# Understanding Photographic Style with Deep Learning

**Jeff T. Sheng**[*]
Department of Sociology
Stanford University
jtsheng@stanford.edu

## Abstract

This project uses deep convolutional neural nets on a unique dataset of over 20,000 photographs by 511 celebrated art photographers from the 20th and 21st centuries to better understand artistic style in art photography. Using a ResNet-18 network with transfer learning in a classification task, we find that deep neural nets are remarkably accurate when used on a smaller dataset of art photographers (~6 artists, 90% accuracy), but are less accurate when the number of photographers increases (~110 artists, 53% accuracy and ~500 artists, 44% accuracy). While this drop in overall accuracy is not surprising, by examining the confusion matrix outputs, we discover valuable insights to the power of deep convolutional neural nets in "learning style". Namely, we see that prediction errors were more likely to be made among artists with overlapping styles rather than those that were only similar in content, while artists whose work have strong visual styles often still had very high accuracy rates. The findings show how deep neural nets are capable of understanding complicated non-linearities involved in aesthetic appreciation and can be powerful tools in helping us understand artistic expression.

## 1   Introduction

In recent years, deep neural nets have become extremely useful in artificial intelligence and machine learning applications due to technical and algorithmic innovations that have allowed researchers to achieve greater accuracy in more complicated tasks such as vision detection and speech recognition. One of these areas where deep learning has been heavily utilized has been in the work on artistic style because the ability to create very deep networks such as ResNets [6] have allowed us to better understand the non-linearities that are highly present in this realm. This paper contributes to the deep learning literature on artistic style by utilizing a unique dataset of art photographers and their work to better understand how deep neural nets can classify and learn based on artistic style.

To the best of our knowledge, this is the first paper in the deep learning literature that examines style in artistic photography, using a dataset of over 20,000 photographs by art photographers whose work is collected by major museums and studied by students in fine art institutions. Utilizing an original dataset assembled, collected and cleaned by the author, we are able to extend and connect some of the previous research done on artist style and deep learning to the medium of artistic photography.

Specifically, we use a ResNet-18 model [6] with transfer learning over a training set consisting of images by various art photographers. Using a softmax classifier, we then predict on a separate test set, the actual artist based on our model's learning. We run various experiments to first build the most accurate classification model, and then we use this version of our model on five variations of our dataset, each containing a different number of artists. From these results, we examine both the accuracy scores and the corresponding confusion matrix to better understand how our model is predicting, and we interpret these findings in the discussion and conclusion.

This paper offers the following contributions: we show that in datasets containing a few artists (6) and a sufficient number of training samples (over 200), our model can achieve test accuracy rates above 90%. Even when the model is expanded to our full dataset of 511 artists and with limited training samples, we achieve an accuracy rate of 44% which is close to an almost even 50/50 chance of being correct despite there being over 500 artists to choose from. In examining the corresponding confusion matrix for each variation, we see that classification errors occur more frequently among artists with similar style rather than over content, and we also observe a few surprises in what the models are able to accurately predict and what they do less well on, verifying some assumptions about how our models are learning.

## 2 Related work

One of the more recent groundbreaking projects in artistic style has been the work of Gatys, Ecker, and Bethge, who show how it is possible to use the different layers of deep neural nets in ways to "transfer" the style of one artist towards another[4][5]. One classic example uses Vincent van Gogh's "The Starry Night" and applies the style of van Gogh to a photographic image [5]. Their algorithmic techniques showed how deep neural nets could learn aesthetic content related to edges, brush strokes, color and visual expression and that these components could be separated from layers and applied to other images. Other scholars have shown how deep learning can be applied to classification tasks over a large dataset of paintings. Applying a ResNet-18 transfer learning model over the Artsy dataset [1], some have found classification test accuracy rates of almost 80% when used on a dataset with over 50 painters [18]. These accuracy rates are much higher than previous methods of painting classification done by SVMs or examining brush stroke features [11][15].

While much of this work has been promising and shows how deep neural nets can be used to understand and classify paintings, this paper extends some of the work done by [17] to see how well it applies to art photography. Previous work on photographic style has often examined photographic datasets of popular photography sites such as Flickr, using style tags applied to the images [9]. These style tags however are not artist specific, namely, they do not assume each photographer has their own unique style, but groups images based on characteristics such as "Vintage", "Romantic" and "Hazy". Other research using the term "photographic style" also apply a generalized notion to this phrase, such as the work done by [14] who categorize the style of photographs as "Sports", "Abstract" or "Landscape." This paper challenges the existing notion of "photographic style" as being a term that only describes generalized categories of photography, to one that is similar in the way we conceive of painters as having a specific style. As such, we aim to synthesize these two bodies of literature on deep learning and artistic style as it has been conceived in the work of [5] and [18] to the medium of art photography.

## 3 Dataset and Features

### 3.1 Overview and Initial Collection and Cleaning

One of the strengths of this project is that it uses an original dataset collected from the archives and teaching materials shared by various photography professors for this research project, including those that have taught photography courses at the Massachusetts Institute of Art, the Rochester Institute of Technology, and Harvard University. One of the interesting aspects of teaching art photography, is that the training of photography students often involves looking at thousands of "great photographs" taken by hundreds of acclaimed art photographers, with the hope that students can absorb the strong aesthetics shown in these images. Because of this, we were able to obtain the teaching materials from various photography professors in the form of a 65GB dataset of over 50,000 images. After examining the raw version of this dataset, we decided to use a

subset of these teaching materials, given the additional time necessary to clean the entire raw dataset. The final cleaned dataset used for the research here is comprised of 511 artists with a total of over 20,000 photographs. Each artist had anywhere from 5 images to over 300, with the mean being around 40 and the median around 25. The author went through each artist folder to examine them for irregularities and to make sure that the format of each image was in jpg form.



**Figure 1:** Photographs from 4 different artists in our dataset. Left to right, top to bottom: August Sander, Man Ray, Walker Evans, and Julia Margaret Cameron. These artists are from the 6-artist experiment where our model achieved almost 91% test accuracy in a classification task.

### 3.2 Preprocessing

After this initial stage of cleaning, our data would be in the form of each artist's name containing all their images from our dataset. Because our implementation utilized the Dataloader class in Pytorch [2], we needed our dataset to be in separate "Train" and "Test" folders. Here, we chose to only split the data into a train and test set as this aligned better for the goals of the project. We preprocessed this step separately using a Python script, and our code randomly split each artist into a 70% "Train" folder and then a 30% "Test" folder, under the artist's name. This set-up allowed the Dataloader to accurately check if the predictions made by our classifier were correct. We also modify the images before passing them into the main training part of our models so that they are most optimal for the ResNet architecture. We first resize each image, zero center them and normalize them. Then, we take a 224x224 crop of each input image and during training, randomly horizontally flip each image with a 50% probability and take a random crop of this, which helps reduce overfit in our training data. In testing, we resize the image and always take a 224x224 crop from the center of the image.

## 4 Methods

### 4.1 Choice of Model Architecture

Researchers have discovered that the number of layers a network architecture has, or how "deep" it goes, is generally very positively correlated with accuracy in image classification tasks [6]. There are two main challenges in increasing the number of layers, one is the problem of vanishing and exploding gradients, and the other is degradation, meaning that accuracy gets saturated and degrades rapidly as the network increases.

One method that has shown to be effective in remedying these issues are residual block architectures which take advantage of skip connections to take the activation from one layer and feed it to another layer much deeper in the network. Figure 2 shows a basic two-layer building block for residual learning. Here, we define $F(x) = W_2\sigma(W_1 x) + x$ where $W_1$ and $W_2$ are the the weights for the convolutional layers and $\sigma$ is the activation function, where we choose a rectified linear unit, or RELU, function. Observe that the operation $F(x) + x$ is realized by the shortcut connection (the $x$ identity skip connection or identity mapping), and we can use element-wise addition followed by another activation function $\sigma$, where our resulting formulation for the residual block is: $y(x) = \sigma(W_2\sigma(W_1 x) + x)$.
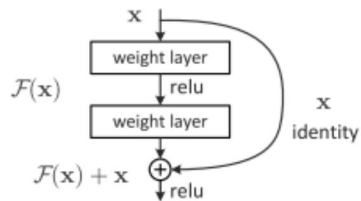


**Figure 2:** Building block for residual learning

This underlying structure is the foundation for the ResNet model by [6], where the usage of such residual blocks and skip connections facilitates the training of much deeper networks without accuracy getting saturated and degrading. Because of their success in image classification tasks, particularly ones that explore high levels of non-linearities and could benefit from many more layers, we chose the ResNet-18 architecture as our main model, which is a ResNet with 18 layers [6].

## 4.2 ResNet-18 with Transfer Learning

The network architecture used is based on a ResNet-18 architecture that starts with pre-trained weights from ImageNet [14]. On the right, we see a diagram of the ResNet-18 architecture based on [6] and [13]. The final fully-connected layer is replaced with a new layer to calculate a score for each artist in our dataset instead of a score for ImageNet classes. We use a softmax classifier with cross-entropy loss:

$$L_i = -log(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}})$$

$L_i$ is the loss for example $i$ in the training minibatch (size 32), $f$ is the score for a particular class calculated by the network, $j$ is one of the possible classes, which depends on the number of artists in the variation of the dataset we were using. This loss function is optimal for ensuring that our network maximizes the score of the correct artist in the training examples relative to other artists, which allows us to see how the model accuracy varies when we increase the size of the number of artists in our various experiments. For training, we first held the weights of our pre-trained network constant for 10 epochs, and then allowed all the weights throughout the entire network to update, training for an additional 10 epochs. We experimented with different regularization techniques including dropout [17], L2 regularization (weight decay), as well as changing the learning rate in the Adam Optimizer [10]. Additional experimentation was done where we allowed the weights in one or more layers to update in the first ten epochs (i.e. the last layer before the fully connected layer), as well as not updating all the weights in the last 10 epochs. These results consistently performed worse than our original model that relied on transfer learning with fine tuning. Additionally, the pre-trained weights on ImageNet performed the best as opposed to a ResNet-18 model trained from scratch. In the results section, we present the best models which are the ones where the weights are pre-trained on Image Net and then fine tuning was achieved where all the weights were allowed to update.



**Figure 3:** The ResNet-18 architecture based on [6]

Accuracy charts were created after the first 10 epochs, and a second accuracy chart was made after the additional 10 epochs to show the effects of fine tuning, and our final accuracy scores reflect this version of our model architecture.
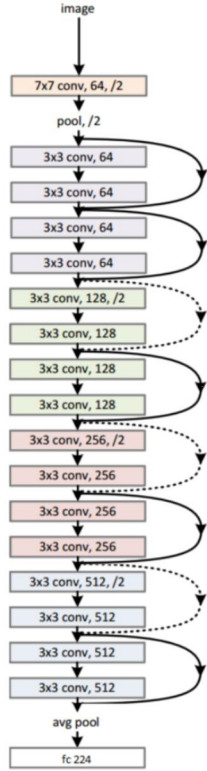
# 5 Experiments/Results/Discussion

## 5.1 Experiment Setup, Implementation Details, and Evaluation Metrics

All of our coding implementation for training and testing was done in Pytorch [12] and used the Dataloader, Dataset, and Imagefolder Classes [2][3][7], and the visualization of our results adapted a confusion matrix presentation from Sci-kit Learn [16]. Original code was written in Python to pre-process the dataset so that it could be used accurately by the Dataloader implementation [2]. The ResNet-18 architecture and the pretrained weigths from ImageNet were obtained from [14][13], and we used [8] as a template for preforming transfer learning by replacing the fully connected layer of the pre-trained ResNet-18 network. Our best models used Dropout [17] with a 80% zero probability, and we used an Adam Optimizer based on [10], with a learning rate of $10^{-3}$ and weight decay value of $10^{-3}$. Experiments were performed on GPU's, specifically an Amazon EC2 p2.xlarge instance and a p2.8xlarge instance donated by AWS.
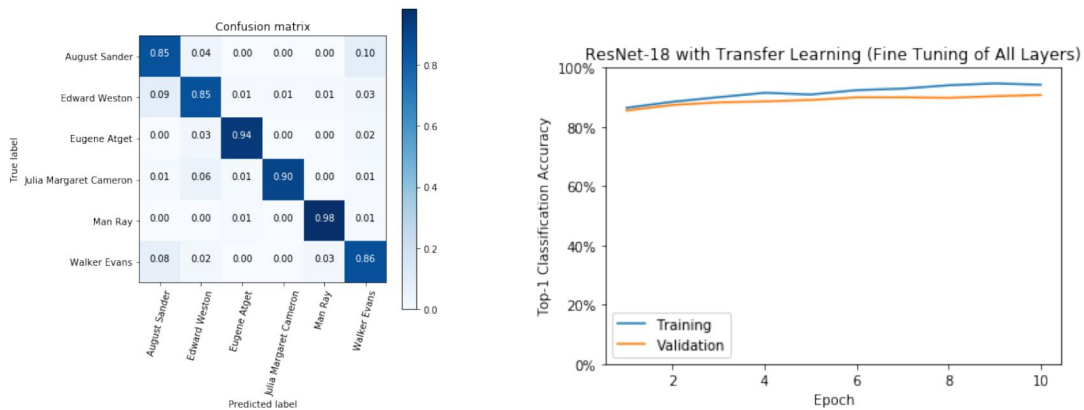
In our results, we report our top train and test accuracy scores (we multiply them by 100 to show the percentage). Our confusion matrix has our y-axis showing the true label and the x-axis showing the predicted label, and the corresponding number is a normalized score.

4

## 5.2 Results

The following table shows our best train and test accuracy scores in the 5 variations of the dataset that we experimented with, as well as noting the sample size range for each artist.

| Best Accuracy Scores (Train and Test) | | | |
|---|---|---|---|
| Number of Artists in Model | Each Artist Sample Size | Best Train (70% of sample | Best Test (30% of sample) |
| Artists: 511 | > 5 Photos | 73.54 | 44.04 |
| Artists: 110 | > 40 Photos | 75.01 | 53.12 |
| Artists: 52 | > 70 Photos | 69.02 | 55.46 |
| Artists: 31 | > 100 Photos | 91.88 | 68.30 |
| Artists: 6 | > 250 Photos | 94.21 | 90.76 |

We also show the corresponding confusion matrix and accuracy plots for each of the various experiments (see appendix). Here we present the confusion matrix and final accuracy chart for the 6 artist variation.



(Note that the number of epochs shown in the above right chart are the number after an initial 10 epochs of training on the pre-trained ImageNet weights [14]. Also, due to space limitations here, we provide additional analyses and illustrations of more cases in the appendix.)

Quantitatively, we see that accuracy increases as the number of artists decreases, and at very low numbers, there is very high accuracy, at over 90%. However, we also see that there is still a decent level of overfit in all the models that have more than 6 artists. Many different regularization techniques were tried to decrease the gap between the training accuracy and the test accuracy, including a very high Dropout rate (80% zero-probability), as well as a fairly large weight decay amount of $10^{-3}$. Despite best efforts to address overfitting, our results with the 6-artist sample give us confidence that increasing the number of training examples to over 250 in future studies should help address some of the overfit issues.

Qualitatively, this is where our results are perhaps even more informative. In a thorough examination of the various confusion matrix outputs, a few trends emerged that highlighted the power of neural nets in learning artistic style. The three main findings include:

1. Certain photographers whose work would be categorized by human experts as having a distinct and strong style, enjoyed very high accuracy rates despite the increase in the number of artists in the dataset. For example, the Dutch artist, Rineke Dijkstra is well known in the art world for having a very distinct and formal style. When classified in the 31-artist experiment, the model predicted her work at a 93% accuracy rate. However, in the 110-artist comparison, the model still had an 88% accuracy rate. What is striking about this observation, is that her work is characterized by a distinctive style where most human experts would be able to recognize one of her images based on aesthetics.

2. When more artists were compared against each other, the errors made by the model often happened in cases where style was similar, not just because the artists had the same content. For example, in the larger 110-artist experiment, the artists Sam Taylor-Wood and Philip-Lorca DiCorcia were confused at a 0.24 normalized error rate, even though when there

5

are fewer artists in the dataset, the neural nets do a relatively good job at predicting each artist individually. This is interesting because both artists come from the same genre of photography called "Documentary Fiction," a style known for expressive moods and emotion, using carefully constructed lighting, actors, and staged imagery to construct a fictionalized image rather than one found in regular documentary photography, and the two artists are often grouped together in stylistic genre by museum curators and other experts.

3. There were unique cases where a human expert would have no trouble in classifying a photograph but the model struggled more due to the lack of stylistic information in the image. The best illustration of this is the artist Cindy Sherman who is well known for her self-portraits that mimic the aesthetic style of other artists. Even though Cindy Sherman herself is in each of the photographs - a human being who knows what she looks like, can easily tell if a photograph is her work. However, in her case, the deep neural nets were often confused and only predicted her work at a 50% accuracy rate despite herself being in every one of her images.

In sum, an extensive qualitative analysis show that deep neural nets utilize a high degree of understanding style in their models. Artists who are well known for having strong distinctive artistic styles have high rates of accuracy regardless of the number of other artists they are being compared against. Some artists who overlap heavily within sub-genres of style and would be curated by art experts as being stylistically very similar, are often confused by our models, especially when there are more artists. And finally, artists who are known for mimicking the style of other artists confuse the neural nets, despite cues that human experts may perceive (such as the same person appearing in all of the images).

# 6  Conclusion/Future Work

Future work on this topic will include experimenting with other architectures such as Triplet Loss and other CNN architectures that can be applied to understanding art photography in different ways than just classification. I also plan on expanding and cleaning more data to add to the current dataset to see if there can be other more robust findings with more data. As previously mentioned, increasing the number of samples for each artist may also help with overfit in some of the models.

This paper is the first of a number of projects using this dataset. We hope that our findings here and future explorations of this topic can have a significant impact towards the understanding of artistic style with deep learning, contribute to the literature in the sociology of culture, and also used by museum curators and art students to better see the connections (and accuracy) between classification done by human "experts" and those now achievable through deep learning techniques. In conclusion, this work shows that deep neural nets are powerful tools in helping us realize how art and artistic expression contain a high degree of non-linearity and subtleness in its creation, appreciation and knowledge, where we can use deep learning to better understand artistic style in photography.

# 7  Contributions and Code

I worked alone on the project which involved assembling and cleaning the dataset, writing code for all the preprocessing of images, and writing/adapting code for the models and the visualizations of outputs. The final best models can be downloaded from this Dropbox link as Jupyter Notebooks showing the final outputs and graphs presented here, as well as the code written to pre-process the images: `https://www.dropbox.com/sh/xt0aowjubhb5tqc/AAAZ-jL52ThGbJhy94gvR5_Aa?dl=0`

# References

[1] Artsy. https://www.artsy.net/about/the-art-genome-project.

[2] Dataloader Class, Pytorch. https://github.com/pytorch/pytorch/blob/master/torch/utils/data/dataloader.py#L262

[3] Dataset Class, Pytorch. https://github.com/pytorch/pytorch/blob/master/torch/utils/data/dataset.py

[4] Gatys, L.A., Ecker, A.S. & Bethge, M., (2015) A neural algorithm of artistic style. arXiv preprint arXiv:1508.06576.

[5] Gatys, L.A., Ecker, A.S. & Bethge, M., (2016) Image style transfer using convolutional neural networks. In Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on (pp. 2414-2423). IEEE.

[6] He, K., Zhang, X., Ren, S. and Sun, J., 2016. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 770-778).

[7] ImageFolder Class. Pytorch. https://github.com/pytorch/vision/blob/master/torchvision/datasets/folder.py

[8] Johnson, J. Neural-style. https://github.com/ jcjohnson/neural-style, 2015.

[9] Karayev, S., Trentacoste, M., Han, H., Agarwala, A., Darrell, T., Hertzmann, A. and Winnemoeller, H., (2013) Recognizing image style. arXiv preprint arXiv:1311.3715.

[10] Kingma, D.P. and Ba, J., 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.

[11] Li, J., Yao, L., Hendriks, E. and Wang, J.Z., (2012) Rhythmic brushstrokes distinguish van Gogh from his contemporaries: findings via automated brushstroke extraction. IEEE transactions on pattern analysis and machine intelligence, 34(6), pp.1159-1176.

[12] PyTorch. https://github.com/pytorch.

[13] Resnet Class for Pytorch. https://github.com/pytorch/vision/blob/master/torchvision/models/resnet.py#L111*

[14] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M. and Berg, A.C., 2015. Imagenet large scale visual recognition challenge. International Journal of Computer Vision, 115(3), pp.211-252.

[15] Saleh, B. and Elgammal, A., 2015. Large-scale classification of fine-art paintings: Learning the right metric on the right feature. arXiv preprint arXiv:1505.00855.

[16] Sci-Kit Learn Confusion Matrix. http://scikit-learn.org/stable/modules/generated/

[17] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. and Salakhutdinov, R., (2014) Dropout: A simple way to prevent neural networks from overfitting. The Journal of Machine Learning Research, 15(1), pp.1929-1958.

[18] Viswanathan, N. (2017) Artist Identification with Convolutional Neural Networks, CS231 Final Projects, Stanford University.

# 8 Appendix and Additional Analysis

In this section, I provide additional analysis, figures and charts that could not fit in the main body of the suggested paper length.

## 8.1 Additional Qualitative Analysis and Images

Here are images from the qualitative highlights from the Results section 5.2, 1, 2, and 3, showing visually what is described there.



**Figures for 5.2.1**: Three different photographs by the Dutch artist, Rineke Dijkstra, whose work is

reknown in the art world for having a very distinct and formal style. CNNs were able to predict her work with extremely high accuracy in all the various experiments done classifying her work. For example, the 31-aritst model predicted her work at a 93% accuracy rate, while the 110-artist model predicted her work at an 88% rate.
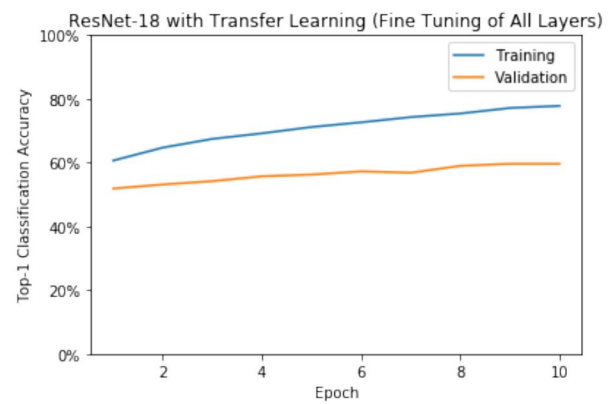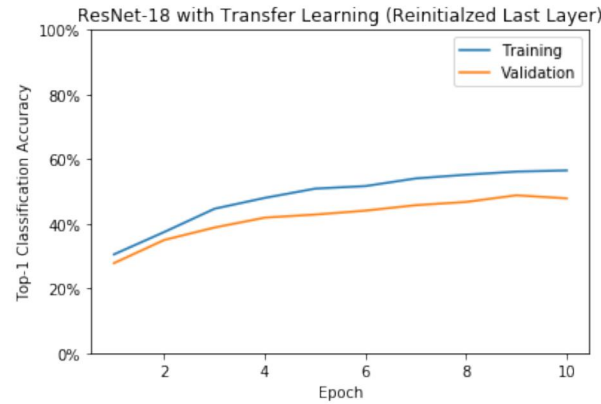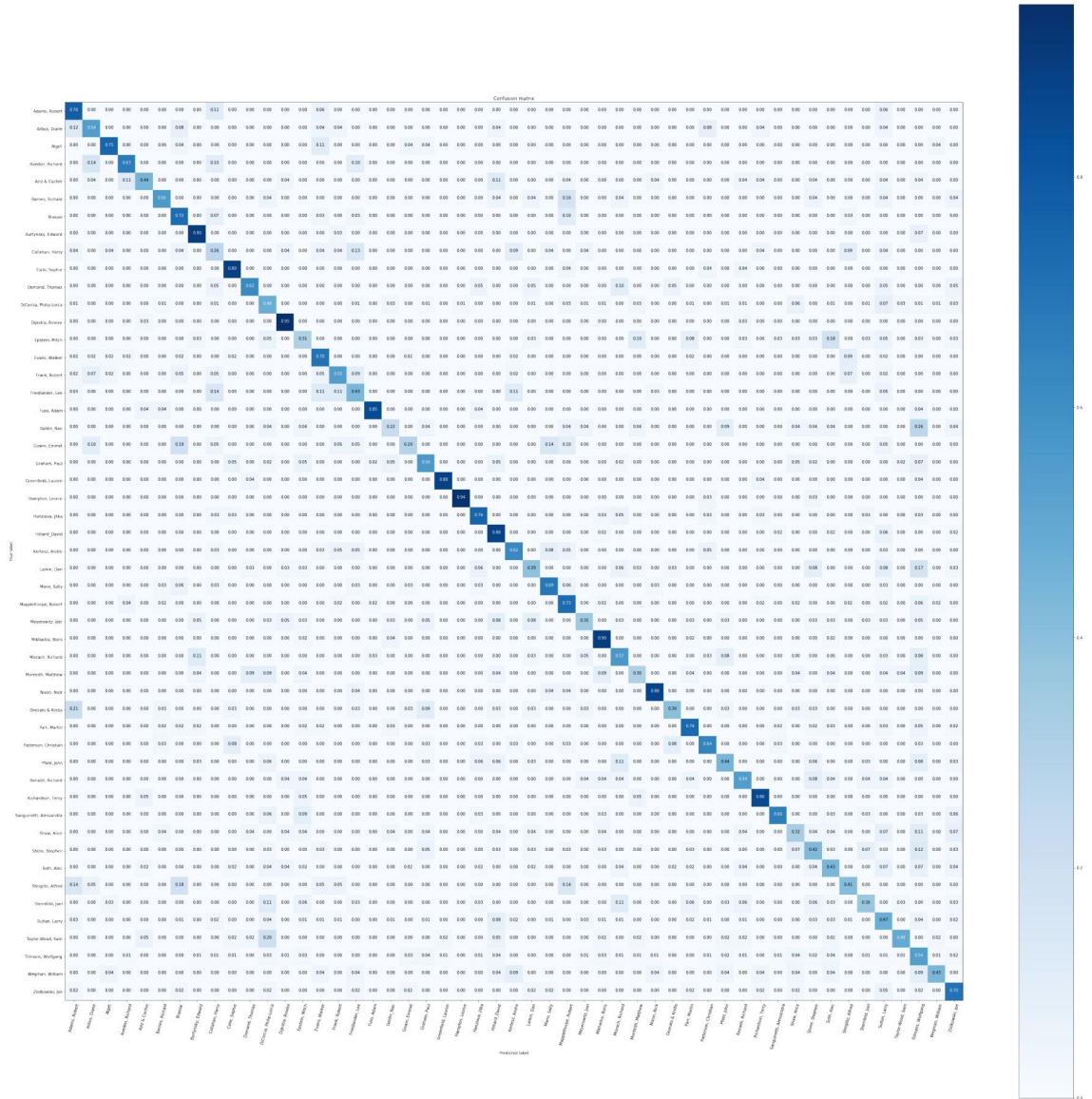


**Figures for 5.2.2**: The above left photograph is by the artist Sam Taylor-Wood while the above right image is by Philip-Lorca DiCorcia. When there are fewer artists in the dataset, the neural nets do a relatively good job at predicting each artist individually. However, as the dataset grows, the model begins to confuse these two artists more (0.24 normalized error in the 110 artist dataset). This is noteworthy because both artists come from the same genre of photography called "Documentary Fiction," a style known for expressive moods and emotion, using carefully constructed lighting, actors, and staged imagery to construct a fictionalized image rather than one found in regular documentary photography, and the two artists are often grouped together in stylistic genre by museum curators and other experts.
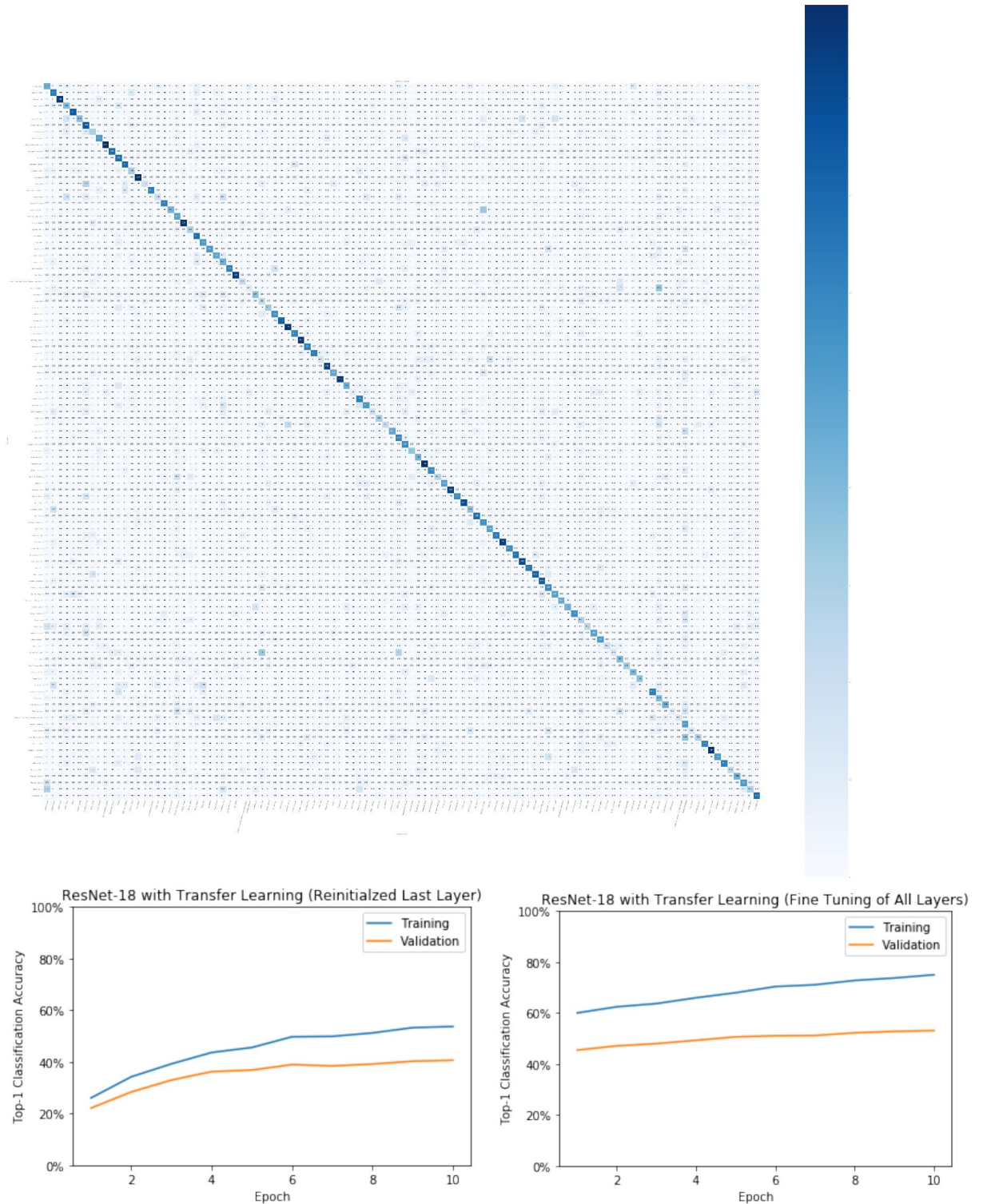


**Figures for 5.2.3**: The above 5 images are all from the artist Cindy Sherman, who mimics the aesthetics and style of various artists and photography genres. However, she is the subject of all her photographs which makes it easy for a museum curator to spot her images (as long as they recognize her and know about her work ). Our model in the 110-artist experiment only predicted her work at a 50% accuracy rate.

## 8.2 Confusion Matrix and Accuracy Charts from 31 Artist Experiment





Note that these two accuracy charts represent a total of 20 epochs. The left chart represents the first 10 epochs with just the pre-trained weights from ImageNet, while the chart on the right represents an additional 10 epochs where we fine tuned the model by allowing all the weights to update.

## 8.3 Confusion Matrix and Accuracy Charts from 52 Artist Experiment





Note that these two accuracy charts represent a total of 20 epochs. The left chart represents the first 10 epochs with just the pre-trained weights from ImageNet, while the chart on the right represents an additional 10 epochs where we fine tuned the model by allowing all the weights to update.

## 8.4 Confusion Matrix and Accuracy Charts from 110 Artist Experiment





Note that these two accuracy charts represent a total of 20 epochs. The left chart represents the first 10 epochs with just the pre-trained weights from ImageNet, while the chart on the right represents an additional 10 epochs where we fine tuned the model by allowing all the weights to update.