

# EE8601 Project 2: Image Style Transfer

Implementation of *Image Style Transfer Using Convolutional Neural Networks* by Gatys et al

Alegre, Edwin Philippe, EIT

Department of Electrical, Computer, and Biomedical Engineering

Ryerson University

Toronto, Ontario, Canada

edwin.p.alegre@ryerson.ca — 500552920

**Abstract**—The aim of this report is to document the procedures, as well as the theory gained from the implementation of Gatys' paper - *Image Style Transfer Using Convolutional Neural Networks*. This paper outlines the motivation, a brief top level theoretical outline on the workings of the above mentioned paper, the implementation (including dependencies, code, etc), as well as outcomes related to the directed studies course, EE8601 - Fundamentals of Computer Vision and Deep Learning.

**Index Terms**—image processing, convolutional neural networks, style transfer learning

## I. INTRODUCTION

The visual arts has been a medium that conveys ones sense of wonder and abstract viewpoint of the universe to its patrons. From Vincent Van Gogh, to Picasso, to the Group of Seven, an artist wields their brush in their own manner, showcasing their own personal genius in the use of colours, strokes, and realism to convey how they see a certain subject. It is no wonder that certain artists have become synonymous with their own personal style. Van Gogh is known for his vibrant colours ad swirled strokes that are evident in his great works, with "Starry Night" as the perfect example. Picasso brings out the abstract use of shapes in his work. Monet captures the vibrant realism of his landscape masterpieces. Artists have tried to copy that style and take years to reach a level that was achieved by these masters. Recent technological advancements have allowed for the use of *convolutional neural networks* (CNN) to achieve great feats in the field of computer vision and engineering. An example of image style transfer can be seen in Figure 1.

This paper will focus on the implementation of *Image Style Transfer Using Convolutional Neural Networks* by Gatys et al. This paper utilizes convolutional neural networks to extract the *style* of one image (the way an artist paints) to an image with a different *content*. In the world of computer vision this is known as *style transfer*. This is done with the use of CNNs that extract the style of one image and combines it with another.



Figure 1. An example of image style transfer. The top left photo showcases Stanford University's Tower and is known as the *content image*. The top right photo is that of Van Gogh's Starry Night and is known as the *style image*. By applying the image transfer technique, a synthesized image of the two input images is formed and shown at the bottom

## II. DEPENDENCIES

In order to successfully implement this project, the following technical dependencies are needed.

### A. Python (ver. 3.7.7, 64-bit)

Python is the primary prototyping language used for machine learning, computer vision, and image processing applications. As such, it is the programming language used to implement this project.

It is important to note that Python will be supplemented by Tensorflow and Keras (both tools that correspond with easily implementing machine learning models within the Python environment).

It is highly recommended that Python is supplemented with *pip*, a Python package manager that makes it easier to install libraries to use within Python. An ideal IDE that the reader is comfortable with using is also suggested.

## B. Tensorflow 2.0

Tensorflow is an open source machine learning library, created by the Google Brain Team, that is capable of interfacing with Python. Tensorflow is widely used in academia and the industry, which makes it a viable platform to use for this project. Specifically, this model relies on a CNN model known as VGG19 which needs to have been trained using a specific data set known as *ImageNet*. Luckily, Tensorflow has already implemented this exact model and set for use.

Tensorflow is capable of training and running deep learning models. The most recommended API to use when dealing with neural networks on Tensorflow is called Keras. It is worth noting that Tensorflow may be used locally (on the host machine) or with a cloud computing service such as Amazon Web Services (AWS) in order to run the models utilized. While it is entirely possible to run Tensorflow locally, it may be time consuming as the vanilla Tensorflow package utilizes the host CPU in order to process the model. This method takes approximately 1 to 2 hours to run the code. The time used to run the model can be optimized by the use of cloud computing or through outsourcing the processing to a local machine's GPU. Typically, GPU outsourcing is achieved through the use of CUDA.

## C. CUDA (optional)

CUDA is a parallel computing platform created by NVIDIA that enables the use of a local machine's GPU for general purpose computing. Specifically, the use of CUDA in this project allows our model to return the results at a much quicker rate due to the use of the local machine's dedicated GPU instead of its CPU. CUDA relies on a variety of dependencies in order to interface with Tensorflow. The following dependencies must be met before CUDA can be utilized:

- NVIDIA GPU card with CUDA Compute Capability of 3.5 or higher
- NVIDIA GPU Drivers - CUDA 10.1 requires 418.x or higher
- CUDA Toolkit - Tensorflow 2.1.0 requires CUDA 10.1
- CUPTI - Installed along with the CUDA Toolkit
- cuDNN SDK

## D. Google Colab (optional)

Google Colab is a notebook service that allows for the creation and execution of Python code within the Google interface. By uploading Python code to a linked Google Drive, code can be run online. Furthermore, Google Colab allows for free GPU utilization, which is extremely useful if CUDA is not a viable solution due to the hardware constraints of the local machine.

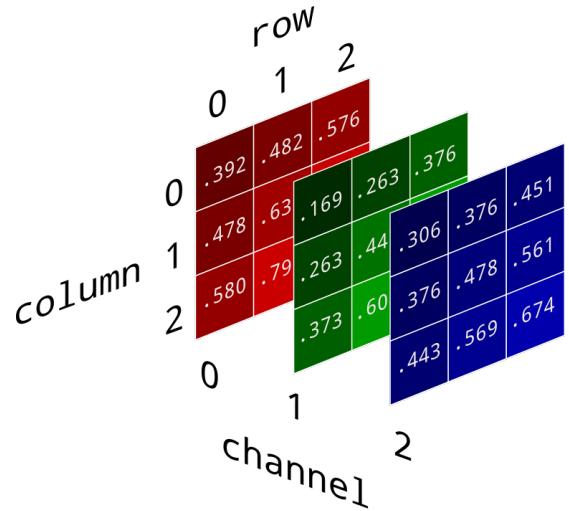


Figure 2. Visual representation of an RGB image

## III. BACKGROUND AND RELATED WORK

In order to properly understand how Gatys' algorithm works, a background in various concepts in Image Processing and Deep Learning is needed. The main concepts are discussed below.

### A. Image Representation

In the realm of engineering, images are defined as a matrix of pixel values of size  $M \times N$ . Typically, the pixel intensity is represented by 1 byte with a range of [0, 255]. Greyscale images are restricted to this range and are represented by a single matrix. Color images are typically represented as a 3 channel matrix that corresponds to red, green, and blue, which is why colour images are known as RGB images of size  $M \times N \times 3$ . An example of an RGB image is show in Figure 2.

### B. The Convolution Operation

The *Convolution* Operation is the backbone on which Convolutional Neural Networks are built on. Convolution is an elementary operation in the field of electrical engineering, specifically in signal processing and image processing. Convolution works by sliding a *filter* or a *kernel* of a pre-specified dimension (typically 3x3) over a matrix. Each element in the matrix is multiplied by the overlaying kernel value and all of the products are summed to produce a single output. Typically, convolution involves flipping the matrix, especially in the field of signal processing. However, in the case of computer vision and CNNs, the matrix remains unflipped. This specific case of convolution is known as *cross correlation*.

Specifically, convolution in CNNs is used to extract features from the input image. By using small filters of input data, convolution maintains the relationship between pixels by

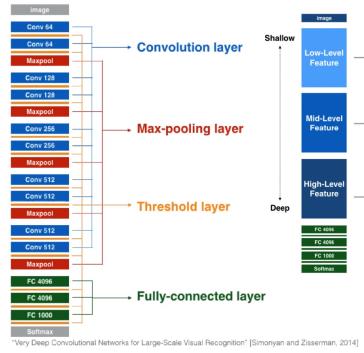


Figure 3. A generic CNN based off the VGG19 network that shows all the typical layers used. The smaller figure on the right side displays what types of features are captured at each convolutional layer in the CNN

identifying image features. The resultant matrix formed by the convolution within a convolutional layer is known as a *feature map*. This feature map is a 3D matrix of size  $N_{\text{width}} \times N_{\text{height}} \times N_{\text{convolution}}$  where  $N_{\text{width}}$  is the width of the input image,  $N_{\text{height}}$  is the height of the input image, and  $N_{\text{convolution}}$  is the number of filters in the layer.

### C. Zero-Padding

*Zero Padding* is another concept widely used in signal and image processing. Zero padding surrounds the input data by a border of zeros such that the borders of the input image can also be captured int eh feature map during the convolution. This also allows for preservation of the dimensions when the input data is convolved with a filter.

### D. Convolutional Neural Network, CNN

The *Convolutional Neural Network* is broken down into a smaller set of functions of made up of layers, where each one of those layers are a neural network in of itself. Specifically, the CNN is comprised of 3 main layers, those being: *the convolution layer*, *the pooling layer*, and *the threshold/non linear layer*. Depending on the CNN architecture, there may be more layers involved in the CNN, but these 3 layers remain constant. The CNN works extremely well for image based problems. Specifically, the architecture of the CNN allows it to learn a variety of features, from low-level features (lines, edges, etc) to mid-level features (combined lower level features that showcase patterns) high-level features (truly sophisticated and complex features/textures at the top level). The initial layers usually learn the low-level features, with each progressive layer learning higher-level features. An example of a CNN is shown in Figure 3, specifically showcasing the CNN used in this project, the VGG19 network. THe VGG19 network contains 3 fully connected layers and 1 softmax layer, which are not utilized in this project. Hence it will not be discussed, however each utilized layer is explained in more detail below.

1) *The Convolution Layer:* The convolution layer is the most important layer in the CNN and is the name sake of

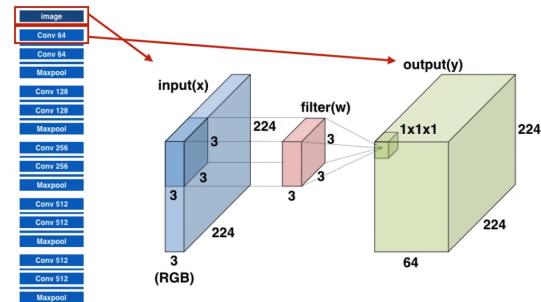


Figure 4. A graphical example of how a convolutional layer in a CNN works. Note that the actual convolutional layer is comprised of a tensor with size 224x224x64

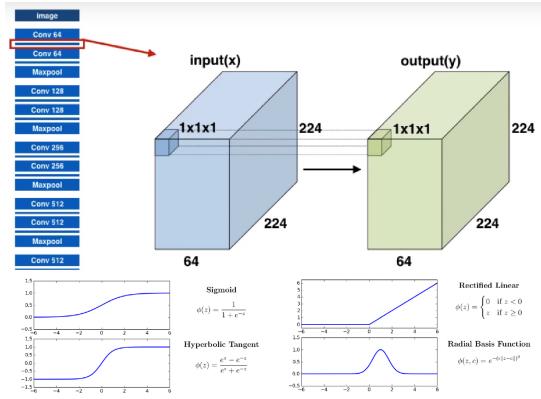


Figure 5. A graphical representation of the threshold layer in a CNN. The threshold is based off the activation function used in the network. The 4 most common activation functions are shown at the bottom of the figure

the CNN. In general with the case of images as the input, the convolution layer is a 3D convolution. The input image is convolved with a kernel where the output is the basis for the feature map. In the case of the VGG19, the first layer generates as 224x224x64 feature map, which is due to 64 different kernels being used to produce the different depths of this layer. The first stage convolution layer of the VGG19 network is shown in Figure 4.

2) *The Threshold/Non Linear Layer:* Following the Convolution Layer is the *threshold layer*. This layer introduces a non linear component to the network, where the most used activation function is normally the *rectified linear unit (ReLU)* but other functions such as the hyperbolic tangent can also be used. This layer is the same size as the previous feature map and each element corresponds to the same element space of the feature map before it. This can be seen in Figure 5

3) *The Pooling Layer:* The *pooling layer* is essentially the layer that downsamples its input layer in space. For example, the first pooling layer in the VGG19 network takes an input from the threshold layer of size 224x224x64 and applies *max pooling* onto it. Usually, pooling is done by sliding over the input layer with a 2x2 kernel of *stride* 2, meaning it moves two pixels over instead of the usual one.

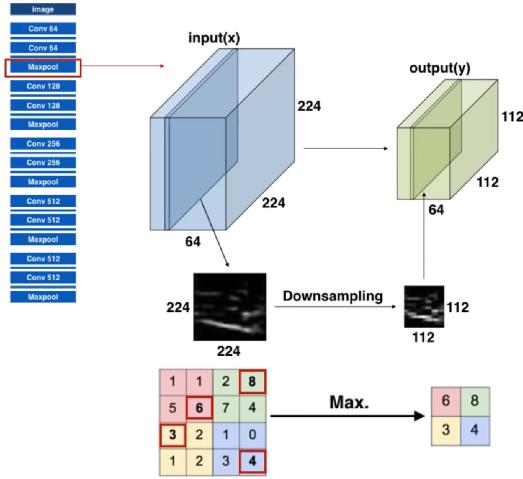


Figure 6. A graphical representation of the pooling layer in a CNN. Note that this layer essentially downsamples the input feature map. The figure at the bottom shows how max pooling is achieved, where a 2x2 kernel with a stride of 2 slides over the input matrix and translates the maximum value in the viewing window as the new downsampled pixel in the resultant pooled feature map.

This kernel performs a max function over its window, choosing the maximum element to represent that feature area in the downsampled feature map. There exists other pooling methods such as the *average pooling* or *sum pooling*. Gatys states that average pooling yields better results, which is why it is utilized as the pooling method for this implementation. As specified earlier, the first pooling layer in VGG19 takes a 224x224x64 input feature map and returns a 112x112x64 downsampled feature map. In all cases, pooling helps make representations become approximately invariant to small translations to the input. Invariance to translation causes the values of most of the pooled output values to stay the same despite the input being translated by a small amount. Invariance to translation is particularly useful if the presence of the feature is more important than the location of it. A sample of the pooling layer can be seen in Figure 6

#### E. VGG19 Network

The Visual Geometry Group Network (VGG) is the winning submission of the 2014 ImageNet challenge. This network has two versions, the VGG16 and VGG19 networks. The number symbolizes the amount of layers in the network, with VGG16 having 16 layers and VGG19 having 19 layers. VGG19 is used for this application and consists of 16 convolutional layers and a rectified linear unit (ReLU) as its non-linear threshold function. The convolutional layers are separated by 5 pooling layers and ends with 3 fully connected layers. As stated previously, the fully connected networks will not be used for this project. Specifically, the content representation will be matched using 'conv4\_2', and the style representation will be matched using 'conv1\_1', 'conv2\_1', 'conv3\_1', 'conv4\_1', and 'conv5\_1'. Each layer will have a weight of

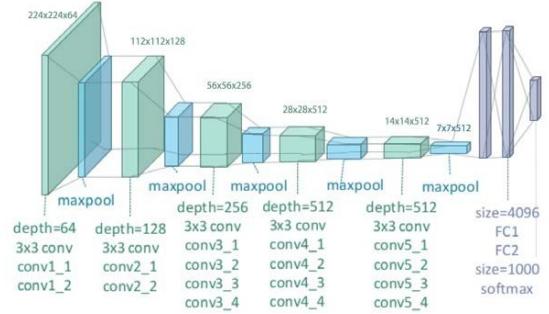


Figure 7. The VGG19 Network. Notice the last 3 layers of the network, known as the fully connected layers. While these layers are part of the VGG19 network, they will not be utilized in this implementation

1/5. The graphical representation of the VGG19 network can be seen in Figure 7

#### IV. DEEP IMAGE REPRESENTATION

The main takeaway from this paper is that the **representations of the content and style are separable and thus, both representations can be manipulated independently to produce an output of the user's liking**. Using this finding, it was determined that the overall loss function is a linear combination of the two individual loss functions. Thus, by regulating the individual loss functions, the emphasis on the resultant image can be based on either the reconstruction of the content or style input. This will be discussed in further detail in the upcoming sections.

##### A. Content Representation

Given an input image  $\vec{x}$ , this input image is passed into each layer of the CNN, providing a feature representation at each layer,  $F^l$ . Gradient descent can be applied to a white noise image in order to visualize the image information that is encoded at every layer of the network. This provides another image that can be used to match the feature responses of the original image. These responses in a layer,  $l$ , can be stored in a matrix with the generalized form of  $X^l \in \mathcal{R}^{N_l \times M_l}$ , where each layer contains  $N_l$  distinctive filters and thus, has  $N_l$  each of size  $M_l$ , where  $M_l$  is the height times the width of the feature map. From this, the activation of the  $i^{\text{th}}$  filter at position  $j$  in layer  $l$  can be symbolized as  $X_{ij}^l$ .

These two images can be called the *content image*,  $\vec{p}$  with feature representation (content features)  $\vec{P}$  generated in Layer 4\_1, and the *generated image*,  $\vec{x}$  with feature representation (content features)  $\vec{F}$  generated in the same layer.

Thus, the squared-error loss between the two feature representations can be defined as

$$\mathcal{L}_{\text{content}}(\vec{p}, \vec{x}, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2$$

From the derivative of this loss function, defined below, the gradient descent with respect to the generated image,  $\vec{x}$ , can be computed using standard back-propagation.

$$\frac{\partial \mathcal{L}_{\text{content}}}{\partial F_{ij}^l} = \begin{cases} (F^l - P^l)_{ij}^2 & \text{if } F_{ij}^l > 0 \\ 0 & \text{if } F_{ij}^l < 0 \end{cases}$$

As stated in earlier sections, CNNs are invariant to translation but are increasingly sensitive to the actual *content* of the image. To reiterate, the lower layers capture the pixel based relationships of the features, such as line and edges, whereas the higher layers capture the actual content of the image. It is for this reason that the feature responses of the higher layers are referred to as the *content representation*.

### B. Style Representation

To generate the *style representation* of an image, a feature space, which can be built on top of the filter responses in any layer, is used to obtain the texture information from this image. Specifically, this feature space contains the feature correlations between different filter responses. These correlations are computed through the *Gram matrix*,  $G^l \in \mathcal{R}^{N_l \times N_l}$ , where  $G_{ij}^l$  is the inner product between the vectorised feature maps  $i$  and  $j$  in the layer  $l$ :

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l$$

Just like the content representation, this gradient descent can be used on a random white noise image to minimise the mean-squared error to achieve a similar visualization. However, this optimization utilizes the Gram matrices from the original image  $\vec{d}$  and the generated image  $\vec{x}$ . The respective style representations in layer  $l$  would then be  $A^l$  and  $G^l$ . The contribution of layer  $l$  to the total loss would then be

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2$$

Thus, the total style loss would then be:

$$\mathcal{L}_{\text{style}}(\vec{d}, \vec{x}) = \sum_{l=0}^L w_l E_l$$

where  $w_l$  are the weights of the contribution of each layer to the total loss. Since there are 5 layers used for the style transfer algorithm, each layer is weighted at  $w_l = \frac{1}{5}$ .

Once again, similar to content representation, the derivative of this mean-squared error, defined below, can be taken with respect to the activation in layer  $l$ . From this, the gradient of  $E_l$  with respect to the generated image  $\vec{x}$  can be computer using standard back-propagation

$$\frac{\partial E_l}{\partial F_{ij}^l} = \begin{cases} \frac{1}{N_l^2 M_l^2} ((F^l)^T (G^l - A^l))_{ij} & \text{if } F_{ij}^l > 0 \\ 0 & \text{if } F_{ij}^l < 0 \end{cases}$$

### C. Style Transfer

The concept of *style transfer* involves the synthesizing of a new image that simultaneously matches the content representation of a photograph  $\vec{p}$  and the style representation of  $\vec{d}$ . The mathematical problem of style transfer essentially boils down to an optimization problem. Specifically, it is the joint minimization of a white noise image's,  $\vec{x}$ , feature representation distance from the photograph's content representation in one layer and the artwork's style representation on a number of layers. The total cost function to me minimized is:

$$\mathcal{L}_{\text{total}}(\vec{p}, \vec{d}, \vec{x}) = \alpha \mathcal{L}_{\text{content}}(\vec{p}, \vec{x}) + \beta \mathcal{L}_{\text{style}}(\vec{d}, \vec{x})$$

where  $\alpha$  and  $\beta$  are weighing factors for the respective content and style representations. The gradient with respect to the generated image  $\frac{\partial \mathcal{L}_{\text{total}}}{\partial \vec{x}}$  can be used for a numerical optimisation strategy. Gatys et al utilized the L-BFGS optimizer, however, this is not available in Tensorflow. Thus, the ADAM optimizer is used as a substitute. The style image is always resized to the same size as the content image before computing the feature representation in order to extract image information on comparable scales.

## V. SUMMARY OF THE ALGORITHM

The summary of the algorithm is shown in Figure 8. The main steps are as follows:

- 1) Provide the input images for the style  $\vec{d}$  and content  $\vec{p}$  and are passed through the network.
- 2) The style representation  $P^l$  is computed on all layers. The style Gram matrix at each layer is also calculated from this feature representation.
- 3) The content representation  $A^l$  is computed on one layer.
- 4) The randomly generated white noise image is passed through the network to generate its style features  $G^l$  and its content features  $F^l$ .
- 5) On each layer included in the style representation, the mean-squared error is calculated between the Gram matrix of the generated image  $G^l$  and the Gram matrix of the original image  $A^l$  to provide the style loss  $\mathcal{L}_{\text{style}}$ .
- 6) On the one layer included in the content representation, the mean-squared error is calculated between the generated content features  $F^l$  and the photograph's content representation  $P^l$  to provide the content loss  $\mathcal{L}_{\text{content}}$ .
- 7) Since the total loss  $\mathcal{L}_{\text{total}}$  is a linear combination of the content and style losses, use these two losses to compute it.

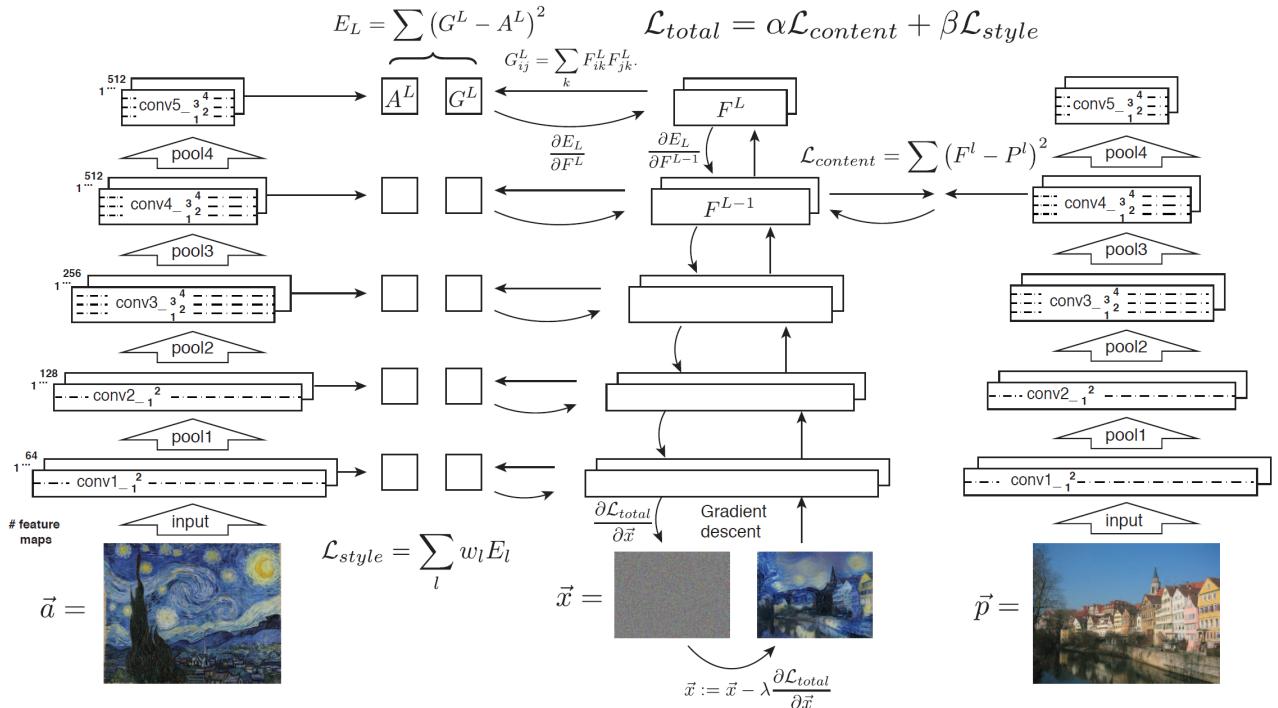


Figure 8. Style Transfer Algorithm. The detailed steps are discussed in Section V - Summary of the Algorithm

- 8) The derivative of the total loss with respect to the pixel values of the generated image  $\frac{\partial \mathcal{L}_{\text{total}}}{\partial \vec{x}}$  is computed through standard back propagation.
  - 9) The gradient computed in the previous step is used as an input into the ADAM optimizer and is iterated until it simultaneously matches the style features of the style image  $\vec{s}$  and content photograph  $\vec{p}$

## VI. PYTHON CODE

The following section details the code used for this project. The first program involves functions that are utilized for the image utilities. This includes loading the image, processing the image for the VGG19 network, deprocessing it back as an image, and saving the image.

The second code is the main function. It outlines the functions used for setting up the VGG19 model, as well as the loss functions and feature extractions. The last part of the code is the actual execution that runs the training steps and saves the resulting image.

Appendix A, which is attached with this report, shows a step by step result of the code as it's run. It was executed using Google Colab and demonstrates the output and flow of how image style transfer actually works based off the written code.

#### A. Image Utilities

```
# -*- coding: utf-8 -*-
"""
Created on Mon Apr  6 13:26:18 2020

@author: ealegre
"""

import tensorflow as tf
import numpy as np
from PIL import Image
import matplotlib as plt

# IMAGE UTILITIES – FUNCTION 1 – LOAD IMAGE

def load_img(img_path):
    # Read the image and convert the computation graph to a image format
    img = tf.io.read_file(img_path)
    img = tf.image.decode_image(img, channels=3)
    img = tf.image.convert_image_dtype(img, tf.float32)

    # Setting the scale parameters to change the size of the image
    # Get the width and height of the image. Cast it to a float so it can be divided
    shape = tf.cast(tf.shape(img)[-1], tf.float32)
    # Set the absolute maximum dimension for the image
    # max_dim = 1024
    max_dim = 512
    # Find which side is the longer side, this will be used to generate our scale
    max_side = max(shape)
    scale = max_dim / max_side
    new_shape = tf.cast(shape * scale, tf.int32)
    img = tf.image.resize(img, new_shape, method=tf.image.ResizeMethod.BILINEAR)

    img = img[tf.newaxis, :]
    return img

# IMAGE UTILITIES – FUNCTION 2 – DEPROCESS IMAGE
```

```

37
38     def deprocess_img(processed_img):
39         processed_img = processed_img*255
40         processed_img = np.array(processed_img, dtype=np.uint8)
41         if np.ndim(processed_img)>3:
42             assert processed_img.shape[0] == 1
43             processed_img = processed_img[0]
44         return Image.fromarray(processed_img)
45
46 # IMAGE UTILITIES – FUNCTION 3 – SAVE IMAGE
47
48 def save_img(best_img, path):
49     img = Image.fromarray(best_img)
50     img.save(path)
51
52 def imshow(image, title=None):
53     if len(image.shape) > 3:
54         image = tf.squeeze(image, axis=0)
55
56     plt.imshow(image)
57     if title:
58         plt.title(title)

```

---

## B. Model Utilities & Main Function

```

1  # -*- coding: utf-8 -*-
2 """
3 Created on Mon Apr  6 13:16:38 2020
4
5 @author: ealegre
6 """
7
8 # DEPENDENCIES, VARIABLES, AND PATHS
9
10 # Libraries
11 import tensorflow as tf
12 from tensorflow.python.keras import models
13 import os
14 import glob
15 import image_utils as IU
16
17 # Define paths
18 content_path = 'input/content/elon.jpg'
19 style_path = 'input/style/ironman.jpg'
20 init_path = 'input/init/512.jpg'
21 lena = 'input/content/lena_test.png'
22 lion = 'input/content/lion.jpg'
23 dog = 'input/content/dog.jpg'
24
25 # Load the style and content images
26 style_img = IU.load_img(style_path)
27 content_img = IU.load_img(content_path)
28 generated_img = tf.Variable(content_img)
29
30 # Define which layers are to be used for this model. These
31 # layers are defined in Section
32 # 3 of Gatys' paper
33 content_layers = ['block5_conv2']
34 style_layers = ['block1_conv1',
35                 'block2_conv1',
36                 'block3_conv1',
37                 'block4_conv1',
38                 'block5_conv1']
39 num_content_layers = len(content_layers)
40 num_style_layers = len(style_layers)
41
42 # MODEL UTILITIES – FUNCTION 1 – VGG19 LAYERS
43
44 def get_vggLayers(style_or_content_layers):
45     """
46     Creates the model with intermediate layer access
47
48     This function will load in the VGG19 model used in
49     Gatys' paper, with access to the
50     intermediate layers. A new model will be generated by
51     using these layers that will take an
52     input image and return an output from the intermediate
53     layers from the VGG19 model
54
55     Returns
56     -----
57     A Keras model that takes inputs and outputs of the
58     style and content intermediate layers
59
60     """
61     # Instantiate the VGG19 model. We are not including
62     # the fully connected layers, instantiate the
63     # weight based off ImageNet training
64     vgg = tf.keras.applications.VGG19(include_top=False,
65                                         weights='imagenet')
66
67     # VGG19 model is already trained off ImageNet, set
68     # trainable to False
69     vgg.trainable = False
70
71     # All we are doing is clipping the model accordingly.
72     # Just set the input of the model to be the same as
73     # its original
74     # The outputs are going to be set to the layers
75     # specified in either the style or content layers,
76     # this will
77     # be set once you've passed one of them into the
78     # function as an argument
79     vgg_input = [vgg.input]
80     vgg_output = [vgg.get_layer(name).output for name in
81                  style_or_content_layers]
82     model = tf.keras.Model(vgg_input, vgg_output)
83
84     return model
85
86 # MODEL UTILITIES – FUNCTION 2 – THE GRAM MATRIX
87
88 def gram_matrix(input_tensor):
89     """
90     Generates the Gram matrix representation of the style
91     features. This function takes an input tensor and
92     will apply the proper steps to generate the Gram matrix
93
94     Returns
95     -----
96     A tensor object Gram matrix representation
97
98     """
99     # Equation (3)
100
101     # Generate image channels. If the input tensor is a 3D
102     # array of size Nh x Nw x Nc, reshape
103     # it to a 2D array of Nc x (Nh*Nw). The shape[-1] takes
104     # the last element in the shape
105     # characteristic , that being the number of channels. This
106     # will be our second dimension
107     channels = int(input_tensor.shape[-1])
108
109     # Reshape the tensor into a 2D matrix to prepare for
110     # Gram matrix calculation by multiplying
111     # all of the dimensions except the last one (which is
112     # what the -1 represents) together
113     Fl_ik = tf.reshape(input_tensor, [-1, channels])
114
115     # Transpose the new 2D matrix
116     Fl_jk = tf.transpose(Fl_ik)
117
118     # Find all the elements in the new array (Nw*Nh)
119     n = tf.shape(Fl_ik)[0]
120
121     # Perform the Gram matrix calculation
122     gram = tf.matmul(Fl_jk, Fl_ik)/tf.cast(n, tf.float32)
123
124     # Generate the Gram matrix as a tensor for use in our
125     # model
126     gram_tensor = gram[tf.newaxis, :]
127     return gram_tensor
128
129 # MODEL UTILITIES – FUNCTION 3 – STYLE & CONTENT MODEL
130 # CLASS
131
132 class StyleContentModel(tf.keras.models.Model):
133     def __init__(self, style_layers, content_layers):
134         super(StyleContentModel, self).__init__()
135         self.vgg = get_vggLayers(style_layers + content_layers)
136
137         self.style_layers = style_layers

```

```

115     self.content_layers = content_layers
116     self.num_style_layers = len(style_layers)
117     self.vgg.trainable = False
118
119     def call(self, inputs):
120         # Expects a float input between [0, 1]
121         inputs = inputs * 255.0
122         preprocessed_input = tf.keras.applications.vgg19.\
123             preprocess_input(inputs)
124         outputs = self.vgg(preprocessed_input)
125         style_outputs, content_outputs = (outputs[:self.\
126             num_style_layers], outputs[self.num_style_layers:\
127             :])
128
129         style_outputs = [gram_matrix(style_output) for \
130             style_output in style_outputs]
131
132         content_dict = {content_name:value for content_name, \
133             value in zip(self.content_layers, content_outputs\
134             )}
135
136         style_dict = {style_name:value for style_name, value \
137             in zip(self.style_layers, style_outputs)}
138
139         return{'content':content_dict, 'style':style_dict}
140
141 feature_extractor = StyleContentModel(style_layers, \
142     content_layers)
143
144 # Load the style and content images
145 style_img = IU.load_img(style_path)
146 content_img = IU.load_img(content_path)
147
148 # Extract the features from the style and content images
149 style_features = feature_extractor(style_img)[‘style’]
150 content_features = feature_extractor(content_img)[‘content’]
151
152 # The Tensorflow variable function initializes our image \
153 # to be used for gradient descent.
154 # This image has to be the same size and type as the \
155 # content image, which means that as long as it’s \
156 # loaded in before being
157 # called upon as the generated image, it should be fine to \
158 # use. This is essentially the image that will be \
159 # optimized.
160 # In case you want to use another image, you can uncomment \
161 # the code below
162 ...
163 init_img = load_img(content_path)
164 generated_image = tf.Variable(init_img)
165 ...
166 generated_img = tf.Variable(content_img)
167
168 # Remember that all of our images are read in as a float32 \
169 # type, so we have to define the range as [0, 1] to \
170 # keep it within 255
171 def clip_range(img):
172     return tf.clip_by_value(img, clip_value_min=0.0, \
173         clip_value_max=1.0)
174
175 # Choose an optimizer. The paper chose L-BFGS but \
176 # Tensorflow doesn’t have that, so we’ll use ADAM
177 optimizer = tf.optimizers.Adam(learning_rate=0.02, beta_1=\
178     0.99, epsilon=1e-8)
179
180 # Remember that we are trying to optimize the Total Loss \
181 # function. However, there are values that correspond \
182 # to the style and content weight
183 # The Alpha value corresponds to content weight and the \
184 # Beta corresponds to the style weight. Let’s define \
185 # these
186 alpha = 1e4
187 beta = 10
188
189 # MODEL UTILITIES – FUNCTION 4 – STYLE, CONTENT, & TOTAL \
190 # LOSS FUNCTIONS
191
192 # Now we define our loss functions
193 def style_content_loss(outputs):
194     style_outputs = outputs[‘style’]
195     content_outputs = outputs[‘content’]
196
197 # Equation 5
198 # Define the style loss function
199 style_loss = tf.add_n([tf.reduce_mean(tf.square(\
200     style_outputs[name] - style_features[name])) for \
201         name in style_outputs.keys()])
202
203 # Multiply the style loss by the weighted variable to \
204 # get the weighted style loss
205 style_loss *= beta / num_style_layers
206
207 # Equation 1
208 # Define the content loss function
209 content_loss = tf.add_n([tf.reduce_mean(tf.square(\
210     content_outputs[name] - content_features[name])) \
211         for name in content_outputs.keys()])
212
213 # Multiply the content loss by the weighted variable to \
214 # get the weighted content loss
215 content_loss *= alpha / num_content_layers
216
217 # Equation 7
218 # Define the total loss function
219 total_loss = style_loss + content_loss
220
221 return total_loss
222
223 # State the total variational weight
224 total_variational_weight = 30
225
226 @tf.function()
227 def train_step(img):
228     # Tensorflow’s GradientTape function performs automatic \
229     # differentiation of the input
230     with tf.GradientTape() as tape:
231         outputs = feature_extractor(img)
232         loss = style_content_loss(outputs)
233         loss += total_variational_weight * tf.image.\
234             total_variation(img)
235
236     # Apply the gradient by passing the loss and the \
237     # generated image
238     grad = tape.gradient(loss, img)
239
240     # The gradient is optimized using the ADAM optimizer we \
241     # declared earlier. This optimizes the \
242     # generated image using the gradient values
243     optimizer.apply_gradients([(grad, img)])
244
245     # The image is rewritten, being sure that the values all \
246     # stay within the viable range of [0, 1]
247     img.assign(clip_range(img))
248
249 import time
250 from IPython import display
251 # Start the run time
252 start = time.time()
253
254 epochs = 100
255 iterations = 100
256
257 step = 0
258 for n in range(epochs):
259     for m in range(iterations):
260         step += 1
261         train_step(generated_image)
262         print(".", end=' ')
263         display.clear_output(wait=True)
264         display.display(IU.deprocess_img(generated_image))
265         print("Train step: {}".format(step))
266
267 end = time.time()
268 print("Total Time: {:.1f}".format(end-start))
269
270 file_name = ‘output/stylized-image.png’
271 IU.deprocess_img(generated_image).save(file_name)

```



Figure 9. Theoretical Results from Gatys et al. The images were created by finding an image that simultaneously matches the content representation of the photograph and the style representation of the artwork. The original photograph depicting the Neckarfront in Tübingen, Germany, is shown (from the top left downwards then across to the right side) in **A** (Photo: Andreas Praefcke). The painting that provided the style for the respective generated image is shown in the bottom left corner of each panel. **B** The Shipwreck of the Minotaur by J.M.W. Turner, 1805. **C** The Starry Night by Vincent van Gogh, 1889. **D** Der Schrei by Edvard Munch, 1893. **E** Femme nue assise by Pablo Picasso, 1910. **F** Composition VII by Wassily Kandinsky, 1913.

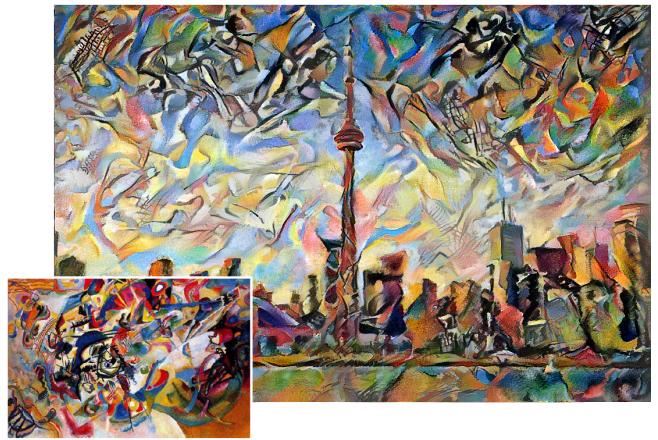


Figure 10. Actual results from implementation. The original photograph depicting the skyline in Toronto, Canada, is shown (from the top left downwards then across to the right side) in A (Photo: Andreas Praefcke). The painting that provided the style for the respective generated image is shown in the bottom left corner of each panel. **B** The Shipwreck of the Minotaur by J.M.W. Turner, 1805. **C** The Starry Night by Vincent van Gogh, 1889. **D** Der Schrei by Edvard Munch, 1893. **E** Femme nue assise by Pablo Picasso, 1910. **F** Composition VII by Wassily Kandinsky, 1913.

## VII. RESULTS, DISCUSSIONS, AND FUTURE WORK

### A. Theoretical Results

The results achieved by Gatys et al is shown in Figure 9. There are a couple of interesting differences to note. First and foremost, while it was very easy to derive the generic steps, there are some nuances that Gatys does not explicitly state in his paper. One of these is how many iterations are ran per result. The number of iterations dictates how much of the "style" is transferred over to the newly synthesized image. However, this can be remedied easily as it is very obvious to note after some iterations when the loss function converges. However, Gatys never stated what size his initial or final images were, despite his suggestion of constraining the image size before passing it through the VGG19 network.

The second is the optimizer. Gatys suggest using the L-BFGS optimizer to optimize the gradient of the total loss function with respect to the generated image. However, as discussed previously, the L-BFGS optimizer is not implemented in Tensorflow. Furthermore, the parameters of the optimizer were not stated.

The third nuance is that Gatys extracted the content image on layer 4\_2. While it is never explicitly stated why this was done, it is assumed that Gatys assumed that extracting the content features from this layer was high enough to allow the more profound "features" of the content image to be displayed in the style transfer.

Finally, while the results from Gatys' style transfer serve as the benchmark, there is no numerical results to quantify how well his method works. As such, the results are evaluated just as art would be in real life: objectively.

### B. Experimental Results

The implemented results from my code can be viewed in Figure 10. This section will discuss the deviations and remedies to the nuances discussed earlier. First and foremost, our images were synthesized with the following benchmark parameters

- Maximum Image Dimension: 1024 pixels
- Iterations/Epoch: 100
- Epochs: 100
- Optimizer: ADAM
- Learning Rate: 0.02
- Beta: 0.99
- Epsilon: 1e-1
- Alpha (Content Weight): 1e4
- Beta (Style Weight): 10
- Content Layers: 5\_2

These parameters are important to note as some deviate from Gatys' original suggestions and others are derived from experimentation due to the lack of information provided in the paper. As stated previously, Gatys never specified the number

of iterations or the image size. The maximum image size of 1024 was chosen as it is a staple size when it comes to images, as well as the fact that it is substantially large enough that the contents of the image are distinguishable when viewed. After some experimentation with the benchmark parameters, it was decided that 100 epochs with 100 iterations provided a result that was satisfactory.

The most notable deviation from Gatys' method was the difference in where the content image was extracted. Instead of extracting the content image on Gatys' suggested layer of 4\_2, our content image was extracted on layer 5\_2. This was due to the difference in the optimizer used in our algorithm. While it was possible to extract our content image on the same layer Gatys did, the results were not always pleasant or consistent. By that, it's meant that the model needs to be constantly adjusted based on the style image and content image used. It was fairly difficult to find the right balance of style/content weights and optimizer parameters that would yield suitable results from this layer. We also used the ADAM optimizer that was tuned with a learning rate of 0.02, a beta value of 0.99, and an epsilon of 1e-1. Using Gatys' Alpha and Beta values also yielded poor results, as we basically had to switch his values. In Gatys' method, the style weight was weighed significantly more than the content weight, ours was the complete opposite. These values were achieved after some investigation online. These were the hyper parameters that contributed to the deviation of our results from Gatys'.

The final deviation that our method has is that it accounts for total variational loss. While this was never mentioned in Gatys' paper, the total variational loss allows for the elimination of high frequency artifacts from our synthesized image, making them appear more natural.

Overall, our method seems to work quite well, although not as well as Gatys'. The majority of the images are comparable to the original results, however, the deviations become apparent when comparing images synthesized from Van Gogh's Starry Night. Gatys' results yielded a comparable image, with Van Gogh's swirled portrayal of the night sky as being more prominent, as well as the actual moon being included in the final image. Our result does a good job of copying this swirl pattern but not to the same degree. Furthermore, the patterns of the stars are not as intelligently placed as Gatys' results are. Surprisingly enough, Starry Night was used as the default image and is what was used to test the parameters found above. Other images were very easy to tune regardless of the parameters but Van Gogh's Starry Night had a very small window of passable results, which prompted the final benchmark values seen above.

### C. Future Work

First and foremost, any future work done on this project should focus on trying to implement Gatys' method as closely as possible. While most modern and up to date tools make it easy to get close to Gatys' method, there are elements that are

lacking and crucial to obtaining the same results. Specifically, the optimizer used for gradient descent plays a larger role than originally understood. Despite my best efforts, implementing the L-BFGS optimizer was out of the scope for this project as it would require a much deeper understanding of numerical analysis than what this project entailed. Furthermore, many online resources have stated that using the ADAM optimizer was "comparable" to using the L-BFGS optimizer, it was never truly elaborated on what that meant. While the results obtained can be called comparable to Gatys', there are still glaring issues that need to be addressed.

Furthermore, the hyper parameters used could also be fine tuned. With the use of a more powerful machine and parallel computing, different ranges and combinations of these parameters can be run in parallel to see what would give the closest result to the original.

There was also an interest in writing a script that produces a high resolution stylized image. The framework of this is to essentially start from a low resolution input, with a maximum dimension of 256 pixels, run that through the algorithm for a certain number of iterations, save the result and then use this result as the generated image for a secondary run with a maximum dimension of 512 pixels. After each run, the process would be repeated, with the result of the previous run being used as the generated initial image and continuously up-scaling the output until the desired high resolution result was achieved. The issue with this is that it would be fairly time consuming as the number of iterations would need to change every time, the hyper parameters tuned after every run, and the run times would increase exponentially as the resolution of the image increased. Current implementations exist and utilize 4 GPUs which result in an image with a maximum dimension of 3620 pixels.

The most impressive improvement to this algorithm would come in the form of either *photorealistic style transfer* or *multiple style transfer*. While the prior was not researched during this project, the latter already has an implementation. With a bit of tweaking, this code could accommodate multiple style images that would be used for the style transfer. The only constraint is optimization of

#### REFERENCES

- [1] L. Gatys, A. Ecker, M. Bethge, "Image Style Transfer Using Convolutional Neural Networks," 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, 2016, pp. 2414-2423.