

Deep Learning for Distinguish AI Generated Images

Chuxuan Yin
Rice University
6100 Main St, Houston, TX 77005
cy48@rice.edu

Abstract—AI painting is an emerging field. Although there is some controversy over how to use it, it is without a doubt one of the most successful applications of generative models in industrial production. How to distinguish whether a image is generated by AI is a topic that is always concerned by people, while it is not fully developed yet. In this paper, we will try to categorize this problem as an image style recognition problem, and use some well-developed models to classify AI generated images. Specifically, we will train ResNet50, fine-tuning GoogleNet and ViT on different kinds of dataset, and compared their performance based on some benchmark. Then we'll combine all the different kinds of images into one dataset and try to train models on it. We'll test our model in this way in order to see whether our model could capture some common features among different categories of AI generated images. The result of our experiments show that current image classification model could really distinguish AI generated images to a certain degree.

I. Introduction

In the last months, AI painting suddenly began to attract people's attention. With the implementation of the stable diffusion framework, AI painting finally showed its great potential on large scale commercial usages. People who know nothing about painting could also easily use this model to generate high quality images in just one minute. While AI painting looked really awesome, it also raised some critiques, one of these is that those generated images are very homogenized, so they are unable to substitute the hand-writing pictures. Although this is a very controversial topic, we could try to think this problem in a different way: do AI generated images really have some features in common? While it might be a hard task for human eyes to distinguish these slight differences, AI itself could be a good classifier. In fact, if we think of AI-generated images as images of a certain style, then this problem is really close to the fine-grained image classification problem which we are familiar with. With the inspiration of this idea, we'll conduct a fine-grained image classification way to solve this problem in this paper.

According to the literature review[1], we know that there are four mainstream methods for fine-grained image classification based on deep learning:

- The first one is based on regular image classification model such as AlexNet, ResNet[2], GoogleNet, etc. Although these models have strong representational capabilities, their performance on fine-grained classification is not very ideal. The common solution is using

some pre-trained weights trained on ImageNet, then fine-tune this model to get a final result.

- The second one is the fine-grained feature learning method. In the paper [3] published in ICCV in 2015, Lin et al. proposed a bilinear convolutional neural network model to achieve a better representation of deep convolutional features. This method uses two networks, VGG-D and VGG-M, as the benchmark network. Without using the Bounding Box (border) label information, it reaches a classification accuracy of 84.1% on the CUB200-2011 dataset; while using the Bounding Box, its classification accuracy is as high as 85.1.
- The third one is object part based detection. The idea of the method based on object part detection is: first detect the position of the target in the image, and then detect the position of the discriminative area in the target, and then combine the target image and the discriminative area the target region blocks are simultaneously fed into a deep convolutional network for classification. One of the representative work is the Part-RCNN method proposed in 2014 ECCV[4].
- The last one is vision attention method. This method is inspired by the vision attention mechanism in human beings, and is widely adopted in computer vision afterwards. Since vision attention model could identify discriminative area in images without labels, it showed a great potential in fine-grained image classification. The recurrent attention convolutional neural network proposed in the 17-year CVPR is a representative work[5].

II. Dataset

Since this is a totally new field, we do not have off-the-shelf dataset, so we need to count on ourselves to collect data. Basically we need two different kinds of images: AI-generated images and hand-drawing images. The former can be generated with models, and the latter could be found on website. But there is a key issue: we need to make sure that AI-generated images and hand-drawing images are basically looked the same for human beings, so that our experiment result would not be affected by other factors.

On the basis of above, we decide to collect three specific genres of images: the avatar of Ganyu(a game character) in manga style, the castle in fantasy tyle, and the sports car

Fig. 2: Castle

III. Models

A. ResNet50

B. GoogleNet

```

graph TD
    Prev[Previous layer] --> C1_1[1x1 convolutions]
    Prev --> C1_2[1x1 convolutions]
    Prev --> C1_3[3x3 max pooling]
    C1_1 --> C2_1[3x3 convolutions]
    C1_2 --> C2_2[5x5 convolutions]
    C1_3 --> C2_3[1x1 convolutions]
    C2_1 --> FC[Filter concatenation]
    C2_2 --> FC
    C2_3 --> FC
  
```

The diagram illustrates a convolutional layer architecture. It starts with a 'Previous layer' (green box) at the bottom. Three arrows lead from it to three parallel processing blocks: '1x1 convolutions' (blue box), '1x1 convolutions' (yellow box), and '3x3 max pooling' (pink box). From the first '1x1 convolutions' block, an arrow points to a '3x3 convolutions' block (blue). From the second '1x1 convolutions' block, an arrow points to a '5x5 convolutions' block (blue). From the '3x3 max pooling' block, an arrow points to a '1x1 convolutions' block (yellow). Finally, arrows from the '3x3 convolutions', '5x5 convolutions', and the second '1x1 convolutions' block all point to a 'Filter concatenation' block (green) at the top.

C. ViT

Vision Transformer (ViT)

The diagram illustrates the ViT architecture. At the bottom, input images are processed by a **Linear Projection of Flattened Patches** layer. Above this, a **Patch + Position Embedding** layer (consisting of learnable class and position embeddings) is added. The resulting sequence is fed into the **Transformer Encoder**. The output of the encoder is passed through an **MLP Head** to produce the final **Class** (e.g., Bird, Ball, Car).

Transformer Encoder (Detailed View):

- The input is **Embedded Patches**.
- The encoder consists of $L \times$ identical layers.
- Each layer contains:
 - Layer Normalization (LN)**
 - Multi-Head Attention** (with residual connection and layer normalization)
 - Feed-Forward MLP** (with residual connection and layer normalization)

IV. Experiments and Results

Due to the limitation of time and computing resources, we only train the model for 300 epochs. This is enough for us to show the tendency of convergence and get a clear

understanding of the ability for each model to classify our dataset. For other hyperparameters, we choose to train models in batch-size 100 to control the stability of our training process. We also use SGD optimizer with learning rate $1e-3$ to improve the training efficiency. Finally we use cross entropy loss function for this classification problem.

A. Training on Dataset of Single Image Category

First we'll look at the performance of ResNet50. Here we display the training loss and accuracy of each model on avatar dataset. From the training loss and accuracy curve we can see that the loss function does not converge very well during the training process. But we can see the tendency that the training loss is decreasing and the accuracy is increasing. The test accuracy of trained ResNet50 model is around 0.7, which means after training, ResNet50 could really identify the AI generated avatar images to a certain degree, but the trembling curve indicates that this model doesn't fit our dataset very well.

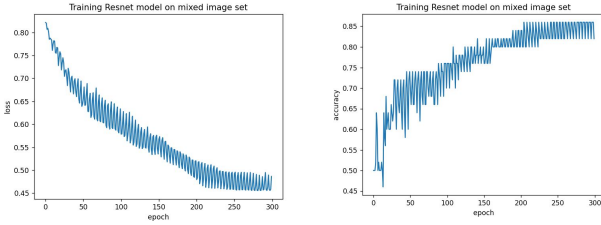


Fig. 7: training performance on avatar dataset for ResNet50

Then we try to train GoogleNet on the same dataset. As we described before, we load the GoogLeNet model with the pre-trained weights on ImageNet, and start to train it on our own dataset. We can see that for the single category image dataset, GoogLeNet model achieved great performance. It's loss function is much more smooth than ResNet50 and it converges very fast. The test accuracy of 0.8 also shows its ability of predicting the AI generated avatar images.

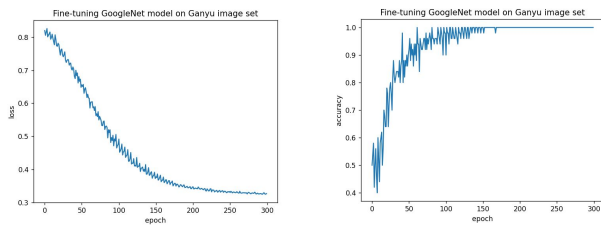


Fig. 8: training performance on avatar dataset for GoogleNet

Finally we turn to see the performance of ViT model. Unfortunately, the ViT model didn't perform well on our dataset. The loss function is not smooth and the accuracy is not increasing. The test accuracy is only 0.6, which

means the model is not able to identify the AI generated avatar images.

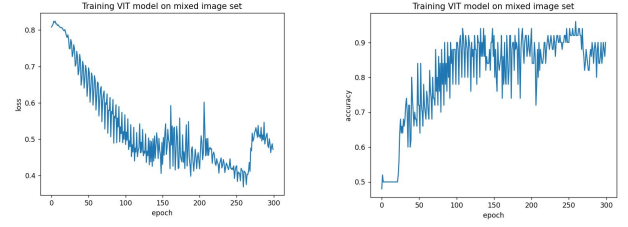


Fig. 9: training performance on avatar dataset for ViT

From the above results, we can see that pre-trained weights do affect the performance of the model. The pre-trained GoogleNet fits very well on our dataset, which means the weight gained from other ImageNet helps a lot in identifying AI generated images. While visual attention mechanism does not work very well on our dataset, this might be explained as we can't recognize the AI generated images by observing the local details.

B. Training on Dataset of Multiple Image Categories

After we train the model on single category image dataset, we try to train the model on the dataset with mixed images. This could show that whether our model could capture some common features of the AI generated images, or our model is just be trained to classify a certain type of images in the previous experiments. Here are the training results of ResNet50, GoogleNet and ViT model on the mixed dataset.

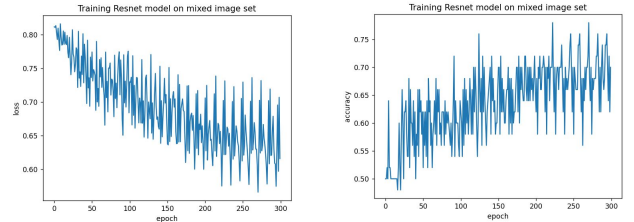


Fig. 10: training performance on mixed dataset for ResNet50

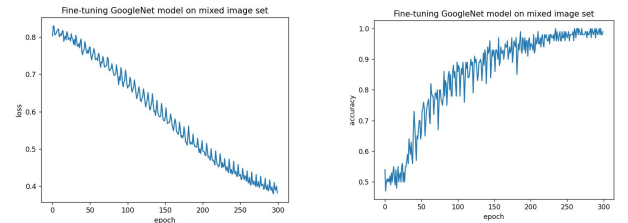


Fig. 11: training performance on mixed dataset for GoogleNet

Unsurprisingly, our models performed much worse on the mixed dataset. ResNet50 and ViT model are not

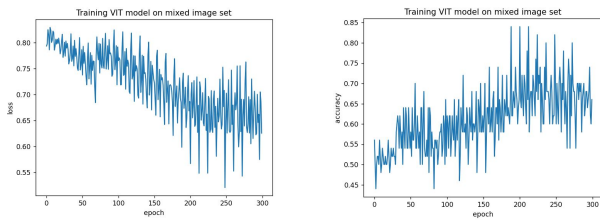


Fig. 12: training performance on mixed dataset for ViT

even converged after 300 epochs of training, and the loss function curve of fine-tuning GoogleNet also trembles a little bit more than the previous experiments. But even in this case, the fine-tuning GoogleNet still showed its great fitness of our dataset. And the test accuracy achieved 0.7, which means our model gives its prediction based on some logic while not just throw a random try.

V. Conclusion

In this paper, we have explored the possibility of using mature fine-grained image classification model to identify AI generated images. We choosed three different kind models to train on the dataset of one category images and the dataset of mixed images. The results show that for current AI generated images, fine-tuning model could classify them with a relatively high accuracy, and it could really caputres some internal features of AI-generated images to a certain degree. This result really help us understand more about AI painting, let us carefully review the difference and connection between AI painting and human painting, which may affect how we use AI painting in the future.

References

- [1] Xiushen Wei, Yizhe Song, et al. Fine-Grained Image Analysis With Deep Learning: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Dec 2022.
- [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun. Deep Residual Learning for Image Recognition. *IEEE Conference on Computer Vision and Pattern Recognition*, June 2016.
- [3] ZHANG, Xiaofan, et al. Embedding label structures for fine-grained feature representation. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, May 2016.
- [4] Zhang N, Donahue J, Girshick R B, et al. Part-Based R-CNNs for Fine-Grained Category Detection. *European Conference on Computer Vision*, Nov 2014
- [5] Volodymyr Mnih, Nicolas Heess, Alex Graves and Koray Kavukcuoglu. Recurrent Models of Visual Attention. *Advances in Neural Information Processing Systems*, June 2014.