10, 5))
64     plt.plot(train_losses, label='Training Loss')
65     plt.xlabel('Epoch')
66     plt.ylabel('Loss')
67     plt.title('Loss vs. Epoch')
68     plt.legend()
69     plt.grid(True)
70     plt.savefig(os.path.join(Config.FIGURE_DIR, 'train_loss_plot.png'),
71                 dpi=300)
72
73     plt.figure(figsize=(10, 5))
74     plt.plot(test_accuracies, label='Test Accuracy')
75     plt.xlabel('Epoch')
76     plt.ylabel('Accuracy')
77     plt.title('Accuracy vs. Epoch')
78     plt.legend()
79     plt.grid(True)
80     plt.savefig(os.path.join(Config.FIGURE_DIR, 'test_accuracy_plot.
81                 png'), dpi=300)
82
83     plt.show()
84
85 # Plot confusion matrix
86 def plot_confusion_matrix(y_true, y_pred, class_names):
87     cm = confusion_matrix(y_true, y_pred)
88     plt.figure(figsize=(8, 6))
89     sns.heatmap(cm, annot=True, fmt='d', cmap='Greens', xticklabels=
90                 class_names, yticklabels=class_names)
91     plt.xlabel('Predicted Labels')
92     plt.ylabel('True Labels')
93     plt.title('ViT Confusion Matrix')
94     plt.savefig(os.path.join(Config.FIGURE_DIR, 'confusion_matrix.png'),
95                 dpi=300)
96     plt.show()
97
98 def train_model(model, dataloaders, criterion, optimizer, num_epochs=
99                 Config.NUM_EPOCHS):
100     best_acc = 0.0
101     best_epoch = 0
102
103     # Lists to store metrics for plotting
104     train_losses = []
105     train_accuracies = []
106     test_losses = []
107     test_accuracies = []
108     best_epoch_labels = []
109     best_epoch_predictions = []
110
111     for epoch in range(num_epochs):
112         print(f'Epoch {epoch+1}/{num_epochs}')
113         print('-' * 20)
114
115         current_epoch_labels = []
116         current_epoch_predictions = []
117
118         for phase in ['train', 'test']:
119             if phase == 'train':

```

```

114         model.train()
115     else:
116         model.eval()
117
118     running_loss = 0.0
119     running_corrects = 0
120
121     total = len(dataloaders[phase])
122
123     #progress bar
124     with tqdm(total=total, desc=f'{phase}') as pbar:
125         for inputs, labels in dataloaders[phase]:
126             inputs = inputs.to(Config.DEVICE)
127             labels = labels.to(Config.DEVICE)
128
129             optimizer.zero_grad()
130
131             with torch.set_grad_enabled(phase == 'train'):
132                 outputs = model(inputs)
133                 _, preds = torch.max(outputs, 1)
134                 loss = criterion(outputs, labels)
135
136                 if phase == 'train':
137                     loss.backward()
138                     optimizer.step()
139
140                 running_loss += loss.item() * inputs.size(0)
141                 running_corrects += torch.sum(preds == labels.data
142 )
143
144                 if phase == 'test':
145                     current_epoch_labels.extend(labels.cpu().numpy
146 )
147                     current_epoch_predictions.extend(preds.cpu().
148 numpy())
149
150                 pbar.update(1)
151                 pbar.set_postfix({'loss': f'{loss.item():.4f}'})
152
153     epoch_loss = running_loss / len(dataloaders[phase].dataset
154 )
155     epoch_acc = running_corrects.double() / len(dataloaders[
156 phase].dataset)
157
158     print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f
159 }')
160
161     # Store metrics for plotting
162     if phase == 'train':
163         train_losses.append(epoch_loss)
164         train_accuracies.append(epoch_acc.cpu().numpy())
165     else:
166         test_losses.append(epoch_loss)
167         test_accuracies.append(epoch_acc.cpu().numpy())
168
169     if epoch_acc > best_acc:
170         best_acc = epoch_acc
171         best_epoch = epoch
172         torch.save(model.state_dict(), 'best_ResNet_model.
173 pth')
174
175         best_epoch_labels = current_epoch_labels.copy()
176         best_epoch_predictions = current_epoch_predictions
177         .copy()

```



```

171         print()
172
173     print(f'Best test accuracy: {best_acc:.4f} at epoch {best_epoch+1}
174     ')
175
176     plot_metrics(train_losses, train_accuracies, test_losses,
177     test_accuracies)
178
179     print(f"Creating confusion matrix from epoch {best_epoch+1} (best
180     accuracy)")
181     plot_confusion_matrix(best_epoch_labels, best_epoch_predictions,
182     class_names=['human', 'ai'])
183
184     return model
185
186 def main():
187     print(f"PyTorch version: {torch.__version__}")
188
189     model = ViTClassifier().to(Config.DEVICE)
190
191     criterion = nn.CrossEntropyLoss()
192
193     optimizer = optim.AdamW(filter(lambda p: p.requires_grad, model.
194     parameters()),
195                               lr=Config.LEARNING_RATE, weight_decay=0.01)
196
197     # Create dataloaders
198     train_loader, test_loader = create_dataloaders()
199     dataloaders = {'train': train_loader, 'test': test_loader}
200
201     print(model)
202
203     model = train_model(model, dataloaders, criterion, optimizer,
204     Config.NUM_EPOCHS)
205
206     torch.save(model.state_dict(), 'final_ViT_model.pth')
207     print('Training complete. Models saved.')
208
209 if __name__ == '__main__':
210     main()
211
212 # Interface for classification using ViT model
213 def predict_image(image_path, model_path='best_ViT_model.pth'):
214     # Load the model
215     model = ViTClassifier().to(Config.DEVICE)
216     model.load_state_dict(torch.load(model_path, map_location=Config.
217     DEVICE))
218     model.eval()
219
220     # Load and transform the image
221     image = Image.open(image_path).convert('RGB')
222     image_tensor = test_transform(image).unsqueeze(0).to(Config.DEVICE
223     )
224
225     with torch.no_grad():
226         outputs = model(image_tensor)
227         _, preds = torch.max(outputs, 1)
228         probabilities = torch.nn.functional.softmax(outputs, dim=1)
229
230     class_names = ['human', 'ai']
231     return {
232         'class': class_names[preds.item()],
233         'confidence': probabilities[0][preds.item()].item(),
234         'human_prob': probabilities[0][0].item(),
235         'ai_prob': probabilities[0][1].item(),

```

```
228         'is_human': preds.item() == 0
229     }
```