

Image Style Transfer		MLOA & NNDL
2	Problem statement	3
3	Background work / literature survey	4
4	Proposed system / implementation	6
4.1	Feature extraction using pretrained model (VGG19)	6
4.2	Loss functions and Gram matrix	9
4.3	Optimizers	11
4.4	Final algorithm	12
5	Evaluation of results	14
6	Important Code Snippets	22
7	Conclusion and future work	24
8	References	25

Chapter 1: Introduction

Style transfer is an area of computer vision that refers to the conversion of visual style of an image or video without changing its content. This technique allows us to combine artistic style from one image with another image's content, resulting in images or videos with visually appealing compositions.

This concept has gained popularity because of its wide-ranging applications across fields like photography, cinematography, and digital art. It is frequently used in camera filters on various social media platforms like Snapchat and Instagram. The style can be instantly transferred in real-time within milliseconds.

Through the advancements in image processing and deep learning, automated style transfer algorithms have emerged. Traditional signal processing methods converted the image data through filters and mathematical operations. However it lacked the ability to replicate complex artistic styles.

In contrast, deep learning models like CNN can analyse vast image datasets to automatically learn diverse artistic styles. It is revolutionising how we transform styles by making the results more realistic. This utilisation of neural networks in style transfer tasks is commonly known as Neural Style Transfer or NST in short.

In this project, we have implemented a NST technique that utilises a pre-trained VGG neural network model for feature extraction. Then we define content, style, and total variation loss functions to optimise a generated image, resulting in visually appealing compositions by combining content and style images.



Figure 1: Famous example of style transfer merging a content image of a dog with a style image, painted by renowned Russian painter Wassily Kandinsky, resulting in a combined image created through NST.

Chapter 2: Problem Statement

The problem statement of the project is to implement NST using a pre-trained VGG neural network model to combine the content image, such as a dog, with the artistic style of another image, such as paintings by Wassily Kandinsky.

This involves defining appropriate loss functions, including style, content, and total variation losses, and iteratively optimising a generated image to minimise these losses (Specifically using different optimization algorithms like Adam, RMSprop and Adagrad) to achieve a balance between retaining the high-level content attributes of the content image and integrating the low-level stylistic traits of the style image.

The process involves a forward pass through the VGG model to extract features, followed by a backward pass to compute gradients and update the generated image.

Here are some of the questions that we aimed to find out through this project:

1. How do different optimization algorithms (Adam, Adagrad, Rmsprop) perform in the context of image style transfer, particularly in terms of convergence speed and quality of output?
2. What are the effects of varying learning rates and epochs on the performance of optimization algorithms in image style transfer?
3. How can feature maps and filter maps from pre-trained convolutional neural networks (VGG) be effectively visualised to aid in image style transfer tasks?

4. How can loss functions be optimised to better capture the difference between content and style images in style transfer tasks?

Chapter 3: Background/Literature review

In paper [3] the author gives a detailed and comprehensive study of the role of the Convolutional Neural Networks (CNNs) in the fundamental perception skills like object recognition, texture detection and segmentation, and synthesis. It involves a comprehensive discussion on the matting of textures, image style representation, and construction of hierarchical models to accurately predict the visual cortex neuronal responses. Additionally, the article explains a new method such that convolutional neural nets can be used to convert artistic styles onto photoimages. The developed algorithm does quite well and by so doing it weakens attachment to image content and brings the two aspects together to create a visually arresting image which is, therefore, amalgamative. It demonstrates the balance between content preservation and the style transfer, explores the role of different network layers and achieves the generation of high-quality stylized for each input image.

Research paper [1] stresses types of the style transfer techniques associated with neural networks for the first class, the second class, the third class, and the generative adversarial networks (GAN) for the last class. These techniques get trained with distinct data sources such as the Microsoft COCO, Imagenet and Cityscapes. Furthermore, it also studies the hardships as well the future course for neural image style transfer discourse comprising the need for the formulation of the standardized assessment criteria and the improvement in image fidelity.

The paper will be researching [2] that involved Neural Style Transfer (NST) in image style transformation. Different models to include the NST model of Gatys et al., Fast NST and the VGG-19 will be compared and contrasted to ascertain what works best in the stylising of the image. It will be examining the development story of NST technology, outline its diverse operational modes, and describe numerous industry applications beyond current use. Then, it will end up identifying core future research opportunities and areas where NST development could be enhanced even further.

The paper [4] describes an approach to image style transfer utilizing VGG-16 and VGG-19 ceiling network models. It evaluates the two models, pointing out that VGG-16 as the best is ideal for style transfer. The model involved in architecture for style-image transfer, feature extraction, calculating content and style losses, and the final weight combination of losses in the generated images are thoroughly explained. The study points to the applicability of the technique in deeper feature extraction to optimize more complex image processing tasks.

Besides, the article consists of an experimental procedure with parameters' settings and outcomes, which can be useful for other researchers who are interested in optimizing it.

Paper [5] discusses a comparative study between two neural style transfer techniques: VGG19 by Gram and fast mathematics way. The paper weighs the performance of two schemes in

production of good samples representation and it found fast-style conversion method as more efficient in computation time and quality of the resulting image. Emphasizing on the uniqueness of balancing style and information preservation in the context of neural style transfer and with Fast-style transfer claimed to favor the combination of fidelity to style and preservation of content, the research points a finger to the most successful balance of the two concerns.

In [6] it is shown how a method of image style transfer is developed which uses an enhanced style loss function which itself is designed using the most advanced Gram matrix feature calculation. This method uses a neural network that is trained by a joint content and style loss functions in order to produce a GAN that can facilitate style-transferred images generation.. The algorithm is expected to improve the image quality by adjusting the style loss function local and adjacent features, so the similarity will be maintained and the styling distortion will be prevented to be more likely to be represented in the details and spatial position. The experimental results showing significant improvement of the Structural Similarity Index also known as SSIM and the Peak to Signal Noise Ratio also known as PSNR compared to the other styles of transfer techniques have been achieved. In addition, these results show that its quality metrics outperform the Gatys algorithm.

Chapter 4: Proposed System / implementation

In our project, we've combined deep learning methods, particularly NST, with optimization techniques to develop a system that can produce artistic compositions. These compositions mix the style of one image with the content of another. We've used a pre-trained neural network to extract features from both images in the form of feature maps. Then, we've created appropriate loss functions to quantify the differences in content and style representations. To improve the resulting image and reduce the losses, we've applied optimization algorithms iteratively with varying number of epochs and learning rates to find the most suitable algorithm and hyperparameters.

4.1. Feature extraction using pretrained model (VGG19):

Our approach relies heavily on a pre-trained model, particularly VGG19. By making use of this pre-trained network, we take advantage of its capability to effectively capture minute patterns and semantic details from images. VGG19 is a convolutional neural network (CNN) architecture comprising 19 layers. It was developed by the Visual Geometry Group at the University of Oxford. It is trained on a vast dataset known as ImageNet which has around 1000 different target classes. Thus, this network has learned to identify a diverse range of visual patterns and objects. We have used pre-trained ImageNet weights from the VGG19 CNN model to generate feature maps of our input images.

The model architecture comprises a series of convolutional layers followed by max-pooling layers. The convolutional layers are crucial in extracting features from the input image, while the max-pooling layers are used to downsample the feature maps. These convolution layers play a very important role in image based applications by extracting features that capture various spatial patterns and structures. VGG uses 3x3 convolution filters and 2x2 pooling filters. VGG19, specifically, consists of 16 convolutional layers. Each convolution layer is followed by a rectified linear activation function or ReLU for introducing non-linearity. These convolutional layers are organised into several blocks, with each containing multiple convolutional layers.

Following the convolutional layers, VGG19 consists of three fully connected layers, succeeded by a softmax layer for classification, particularly in image recognition tasks. These fully connected layers merge the features extracted by the convolutional layers to generate predictions regarding the input image.

The architecture of VGG19 is depicted in figure 2 which highlights different dimensions of the layers and layer names.

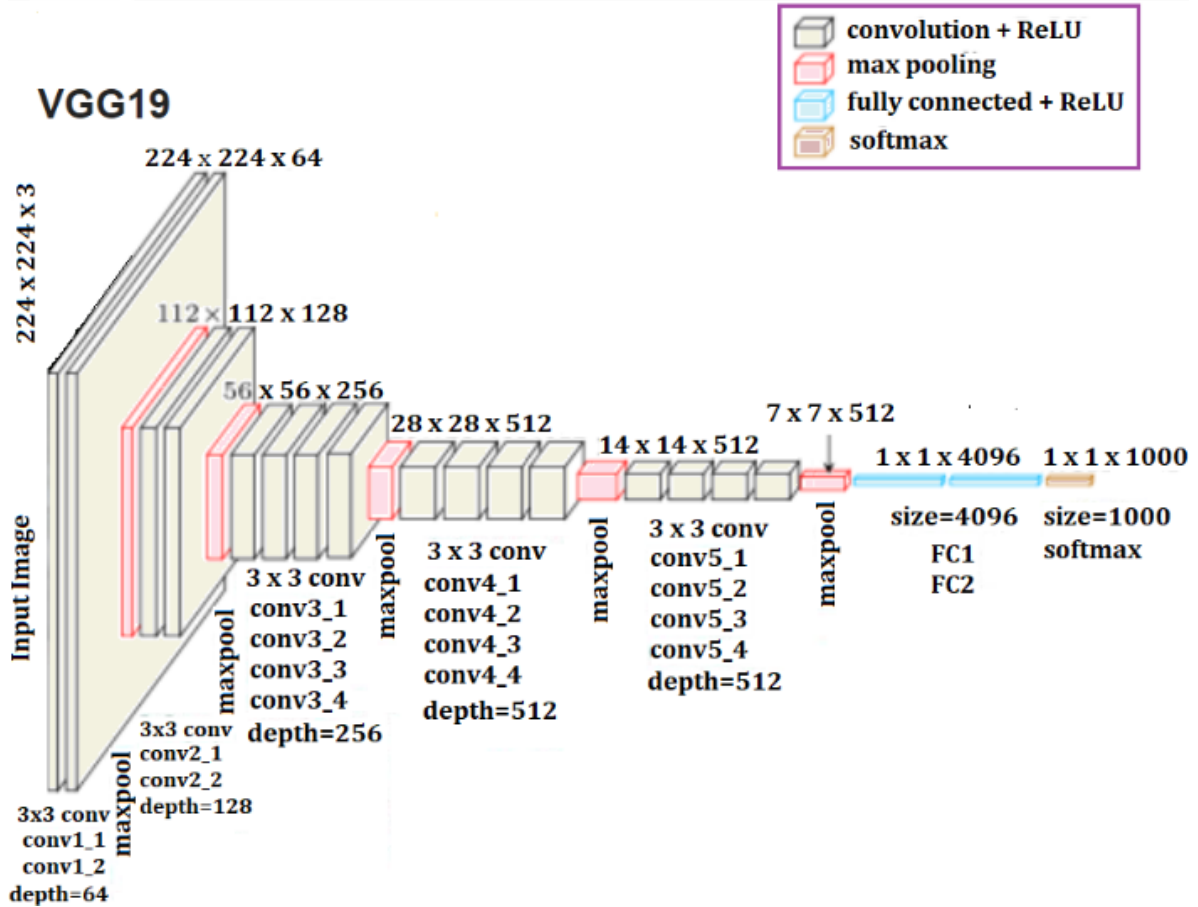


Figure 2: VGG-19 architecture Source: [7]

One important thing to note is that, In NST, the fully connected layers (last 2 layers along with softmax) are removed, and only the convolutional layers are utilised for feature extraction. By utilising VGG, NST can extract high-level content and style features from images. The content features show what the content image looks like, while the style features show the textures, colours, and patterns of the artistic style in the style image. Figure 3 represents a basic image transfer style network.

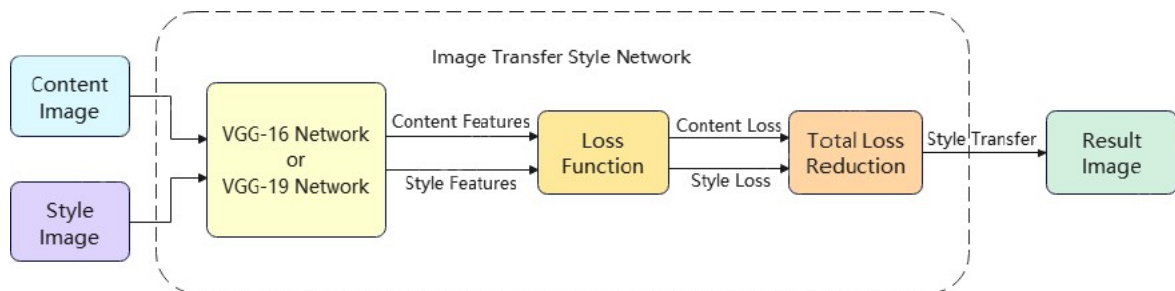


Figure 3: Image transfer style network Source: [4]

In the VGG19 network, earlier in the network is fine-tuned tight details like horizontal lines, vertical lines, diagonal lines, brightness, etc. and as we move through the network, we build up larger and larger sorts of features and also more abstract.

Because deeper features are more effective, the "block5_conv2" or "conv5_2" convolutional layer is selected to extract content information. On the other hand, for the style image, multiple layers from different depths of the neural network are often utilised to capture a comprehensive representation of style features. This approach allows the NST algorithm to extract style information ranging from low-level textures to high-level patterns present in the style image. We have used the combination of "block1_conv1," "block2_conv1," "block3_conv1," "block4_conv1," and "block5_conv1," to capture various aspects of style.

For example the following figure 5 & 4 show the 64 feature maps of the first layer along with the filters of the first layer.

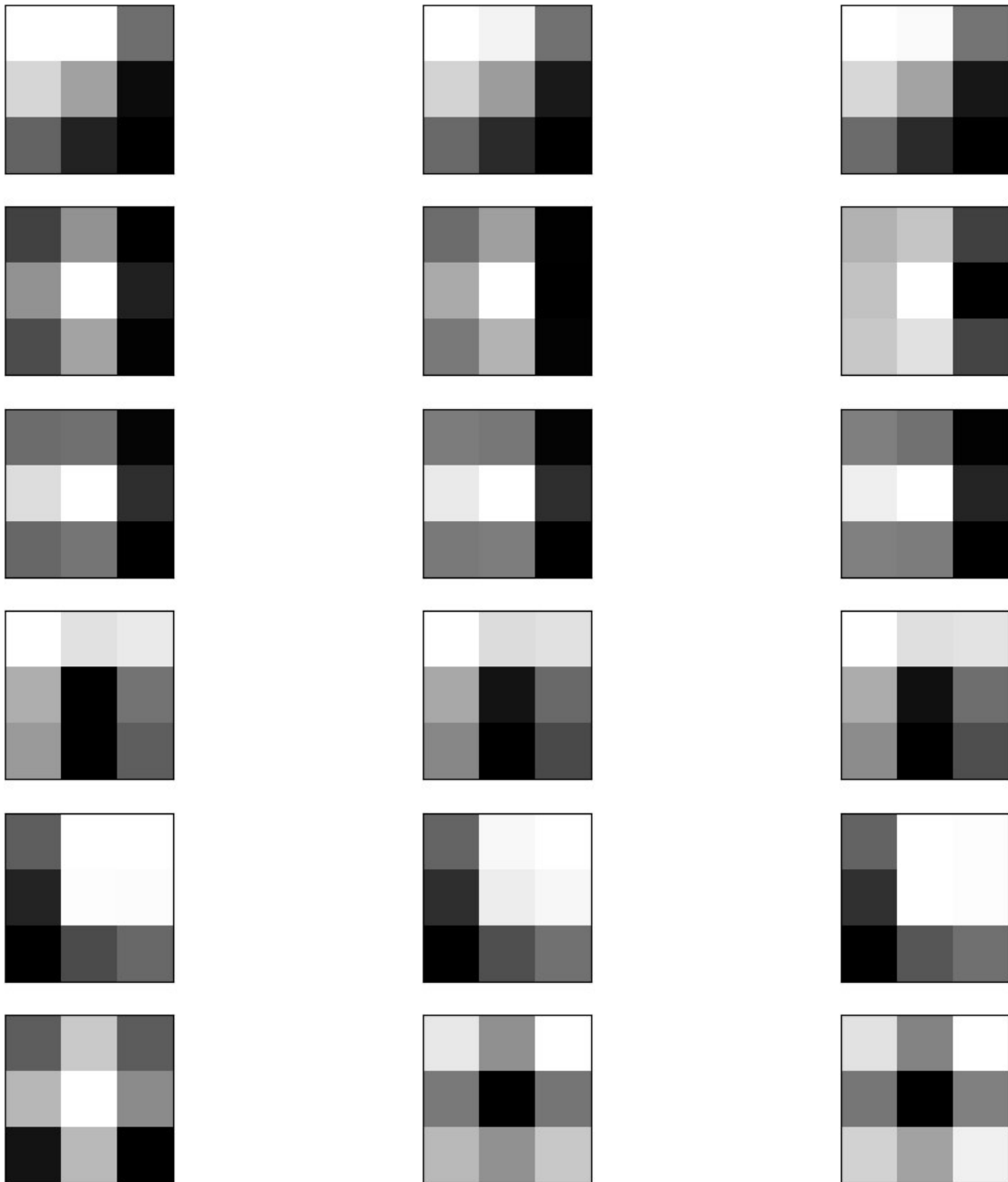


Figure 4: Filters of the first layers



Figure 5: 64 feature maps of the first hidden layer for an image containing a portrait

4.2. Loss functions and Gram matrix:

We can measure the differences in content and style between images by establishing appropriate loss functions derived from the representations mentioned above. Two main loss functions are defined: content loss and style loss.

Important note: Before moving forward consider the following notations - Let a represent the style image, x represent the generated image, weights are represented by α and β each corresponding to content and style image losses. l denotes the layer l of the VGG19 network. F_j^l and P_j^l denotes the style image activation value j and content image at the feature map i .

Content loss: The content loss is computed by comparing the feature maps of the combined image with those of the content image. This loss encourages the generated to retain the image

structure and content of the content image. It is typically computed as the mean squared error (MSE) between feature maps extracted from the convolutional layers of a pre-trained neural network. N represents the l later N feature mapping where each size is M . G^l and A^l denotes the representation of style in style image and content image.

$$L_{content}(p, x, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2$$

Gram matrix: To compute the style representation, we calculate the Gram matrix for each layer of interest. It is used to capture the style information of an image by examining the correlations between different features. It is calculated by computing the inner product of the feature maps generated by a particular layer and then averaging over all spatial locations. The inner product is between the vector feature map denoted as i and j in l layer. Then the

style of the gram matrix is defined as

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

Style loss: To measure the difference in style between the style image and the generated image, we use the MSE between the Gram matrices of the feature maps for selected layers in the VGG19 network.

$$L_{style}(a, x) = \sum_{l=0}^L w_l E_l$$

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

Total loss: The total loss used in NST is a weighted addition of the style and the content loss. The weights for each loss term are typically chosen empirically to balance the contributions of content and style to the generated image.

$$L_{total}(p, a, x) = \alpha L_{content}(p, x) + \beta L_{style}(a, x)$$

In this equation, weighted alpha and beta are coefficients used to balance the contributions of the style and the content loss to the overall loss function. These coefficients decide how much importance to give to maintaining the content of the content image and adding the style of the style image. These coefficients can be fine tuned to get the right mix between keeping the content and transferring the style. The values of alpha and beta depend on what we want the

final image to look like. For our project, we've chosen the style weight to be 1e-2 and the content weight to be 1e4.

4.3. Optimizers:

To improve the generated image iteratively and reduce the losses discussed above, we've used optimization algorithms like Adam, Adagrad, and RMS-prop. These algorithms help in adjusting the parameters of the generated image in a way that decreases the overall loss. The goal is to move towards a final result that combines the style and content of the input images in a visually pleasing way.

Updation algorithm in NST:

1. Let I_{comb} be the combination image
2. Initialize I_{comb} randomly
3. For e in range (epochs)
 - a. Calculate $Loss(I_{comb})$
 - b. Calculate $\frac{\partial Loss}{\partial I_{comb}}$
 - c. $I_{comb} -= \lambda \frac{\partial Loss}{\partial I_{comb}}$ (General update rule e.g. SGD)

Adam Optimizer:

Adam, short for Adaptive Moment Estimation, is an adaptive optimization technique that combines the advantages of AdaGrad and RMSprop optimizers. It maintains two moving averages of gradients: the second moment (uncentered variance) and the first moment (mean). Using these estimates, the algorithm updates parameters while adjusting the learning rate for each parameter independently. The update formula for Adam optimizer can be expressed as follows:

Here m_t and v_t are the first and second moment estimates of the gradients. g_t is gradient at time t . β_1 and β_2 are decay rates for the estimates of moments. η is the learning rate ϵ is

a constant. It is very small and is used to prevent division by zero.

- $m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$
- $v_t = \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$
- $\hat{m} = \frac{m}{1 - \beta_1^t}$
- $\hat{v} = \frac{v}{1 - \beta_2^t}$
- $\theta_{t+1} = \theta_t - \frac{\eta \cdot \hat{m}}{\sqrt{\hat{v} + \epsilon}}$

RMSprop Optimizer:

RMSprop (Root Mean Square Propagation) is an adaptive learning rate optimization technique designed to beat the problem of diminishing learning rates seen in AdaGrad. It works by

maintaining a moving average of squared gradients and adjusting learning rates for each parameter independently. The update formula for RMSprop optimizer is as follows:

Here v_t is the exponentially decaying average of squared gradients. β is the decay rate. η is the learning rate. ϵ is a constant to avoid division by zero.

- $v_t = \beta \cdot v_{t-1} + (1 - \beta) \cdot g_t^2$
- $\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t} + \epsilon} \cdot g_t$

AdaGrad Optimizer:

AdaGrad (Adaptive Gradient Algorithm) is an adaptive learning rate optimization method. It is designed to give more weight to infrequent features by adapting the learning rate for each parameter using historical gradients. AdaGrad accumulates squared gradients over time. This is used to adjust the learning rate for each parameter independently. Over time, the denominator of the learning rate update can become too large which can cause the learning rate to become too small. The update rule for AdaGrad optimizer is as follows:

Here g_t is the gradient at time step t . η is the learning rate. ϵ is a constant to avoid division by zero. G_{t+1} is the accumulated squared gradient.

- $G_{t+1} = G_t + g_t^2$
- $\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_{t+1}} + \epsilon} \cdot g_t$

4.5. Final algorithm:

Considering the steps mentioned earlier, we can outline the algorithm for Neural Style Transfer as follows:

Start

1. Load Content Image
2. Load Style Image
3. Preprocess Images
4. Load Pretrained VGG Model
5. Define Content and Style Layers
6. Extract Content and Style Features
7. Initialize Generated Image
8. Define Loss Functions

Loop until convergence or desired number of iterations:

9. Forward Pass:

9.1 Compute Content Loss:

9.1.1 Extract Content Features of Generated Image

9.1.2 Compute Mean Squared Error (MSE) between Content Features of
Content

and Generated Image

9.2 Compute Style Loss:

9.2.1 Extract Style Features of Generated Image

9.2.2 Compute Gram Matrix of Style Features for Content and Style Images

9.2.3 Compute Mean Squared Error (MSE) between Style Features Gram Matrices

9.3 Compute Total Variation Loss:

9.3.1 Compute Image Gradient along x and y directions 9.3.2

Compute Total Variation Loss using Image Gradients

9.4 Compute Total Loss:

9.4.1 Combine Content, Style, and Total Variation Losses with respective weights

10. Backward Pass:

10.1 Compute Gradient of Total Loss with respect to Generated Image Pixels

11. Update Generated Image using Optimizer

12. Record Loss Values

13. Post-process Generated Image

14. Save Generated Image

15. Save Loss Values

16. Visualize Loss Curve

17. Visualize Content, Style, and Generated Images

End

Chapter 5: Evaluation of results

We experiment with three optimizers: Adam, RMSprop, and AdaGrad, each offering unique advantages in terms of convergence speed and stability. We evaluate the performance of our system using both quantitative metrics, such as loss values, and qualitative assessments based on visual inspection of the generated images. Additionally, we fine-tune hyperparameters, including learning rates and the number of optimization iterations, to achieve optimal results. Fine-tuning involves experimenting with different parameter settings and evaluating their impact on the quality of the generated images.



The NST algorithm is provided with two input images: a content image featuring a portrait and a style image depicting the renowned painting. Both images are of identical size as seen in figure 6.






Figure 6: Content and Style Image We Provided To The NST Algorithm


Table 1: Output image and total time taken for NST using different optimizers with different epochs and train steps.

Optimize r	Learning rate	Epochs	Train step	Tota l time	Output

Adam	0.02	50	250	30.6	
Adagrad	0.02	50	250	30.4	
Adagrad	0.1	50	250	30.2	

Adam	0.02	100	500	58.0	
Adagrad	0.02	100	500	58.1	
RMSprop	0.02	100	500	57.2	

Adam	0.02	500	2500	151.6	
Adagrad	0.02	500	2500	151.8	
RMSprop	0.02	500	2500	149.6	

RMSprop	0.001	500	2500	152.9	
---------	-------	-----	------	-------	--

Observations:

1. Adagrad with a learning rate of 0.02 consistently shows slightly lower total time compared to Adam and RMSprop across different epochs and train steps.
2. Adam with a learning rate of 0.02 generally shows competitive performance with Adagrad in terms of total time.
3. RMSprop tends to have slightly longer training times compared to Adam and Adagrad in most cases, especially with higher epochs and train steps.
4. Varying the learning rate within an optimizer (e.g., RMSprop with learning rates of 0.02 and 0.001) affects the total training time, with lower learning rates generally resulting in longer times.
5. Increasing epochs and train steps generally leads to longer training times, as expected.

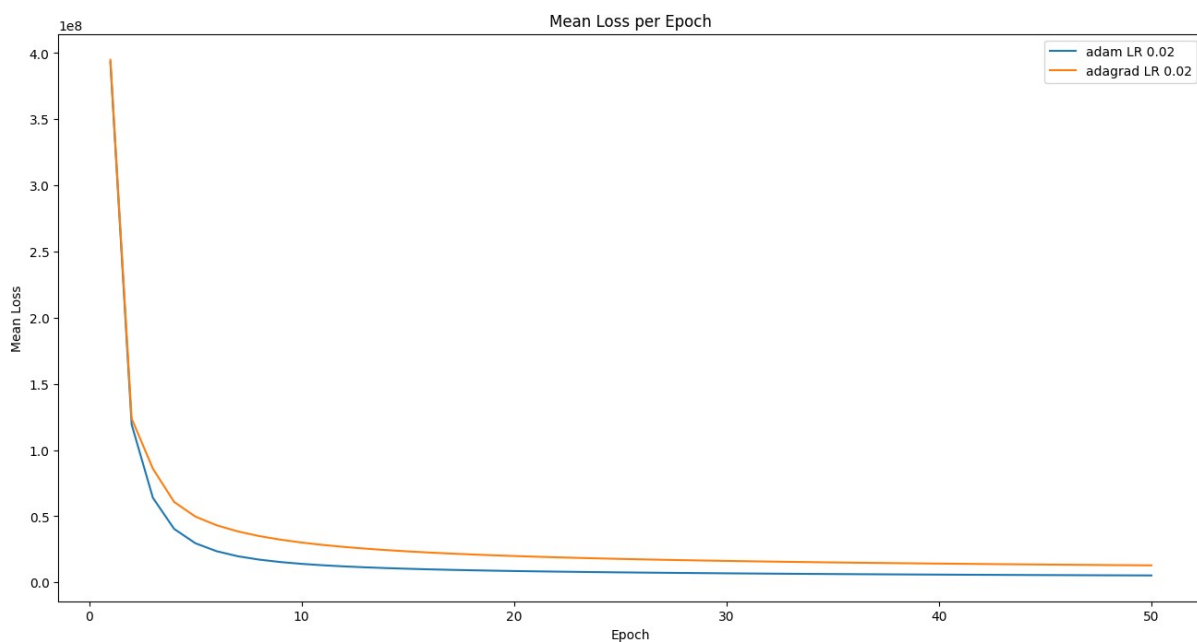


Figure 7: plot loss for Adagrad and Adam with LR=0.02 and epochs=50

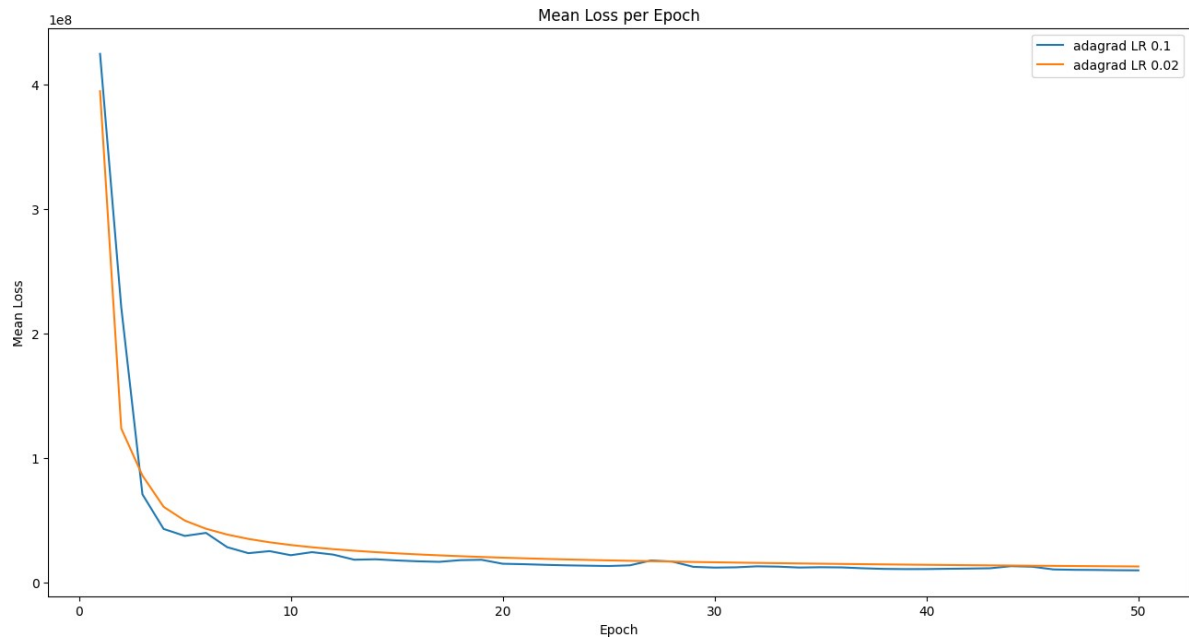


Figure 8: plot loss for Adagrad and Adam with LR=0.1 and epochs=50

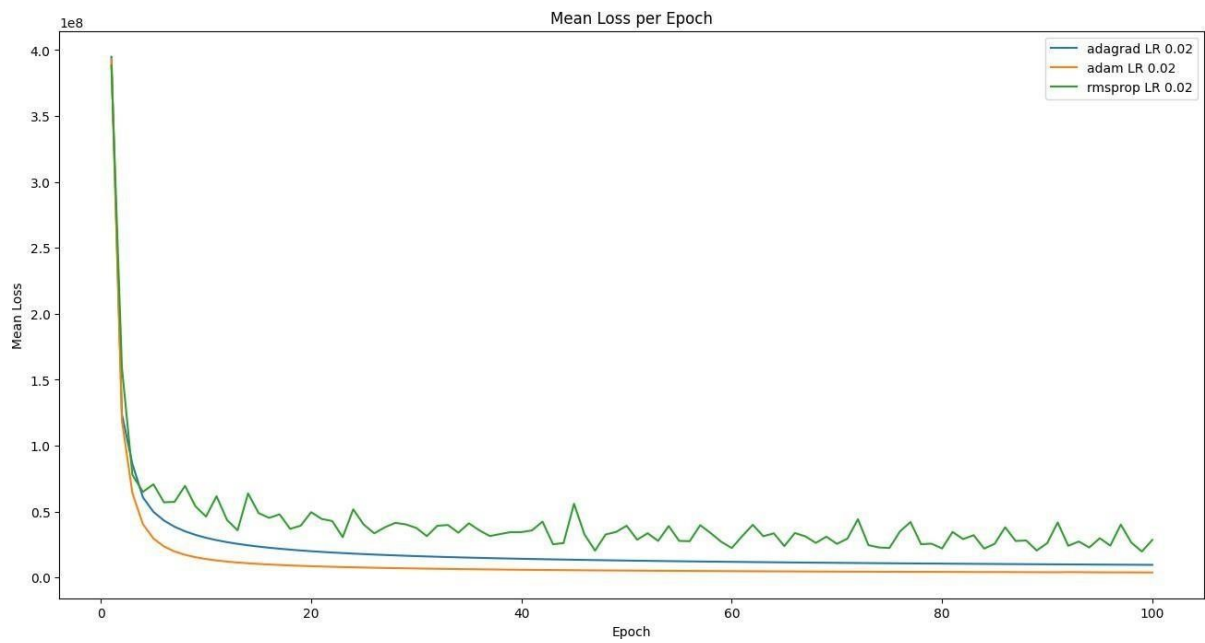


Figure 9: plot loss for Adagrad, RMSprop and Adam with LR=0.02 and epochs=100

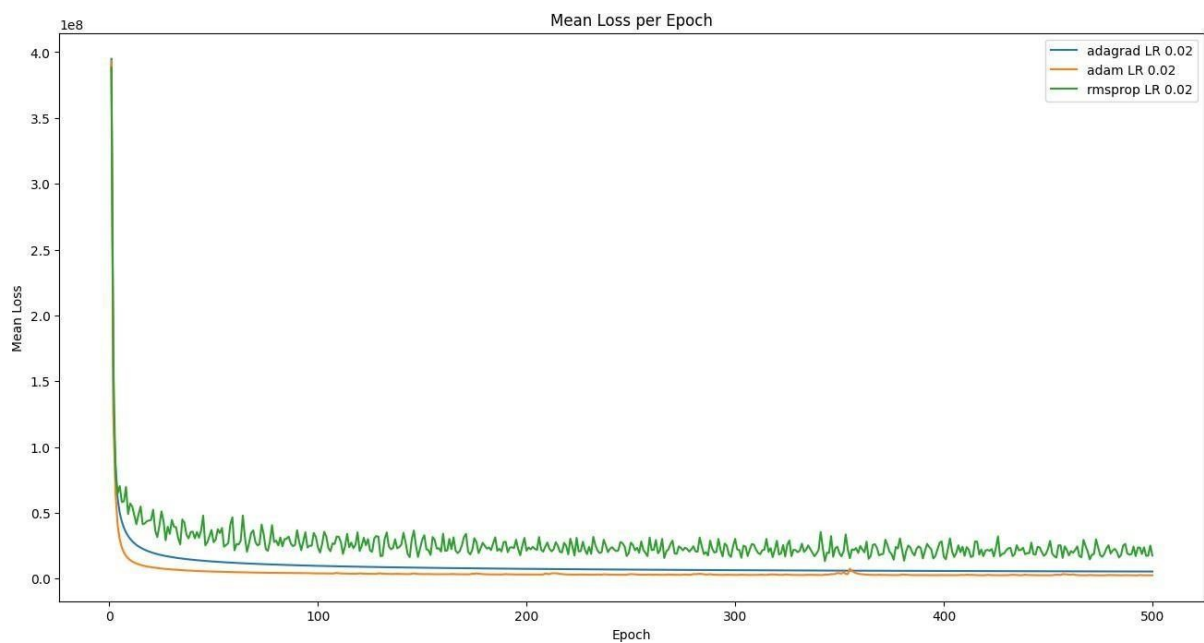


Figure 10: plot loss for Adagrad, RMSprop and Adam with LR=0.02 and epochs=500

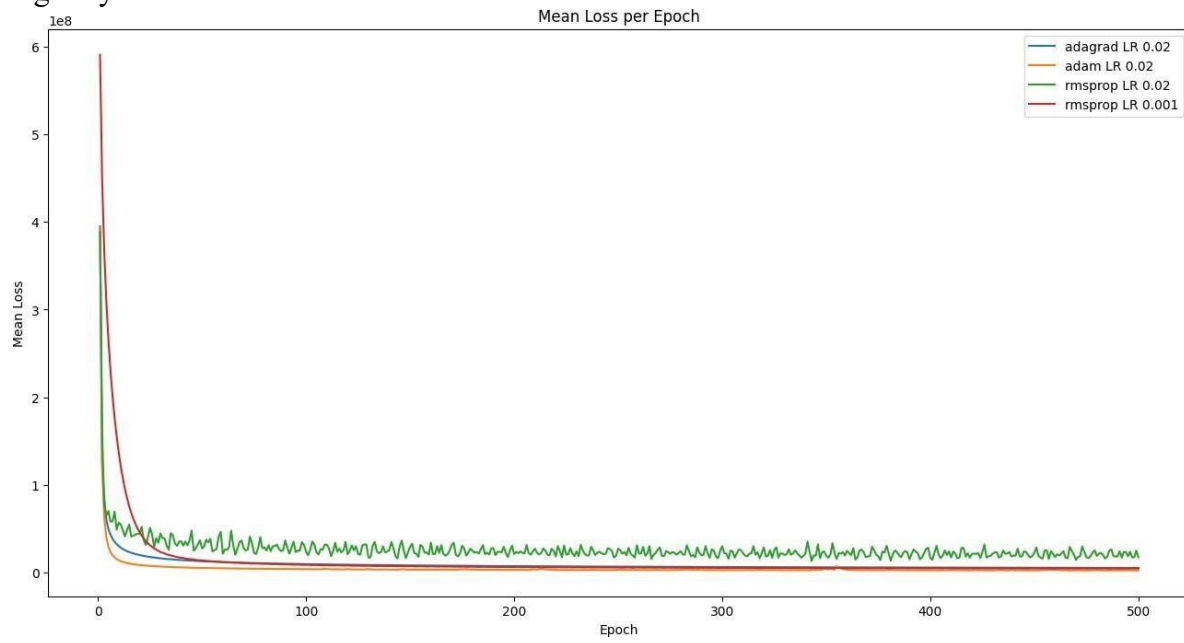


Figure 10: plot loss for Adagrad, RMSprop and Adam with different LR's and epochs=500

Chapter 6: Important Code Snippets

VGG layers:

```
def vgg_layers(layer_names):
    print(layer_names)

    vgg = tf.keras.applications.VGG19(include_top=False,
                                      weights='imagenet')
    vgg.trainable = False

    outputs = [vgg.get_layer(name).output for name in layer_names]
    model = tf.keras.Model([vgg.input], outputs)
    print(model.summary())

    return model
```

Gram matrix:

```
def gram_matrix(input_tensor):
    result = tf.linalg.einsum('bijc,bijd->bcd', input_tensor,
                              input_tensor)
    input_shape = tf.shape(input_tensor)
    num_locations = tf.cast(input_shape[1]*input_shape[2], tf.float32)
    return result/(num_locations)
```

Train step function:

```
@tf.function()
def train_step(image, opt):
    with tf.GradientTape() as tape:
        outputs = extractor(image)
        loss = style_content_loss(outputs)

    grad = tape.gradient(loss, image)
    opt.apply_gradients([(grad, image)])
    image.assign(clip_0_1(image))

    return loss, image
```

Run optimizer function:


```
def run_optimizer(opt, image, epochs, steps_per_epoch):
    losses = []
    result_image= None
    start = time.time()
    step = 0
    for n in range(epochs):
        epoch_losses=[]

        for m in range(steps_per_epoch):
            step += 1
            loss, result_image = train_step(image,opt)
            epoch_losses.append(loss)
            print(".", end='', flush=True)
            display.clear_output(wait=True)
            display.display(tensor_to_image(image))
            print("Train step: {}".format(step))
            epoch_mean_loss = tf.reduce_mean(epoch_losses)
            losses.append(epoch_mean_loss)

        end = time.time()
        print("Total time: {:.1f}".format(end-start))
    return losses, result_image
```

Running Adam optimizer with 50 epochs and LR=0.02:

```
optimizer = tf.keras.optimizers.Adam(learning_rate=0.02)

content_image = image
image_variable = tf.Variable(content_image)

adam_losses_50_lr1, adam_50_lr1=run_optimizer(optimizer,
image_variable, epochs=50, steps_per_epoch=5)
```

Chapter 7: Conclusion and Future Work

In conclusion, Neural Style Transfer represents a powerful and versatile technique for creating artistic images by combining content and style features from different images. Optimizers played a crucial role in training neural networks for Neural Style Transfer. By adapting the learning rates and updating parameters efficiently, optimizers help in accelerating convergence and improving the quality of stylized images.

In conclusion, this report has provided a detailed exploration of Neural Style Transfer, covering its background, proposed system architecture, code implementation, optimization techniques, model architecture explanation, and performance comparison. This detailed discussion provides insights into the working of optimizers in the context of NST, which can help people interested in NST make informed decisions when designing and optimising neural style transfer systems. We found VGG to be faster specifically for image style transfer. Using VGG19 we were able to get more features.

Future research directions for NST include exploring novel loss functions, investigating alternative architectures, improving computational efficiency, and extending the application domain to video and three-dimensional (3D) content where there are more than three loss functions.

Chapter 8: References

- [1] J. Liao, "A Study on Neural Style Transfer Methods for Images," 2022 2nd International Conference on Big Data, Artificial Intelligence and Risk Management (ICBAR), Xi'an, China, 2022, pp. 60-64, doi: 10.1109/ICBAR58199.2022.00019. keywords: {Deep learning;Training;Big Data;Rendering (computer graphics);Generative adversarial networks;Image filtering;Risk management;neural style transfer;deep learning;convolutional neural network;image rendering;generative adversarial nets;texture synthesis}

- [2] Deshmane, Aishwarya. (2023). Image Style Transfer using Neural Network.

- [3] L. A. Gatys, A. S. Ecker and M. Bethge, "Image Style Transfer Using Convolutional Neural Networks," 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, USA, 2016, pp. 2414-2423, doi: 10.1109/CVPR.2016.265. keywords: {Image reconstruction;Neural networks;Image representation;Semantics;Neuroscience;Feature extraction;Visualization}

- [4] Y. Tao, "Image Style Transfer Based on VGG Neural Network Model," 2022 IEEE International Conference on Advances in Electrical Engineering and Computer Applications (AEECA), Dalian, China, 2022, pp. 1475-1482, doi: 10.1109/AEECA55500.2022.9918891. keywords: {Training;Electrical engineering;Computer vision;Computational modeling;Image processing;Neural networks;Production;Image Style Transfer;Convolutional Neural Network;VGG-16;VGG-19}

- [5] A. Ratra, A. Agarwal, V. Sharma, S. Vats, S. Singh and V. Kukreja, "A Comparative Analysis of Fast-style Transfer and VGG19-GramMatrix Approach to Neural Style Transfer," 2023 Second International Conference on Augmented Intelligence and Sustainable Systems (ICAISS), Trichy, India, 2023, pp. 25-30, doi: 10.1109/ICAISS58487.2023.10250595. keywords: {Geometry;Analytical models;Neural Style Transfer;Visual Geometry Group;Gram Matrix;Fast-style transfer;TensorFlow-hub model;Image stylization;deep learning;Image processing;Image analysis;Image synthesis;Image transformation;Comparative study;Performance evaluation}

- [6] H. Ye, W. Liu and Y. Liu, "Image Style Transfer Method Based on Improved Style Loss Function," 2020 IEEE 9th Joint International Information Technology and Artificial Intelligence Conference (ITAIC), Chongqing, China, 2020, pp. 410-413, doi: 10.1109/ITAIC49862.2020.9338927. keywords: {Image quality;Image texture;PSNR;Neural networks;Feature extraction;Distortion;Information technology;image style transfer;gram matrix;neural network;feature extraction}

- [7] Dey, Sandipan. (2018). Hands-on Image Processing in Python, Publisher(s): Packt Publishing, ISBN: 9781789343731