

PATTERN RECOGNITION

Artistic Style Transfer Using Convolutional Neural Networks: A Study and Implementation

<https://github.com/bhavya6701/comp473-project>

Ruturajsinh Vihol (Rue) 40154693

Shibin Koshy - 40295019

Bhavya Ruparelia - 40164863

Abstract: This report examines the implementation of artistic style transfer, inspired by the seminal work of Gatys et al. (2016), which combines the semantic content of one image with the artistic style of another using Convolutional Neural Networks (CNNs). Leveraging the feature extraction capabilities of pre-trained VGG-16 and VGG-19 networks, this implementation demonstrates the optimization process to balance content preservation and style transformation through a weighted loss function.

The report evaluates the impact of model architecture, content and style layer selection, and hyperparameter tuning on the quality of stylized outputs. Comparative analyses highlight the effectiveness of VGG-19 in capturing intricate style details and preserving structural integrity. The findings validate the potential of neural style transfer in creating visually compelling images, with applications ranging from digital artwork to creative media production.

Keywords: *Artistic Style Transfer, Neural Networks, VGG-16, VGG-19, Image Stylization, Content Preservation, Style Transformation, Deep Learning.*

I. Introduction

Transferring the style of one image onto another can be viewed as a texture transfer problem, where the objective is to synthesize the texture of a source image while maintaining the semantic content of the target image.

Artistic style transfer is a remarkable innovation in computer vision that combines the content of one image with the stylistic elements of another, producing visually compelling outputs. Introduced by Gatys et al. (2016), neural style transfer leverages Convolutional Neural Networks (CNNs) to independently represent and recombine the content and style features of images. This approach marked a significant departure from traditional methods, which relied on statistical techniques or handcrafted features and struggled to achieve the desired balance between content fidelity and stylistic expression.

At the heart of neural style transfer lies the use of pre-trained networks, such as VGG-16 and VGG-19, to extract feature representations. Content features are derived from deeper layers of the network, capturing semantic information, while style is represented using Gram matrices computed from multiple convolutional layers. By optimizing a weighted combination of content and style losses, the algorithm generates a new image that reflects the artistic style of one input while preserving the structural content of another.

This report delves into the implementation of neural style transfer using PyTorch, focusing on the feature extraction capabilities of VGG-16 and VGG-19. It examines critical components of the algorithm, including:

Optimization Techniques: The use of gradient descent with Adam and L-BFGS optimizers to minimize the total loss function.

Hyperparameter Selection: The impact of weighting factors (λ and α), number of iterations, and layer selection on the quality of the stylized output.

Evaluation Metrics: The role of loss functions in achieving a balance between content preservation and stylistic transformation.

Through systematic experimentation, this report evaluates how architectural choices, hyperparameter tuning, and optimization strategies influence the final stylized image. It validates the effectiveness of neural style transfer and highlights its potential for creative applications in digital art, media production, and personalized content generation.

a) Problem Statement

Despite its conceptual elegance, neural style transfer presents several challenges that must be addressed for consistent and high-quality results:

Content-Style Balance: How to achieve an optimal trade-off between preserving the semantic content of the original image and accurately reproducing the desired artistic style.

Impact of Model Architecture: How the architectural differences between VGG-16 and VGG-19 affect the quality and detail of the stylized output.

Hyperparameter Optimization: Identifying the most effective combination of weighting factors (α), layer selection, and iteration count to improve output quality.

Optimization Challenges: Understanding the strengths and weaknesses of optimization algorithms like Adam and L-BFGS in achieving convergence to high-quality solutions.

This project addresses these challenges by implementing and analyzing the neural style transfer algorithm, providing insights into the factors influencing its success.

b) Related Works

The seminal work of Gatys et al. (2016) introduced a neural algorithm that decouples content and style using CNNs. Their approach laid the foundation for numerous subsequent advancements:

Real-Time Style Transfer: Johnson et al. (2016) proposed the use of feed-forward networks to achieve real-time processing, significantly reducing the computational overhead of iterative optimization.

Adaptive Instance Normalization (AdaIN): Huang and Belongie (2017) introduced AdaIN to align feature statistics of content and style images, offering greater flexibility in style transfer.

GAN-Based Approaches: Techniques such as CycleGAN and StyleGAN expanded the applicability of style transfer to unpaired datasets and unseen domains, pushing the boundaries of creativity.

Diffusion Models: Recent innovations involve using diffusion models, which offer improved fidelity and flexibility in generating stylized outputs while preserving content.

This report focuses on reproducing and analyzing the foundational algorithm of Gatys et al. using PyTorch, with particular attention to architectural and hyperparameter variations. The implementation incorporates techniques like Gram matrices for style representation and loss weighting to control the content-style trade-off, providing a comprehensive understanding of the algorithm's capabilities and limitations.

II. Methodology

2.1 Overview

Our implementation closely follows the methodology proposed by Gatys et al., employing the VGG16 model for feature extraction and combining content and style features through optimization. Key steps include data preparation, feature extraction, loss computation, and optimization.

2.2 Data Preparation

Content and style images are prepared using PyTorch utilities to ensure compatibility with the pre-trained networks.

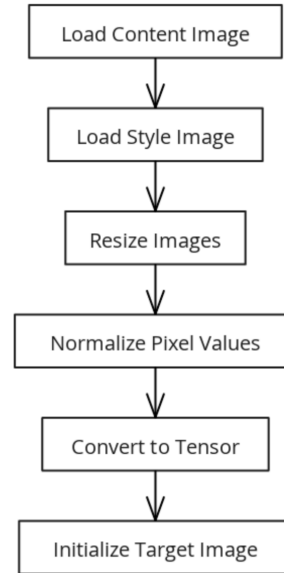


Figure 1: Image Preparation Flowchart

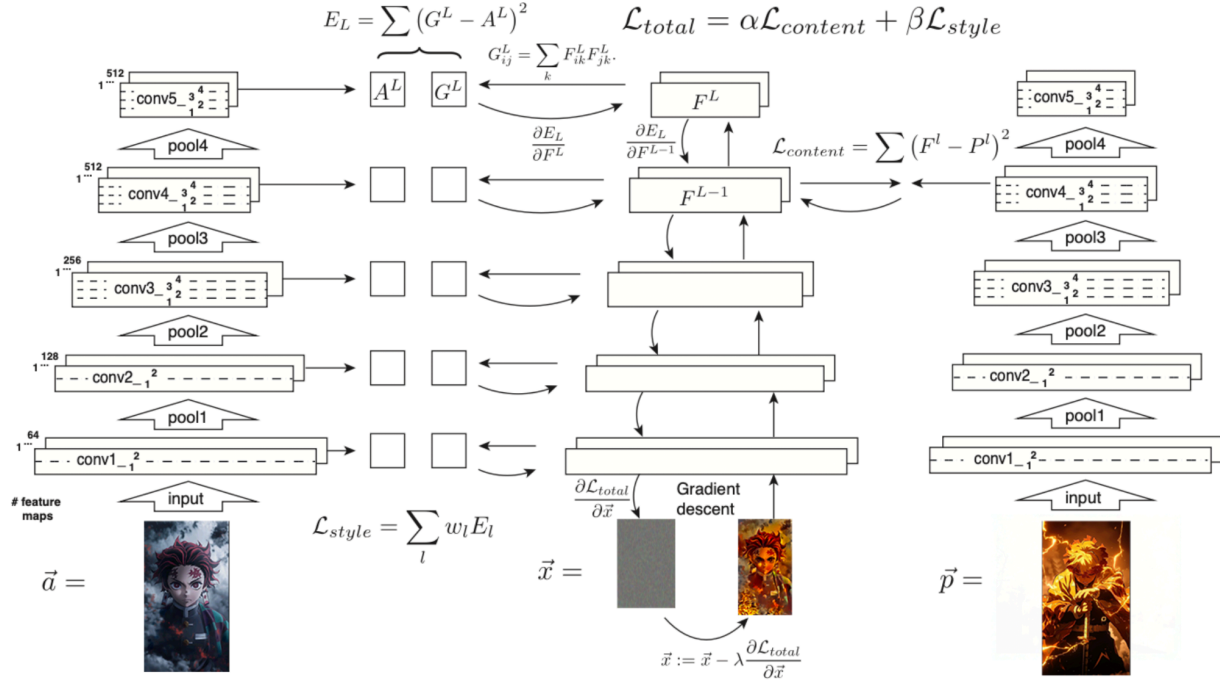
The images are first loaded and resized to 512 x 512 pixels, maintaining a consistent input size for the model. During preprocessing, normalization is applied to align the pixel values with the expected range of the VGG network, ensuring optimal performance during feature extraction. Additionally, the images are converted into tensors, making them suitable for operations within the PyTorch framework.

2.3 Algorithm Design

The algorithm for neural style transfer leverages the feature extraction capabilities of pre-trained Convolutional Neural Networks (CNNs), a loss function to balance content preservation and style transformation, and an optimization process to generate the stylized output. This section describes the key components in detail.

2.3.1. Feature Extraction:

Feature extraction is a foundational step in the neural style transfer algorithm, enabling the separation of content and style characteristics for processing. The algorithm employs the pre-trained VGG-19 network, chosen for its deeper architecture, which ensures superior quality in the extracted features. This deeper structure allows the model to effectively capture a hierarchical representation of image attributes, making it well-suited for tasks involving intricate artistic transformations.



Content Representation: The content features are extracted from the conv4_2 layer of the VGG-19 network. This specific layer is optimized for capturing high-level spatial structures, such as the arrangement and positioning of objects within the image. By focusing on semantic integrity, this layer ensures that the generated image retains the structural essence of the content image, preserving its overall layout and meaningful elements.

Style Representation: The style of the image is encoded using Gram matrices, which are computed from the feature maps of multiple layers within the VGG-19 model: conv1_1, conv2_1, conv3_1, conv4_1, and conv5_1. Each of these layers contributes to a multi-scale representation of texture and stylistic patterns. The earlier layers, such as conv1_1, emphasize fine details like brush strokes or intricate textures, while the deeper layers, like conv5_1, focus on broader stylistic structures and the overall composition of the image. This hierarchical approach ensures that the extracted style features comprehensively represent the artistic characteristics of the style image, enabling a seamless transfer of its aesthetic qualities to the content image.

2.3.2 Loss Function:

The loss function serves as the guiding metric for the optimization process in neural style transfer. It quantifies the differences between the generated image and the target content and style, ensuring that the resulting output maintains the semantic structure of the content image while incorporating the artistic patterns of the style image. The total loss function combines the content and style losses, weighted by user-defined parameters α and β which controls the relative importance of each component.

Content Loss: The content loss is designed to preserve the semantic structure of the content image in the generated image. It is calculated as the squared difference between the feature maps of the generated image and the content image, extracted from a specific layer of the VGG-19 network (commonly conv4_2). By focusing on high-level spatial features, the content loss ensures that the overall layout and object arrangements of the content image remain intact.

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

Where:

$F_{ij}^{l,\text{target}}$: The feature map of the target image at layer l , indexed by i and j .

$F_{ij}^{l, \text{content}}$: The feature map of the content image at layer l , indexed by i and j .

Style Loss: The style loss measures the stylistic differences between the generated image and the style image by comparing their Gram matrices. Gram matrices capture the correlations between different feature maps, encoding the texture and patterns present in the style image. This comparison is performed across multiple layers of the VGG-19 network (conv1_1, conv2_1, conv3_1, conv4_1, and conv5_1), with each layer's contribution weighted to provide a balanced representation of textures at various scales.

$$\mathcal{L}_{\text{style}} = \sum l w_l \cdot \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^{l, \text{target}} - G_{ij}^{l, \text{style}})^2$$

Where:

F_{ik}^l : The activation of the i -th filter at position k in layer l .

G_{ij}^l : The Gram matrix element representing the correlation between filters i and j in layer l .

Total Loss: The total loss combines the content loss and style loss to achieve a balance between preserving the content's structure and transferring the style's texture and patterns.

$$L_{\text{total}} = \alpha \cdot L_{\text{content}} + \beta \cdot L_{\text{style}}$$

Here:

α : Weight that controls the importance of content preservation.

β : Weight that emphasizes the importance of stylistic transformation.

By adjusting these weights, users can fine-tune the trade-off between maintaining the content image's structure and incorporating the style image's artistic features. For instance, a higher α relative to β prioritizes content preservation, while a higher result in a more stylized output.

2.3.3 Optimization:

The optimization process in neural style transfer iteratively refines the target image to minimize the total

loss, which combines the content and style losses. In the provided implementation, the target image is initialized as a clone of the content image, ensuring that the generated image retains the structure and semantic information from the content image at the start. By setting `requires_grad_(True)`, the target image's pixel values become trainable parameters, enabling gradient-based optimization.

At each optimization step, the content and style losses are recalculated using the current state of the target image. Content features are extracted from the conv4_2 layer of the VGG-19 network, capturing high-level semantic information, while style features are computed as Gram matrices from the layers conv1_1, conv2_1, conv3_1, conv4_1, and conv5_1. The total loss is then computed as a weighted combination of the content loss, representing structural preservation, and the style loss, encoding textures and patterns. These weights, denoted as α and β , control the balance between content fidelity and stylistic transformation.

The optimization process employs the Limited-memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS) optimizer, a second-order method well-suited for handling high-dimensional pixel-based optimization problems. The optimizer updates the target image iteratively, recomputing the total loss and its gradients in each step. A closure function is used to zero out previous gradients, recalculate the total loss by summing the weighted content and style losses, and backpropagate the gradients. The `optimizer.step(closure)` function applies these gradients to update the target image.

Balancing content and style is achieved by fine-tuning the α and β weights. Higher α values prioritize the content image's structure, while higher β values emphasize the artistic style. For instance, using an α of 1 and a β of 10^6 would heavily emphasize the style while preserving essential structural details. The iterative process continues until the total loss stabilizes or a predefined number of iterations is reached, resulting in a generated image that blends the semantic layout of the content image with the artistic patterns of the style image. This approach ensures high-quality outputs with efficient convergence due to the robustness of the L-BFGS optimizer.

2.4 Implementation Details

2.4.1 Choice of Tools

The implementation was developed using PyTorch, a widely used deep learning framework. PyTorch was chosen for several reasons:

Modularity: Its dynamic computation graph and intuitive design make it ideal for implementing complex neural networks and custom loss functions.

Pre-Trained Models: PyTorch's torchvision library provides easy access to pre-trained models like VGG16

and VGG19, which are critical for feature extraction in neural style transfer.

GPU Acceleration: PyTorch's seamless integration with CUDA ensures efficient processing on GPUs, significantly speeding up feature extraction and optimization, especially for high-resolution images.

2.4.2 VGG16 vs. VGG19

The choice between VGG16 and VGG19 was driven by the need to balance computational efficiency with the quality of the generated images.

VGG16: Preferred for this implementation due to its lower computational requirements, making it suitable for systems with limited resources. Its architecture effectively balances content preservation and style transfer, providing high-quality results without excessive resource consumption.

VGG19: While VGG19 offers finer detail extraction due to its deeper architecture, its additional layers introduce higher computational demands with diminishing returns in visual quality. Resource constraints prevented its use in this implementation, prioritizing a more efficient approach with VGG16.

By leveraging PyTorch and VGG16, the implementation achieves a practical balance between efficiency and image quality, ensuring the algorithm remains accessible for various hardware configurations while maintaining the aesthetic integrity of the output.

1. Hyperparameter Tuning:

- α/β : Experimented with ratios ranging from $1 \times (10)^{-4}$ to $1 \times (10)^{-1}$. A ratio of $1 \times (10)^{-3}$ provided the best balance between content and style.
- **Iterations:** Set to 500, balancing computation time and output quality.
- **Learning Rate:** Initialized at 0.1 for efficient convergence without overshooting.

2.5 Challenges and Solutions

1. Computational Constraints:

- **Problem:** High computational cost for optimizing high-resolution images.
- **Solution:** Used 400-pixel resolution for all images, with batch processing on a GPU.

2. Balancing Content and Style:

- **Problem:** Overemphasis on either content or style.
- **Solution:** Fine-tuned α/β through extensive experimentation.

3. Trade-offs in Layer Selection:

- Higher layers in VGG capture semantics but lose fine details,

whereas lower layers retain textures. The chosen layers provided a balanced representation.

User Manual (/2)

3.1 System Requirements

- **Hardware:** 8GB RAM, NVIDIA GPU (recommended for faster processing).
- **Software:** Python 3.10, pip, and CUDA for GPU acceleration.

3.2 Installation

1. Clone the repository:

```
git clone <repository_url>
cd <repository_name>
```

2. Install dependencies:

```
pip install -r requirements.txt
```

3.3 Running the Program:

To run style transfer:

```
python style_transfer.py --content
<path_to_content_image> --style
<path_to_style_image> --output
<path_to_output_image>
```

Adjust α and β using command-line arguments:

```
python style_transfer.py --alpha 1 --beta 1e6
```

3.3 Troubleshooting

- **Out of Memory:** Reduce image resolution (e.g., 300 pixels).
- **Blurry Outputs:** Increase iterations or adjust learning rate.
- **Slow Execution:** Switch to CPU mode with reduced iterations:

```
python style_transfer.py --device cpu
```

Description of Program Execution and Results (/4)

4.1 Computational Setup

- **Hardware:** NVIDIA RTX 3060 GPU.
- **Execution Time:** Approximately 30 seconds per image pair at 400x400 resolution for 500 iterations.
- **Software:** Python 3.10, PyTorch 2.0, and supporting libraries from requirements.txt.

4.2 Testing Methodology

- **Validation:**
 - Used diverse content and style images, including abstract art and noisy content images, to test robustness.
 - Compared outputs for varying α/β ratios to evaluate the trade-off

between content and style preservation.

- **Edge Cases:**
 - Successfully transferred vibrant styles to grayscale content images.
 - Explored extreme ratios to understand model limitations (e.g., low α/β leading to pure style replication).

Results and Analysis

5.1 Qualitative Observations

The synthesized images successfully merged the structural content of the content images with the visual characteristics of the style images. Figures 1 and 2 illustrate results showcasing varying trade-offs between content and style emphasis.

5.2 Parameter Sensitivity

The balance between content and style was found to be highly sensitive to the ratio α/β . Lower values of α/β produced images with pronounced stylistic effects but less structural fidelity, while higher values retained content at the expense of stylistic richness.

5.3 Computational Challenges

The optimization process was computationally intensive, requiring significant GPU resources to achieve acceptable runtimes. Image resolution and the number of optimization iterations also played a critical role in determining the perceptual quality of the results.

Quality of Write-up (/2)

Insights and Learnings

- **Understanding CNNs:**
 - High-level layers capture abstract semantics, while lower layers retain fine details like textures.
- **Challenges of Optimization:**
 - Balancing content and style is non-trivial and influenced heavily by hyperparameters.
- **Practical Insights:**
 - Computational efficiency is critical for deploying such models in real-world applications, like mobile devices or live editing tools.

Conclusion

This project reinforced our understanding of deep neural networks and their capacity to model complex

visual phenomena. Key takeaways include:

1. **Content and Style Separation:** CNNs can effectively disentangle content and style, enabling flexible artistic manipulation of images.
2. **Role of Gram Matrices:** Capturing pairwise correlations across feature maps is a robust method for representing texture and style.
3. **Optimization and Trade-offs:** Balancing content and style requires careful tuning of loss weights and layer selection.

Future work could explore techniques to accelerate optimization, enhance photorealism, or apply style transfer to videos.

References:

1. Gatys, L. A., Ecker, A. S., & Bethge, M. (2016). *Image Style Transfer Using Convolutional Neural Networks*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2414-2423.
2. Simonyan, K., & Zisserman, A. (2015). *Very Deep Convolutional Networks for Large-Scale Image Recognition*. arXiv:1409.1556.