
Interpretable Neural Style Transfer using VGG Architectures

Siddhant Joshi¹ William Hu¹

Abstract

Rendering the semantic content of an image in various styles poses a unique image processing challenge that has recently garnered significant attention in the field of computer vision (CV). Gatys *et al.* published one of the first papers concerned with this application of deep learning, dubbed Neural Style Transfer (NST), opening many avenues for further research. In this paper, we will primarily focus exploration into the algorithm's interpretability, providing our own survey of experiments with varying degrees of key parameters involved in a classic VGG-19-based NST algorithm. Our objective is twofold: to gain deeper insights into the algorithmic NST decision-making process and to design a new algorithm that enhances fidelity and content preservation in stylized images. Not only is our proposed NST algorithm able to achieve lower total losses than the baseline, but also produce more visually appealing stylized images with a simple feed-forward neural network.

1. Introduction

Gatys *et al.* recently approached this task using deep neural networks, showing that they can not only encode the content but also the *style* information of images^[7]. Using this concept, they published their cornerstone work on NST utilizing VGG-19 to effectively separate content and style information from images and combine these representations into new, stylized images. Since then, subsequent research focused on improving the efficiency and speed as well as enabling more customization into the stylistic focuses of the algorithms^{[1][4][5][9]}.

Before understanding the goal of our study, we must establish an intuition behind the original NST algorithms. As described in [7], NST can be reduced down to an optimization problem involving four components: a content image (I_C), a style image (I_S), an input image (I_i), and a pre-trained

feed-forward convolutional neural network (CNN). All three images are fed into the CNN, such as VGG-19, to extract characteristic features. Generally, successive convolutional layers within a network capture increasingly complex features, abstracting pixel-wise information to representations that capture the high-level content of an image. Stylistically, various relevant information are encoded throughout the feature maps of each layer. To effectively capture the dispersed style information, [7] suggest considering the feature maps of every convolutional layer and combining their information into Gram matrices to allow for holistic consideration of lower and higher-level styles. As such, extracting content information from deeper convolutional layers would be useful for content preservation whereas consulting Gram matrices would be useful for extracting style. These "feature summaries" are then compared to those of I_i via a squared-loss calculation, resulting in a content loss (\mathcal{L}_C) and a style loss (\mathcal{L}_S) that will be used in gradient descent to update the pixels of the I_i to better represent the style and content information. This entire process is repeated over multiple transfer iterations, eventually converging toward a final stylized output image (I_O). Fig. 1 presents a full diagram of the algorithm.

Our study is concerned with the composition of this algorithm, specifically in unveiling the relationships between adjustable parameters and their individual impact on the quality of style transfer. Additionally, we explore error metrics against subjective visual appeal of the output images, seeking potential for a reconciliation between the quantitative and qualitative aspects of this application.

2. Methodology

As in Gatys *et al.*'s original study, our experiments will be generated with the VGG-19 network, which was trained to perform object recognition on the ImageNet dataset^[2]. We will remove the fully-connected (FC) layers and final softmax activation layer from the model, isolating the pre-trained convolutional blocks for the algorithm. Additionally, we convert every MaxPool2d layer into AvgPool2d layers as average pooling was found to produce higher quality outputs in previous studies^[7].

¹Department of Cognitive Science and Halıcıoğlu Data Science Institute, University of California, San Diego, La Jolla.

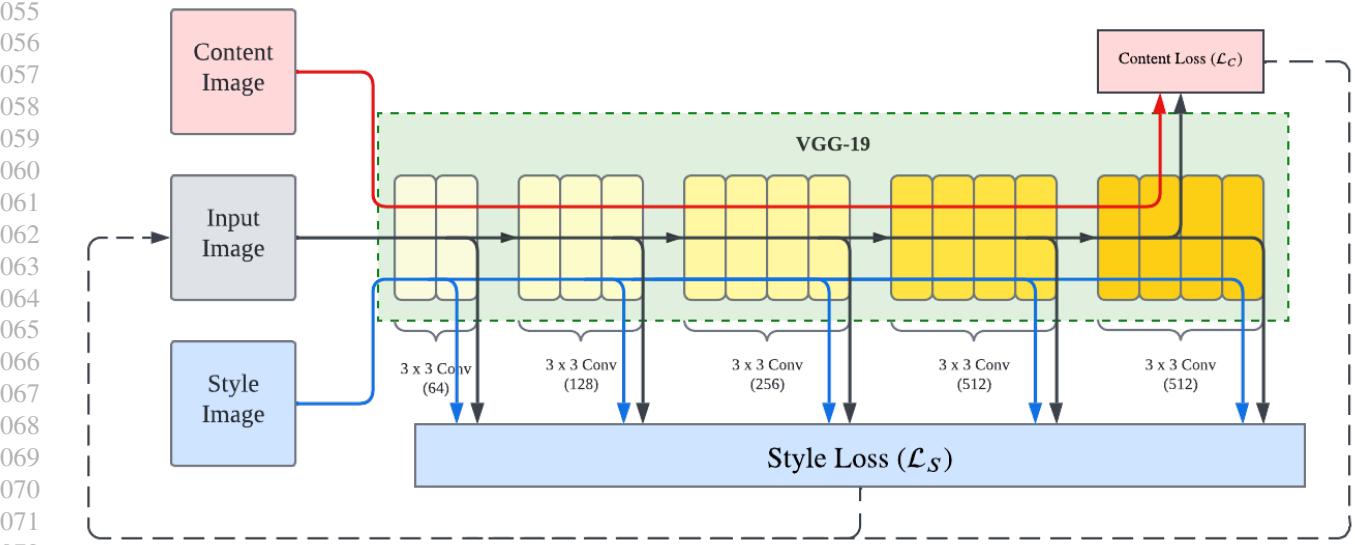


Figure 1. Input, content, and style images are all passed into the VGG-19 network where feature maps from convolutional layers are pulled out and converted into their respective losses. Both \mathcal{L}_C and \mathcal{L}_S are fed into the input image for pixel-wise updating via gradient descent.

The following table describes the parameters of the baseline model for this study.

Parameters	
Number of Iterations	300
Content Error Layer	conv_4
Style Error Layers	conv_1
	conv_2
	conv_3
	conv_4
	conv_5
Content Weight	1
Style Weight	1000000
Loss function	L2
Optimizer	L-BFGS

Table 1. Parameter Settings

To limit computational complexity, we will study the effect of each parameter of interest in isolation of the others, assessing performances across each independent parameter and recording total loss of the stylized image. Furthermore, we will be maintaining the same content and style images throughout our experiments to ensure results are quantitatively and qualitatively comparable (Fig. 2). The following subsections explain the parameters we will be testing and provide rationale for design choices.

2.1. Number of Transfer Iterations

Current literature does not seem to have specific preferences for the number of transfer iterations. As such, we will



Figure 2. Baseline I_C (left) and I_S (right).

arbitrarily assess performance on transfer iterations from the range {100, 200, 300, 600, 1000, 2000}.

2.2. Error Layers

When it comes to assessing transfer quality, the choice of layer with which to construct Gram matrices is critical. Aside from convolutional layers, we will also be assessing representation quality in ReLU layers within the VGG-19 network. This would be the most computationally intensive aspect of our study, as we will be assessing all 31 combinations of convolutional and ReLU layers between both style and content, resulting in 124 total trials.

2.3. Error Weights

Mathematically, the baseline losses for style are scaled to be comparable to that of the content loss. Since the key behind the "learning" in NST relies on the magnitudes of

110 these losses, adjusting their scaling factors may prioritize
 111 different degrees of content and style in the stylized image.
 112 We assessed styles weights in the range {1000000, 100000,
 113 1000, 100} and content weights from the range {1, 5, 10,
 114 20, 100}.

115 2.4. Stochastic Focus

116 We also decided to delve into algorithm generalizability by
 117 leveraging some common image sampling methods used in
 118 image classification to prevent overfitting. Currently, the
 119 baseline algorithm considers the entire style image to update
 120 a copy of the content image. We introduce randomness into
 121 the algorithm to curb these potential issues via the following
 122 methods:
 123

124 2.4.1. RANDOM STYLE SHIFTING

125 In our current process, each transfer iteration is followed by
 126 a comparison of pair-wise similarities between corresponding
 127 pixels. As such, each pixel is at risk of overfitting to
 128 single-pixel value information it receives. We add variability
 129 to this process by randomly cropping and rotating the I_i
 130 with Kornia library before processing similarity. In doing
 131 this, we aim to have the algorithm fit each output pixel to
 132 a neighborhood of style pixels, thereby generalizing style
 133 throughout the I_i . Since these modifications are to be slight,
 134 we decided to consider transformation dimensions of at least
 135 95% of the size of the output image. More specifically, from
 136 the range {0.95, 0.96, 0.97, 0.98, 0.99}.

137 2.4.2. I_i INITIALIZATION

138 We also believe the choice of the initial I_i input plays a role
 139 in the quality of the output image. We will test the following
 140 three options for I_i :
 141

- The content image (I_C) [Baseline]
- 512×512 image of Gaussian-random pixel values
- 512×512 image with all pixel values set to 0 (black)

142 2.5. Style Image Type

143 All images contain style and content aspects, so gauging
 144 the separability of these components seems logical. We can
 145 broadly classify major art works into one of landscapes or
 146 portraits. As such, we will test the baseline on combinations
 147 of images from these classes to get a better understanding
 148 of how NST style extraction ignores content and vice
 149 versa. We introduce the following 4 images, two landscapes
 150 and two portraits, and run the algorithm on every pairwise
 151 combination (Fig. 3 & 4).

152 2.6. Loss Function

153 Aside from the baseline L2 error, we also assessed the
 154 performance of L1 error as it is generally preferred for its
 155



Figure 3. I_C Images



Figure 4. I_S Images

robustness to overfitting. In this context, we anticipate it to counter the output image overfitting to either style or content information.

2.7. Optimizer

The baseline algorithm uses a Limited-memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS) optimizer. We will conduct tests using Adam, Adagrad, and standard SGD optimizers from PyTorch, as they are preferred by research when handling image processing tasks.

3. Experimental Results and Discussion

3.1. Baseline Performance

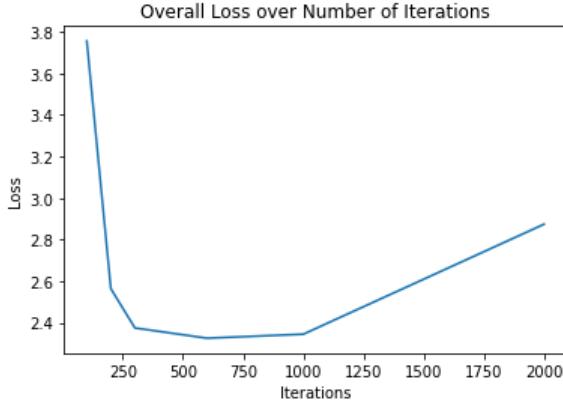
For reference, we generated a baseline output image that achieved the following metrics.



Figure 5. $\mathcal{L}_C = 2.227, \mathcal{L}_S = 0.148$

165 3.2. Number of Transfer Iterations

166 After running the baseline model with the proposed transfer
 167 iterations, we yielded the following loss curve (Fig. 6).
 168



184 *Figure 6.*

185 This loss curve is convex, achieving a minimum around
 186 600 transfer iterations. The parabolic shape seems to be
 187 indicative of the algorithm overfitting to the content and
 188 style images, achieving convergence at 600 iterations.
 189

190 3.3. Error Layers

191 Based on loss metrics alone, the ReLU combinations for
 192 both \mathcal{L}_C and \mathcal{L}_S layers consistently outperformed their
 193 convolutional counterparts. Visually, extracting content and
 194 style information from ReLU layers generated extremely
 195 similar results to that of the convolutional layers.
 196



200 *Figure 7. Conv_1 error layer (left) and ReLU_1 error layer (right)*

201 One interesting thing of note however, was that the earlier
 202 ReLU layers (ReLU_1 and ReLU_2) achieved the lowest
 203 overall content loss of all ReLU trials. In theory, We
 204 expected that deeper ReLU layers (ReLU_3 and ReLU_4)
 205 would have performed better as their feature maps would
 206 be capturing more abstract details from I_C , however that
 207 does not seem to be the case, warranting further tuning and
 208 discussion.
 209

210 3.4. Error Weights

211 Below are some figures of the output with differing style
 212 weights (Fig. 8).
 213



214 *Figure 8. Output images of decreasing style weights: 1000000
 215 (top-left), 100000 (top-right), 1000 (bottom-left), 100 (bottom-right)*

216 Below are some figures of the output with differing content
 217 weights (Fig. 9).
 218



219 *Figure 9. Output images of decreasing content weights. From left
 220 to right: 100, 10, 1.*

221 In both series of experiments, decreasing one weight would
 222 increase the prevalence of the features captured by the other
 223 weight. For example, decreasing the content weight brought
 224 out more of the soft, geometric style of I_S in the output
 225 image. With lower relative style weights, the realism of
 226 the dancer was prioritized, even resulting in an output that
 227 maintained the color scheme of I_C (Fig. 8, bottom-right)

228 Quantitatively, lower weights in either style or content re-
 229 sulted in decreased overall loss. This makes intuitive sense
 230 as the algorithm is focused on optimizing one aspect more
 231 than the other, which is an easier problem to solve than
 232 attempting to minimize the loss in both content and style
 233 simultaneously. However, this behavior highlights the dis-

220 parity between visual quality and quantitative improvements,
 221 specifically that decreasing loss does not necessarily corre-
 222 late to well-stylized outputs. We will discuss this further in
 223 following sections.
 224

225 3.5. Stochastic Focus

226 3.5.1. RANDOM STYLE SHIFTING

228 In testing on the range of values for style shifting, we
 229 recorded the following losses:
 230

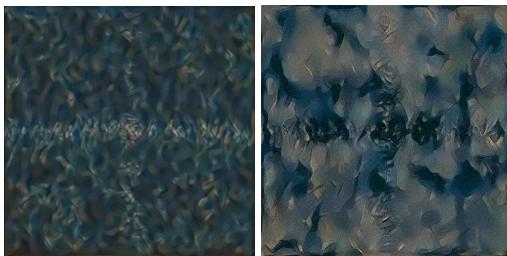
Factor	Total Loss
0.95	36.681
0.96	85.760
0.97	25.836
0.98	131.691
0.99	196.711

232 Table 2. Loss Across Shifting Factors
 233
 234
 235
 236
 237
 238

241 From the table, we can see that the losses varied across
 242 scaling factors with 97% achieving the lowest.
 243

244 3.5.2. I_i INITIALIZATION

246 Aside from the baseline initialization of the input image
 247 with the content image, we assessed an input of Gaussian-
 248 random pixel values and an input of all pixel values set to
 249 zero.
 250



261 Figure 10. Gaussian-random (left) and zero-pixel initialization
 262 (right)
 263

265 Clearly, these output images preserved little to no content
 266 information over the course of 300 transfer iterations. How-
 267 ever, there seems to be an interesting blend of the style and
 268 input images present in either output. This could hint at
 269 a possibility of this particular treatment relying on other
 270 algorithmic changes to be more effective in blending style
 271 with content. Therefore, we will further explore this as a
 272 hyperparameter to assess the viability of our belief in inter-
 273 dependence.
 274

3.6. Loss Functions and Optimizers

We recorded the following results for the different loss functions and optimizers that were tested.

Loss Functions	
L2 (MSE)	2.375
L1 (MAE)	588.243
Optimizers	
L-BFGS	2.375
Adam	33.714
Adagrad	33.184
SGD	37.368

Table 3. Total loss of various loss functions and optimizers

Surprisingly, L1-loss performed significantly worst than L2-loss. Yet, as we saw in Section 3.4, the quality of L1-loss outputs kept/retained/captured more style than that of the L2-loss outputs (Fig. 11) Most notably, the fine details of the figure’s hands and face are less prioritized in the L1 output, an arguably good sign of effective stylization.



Figure 11. L1-loss output (left) compared to L2-loss output (right)

Regarding choice of optimizer, L-BFGS outperformed all other optimizers based on loss alone. However, upon inspection of the output images, L-BFGS appears to have a lower quality transfer (Fig. 12).





275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329

Figure 12. Output images of the different optimizers: L-BFGS (top-left), Adagrad (top-right), Adam (bottom-left), and SGD (bottom-right)

Similar to our interpretations of L1 and L2-loss outputs, the alternative optimizers appear to smooth out details in I_C and bring focus to the overall structure of the dancer, which is preferred given the abstract style image.

3.7. Style Image Type

When assessing the impact of the contents of style images, we analyzed style quality on four different combinations; transferring the style of a cubist portrait onto another portrait and landscape:



Figure 13. Applying the cubist portrait style onto another portrait (left) and a landscape (right)

And transferring the style of an impressionist landscape onto the same portrait and landscape:



Figure 14. Applying the impressionist landscape style onto a portrait (left) and another landscape (right)

Overall, contrasting image types does not seem to negatively impact the visual appeal of the output images. In Fig. 13, we can see textures and colors from the cubist portrait were represented nicely in Messi's tattoo and the greenery beside the Golden Gate Bridge. Furthermore, the fuzziness of the impressionist landscape was applied effectively on both output images, especially the hints of the pink flowers scattered throughout the horizon of the Golden Gate Bridge.

We acknowledge that these outputs could appear more as an application of a filter over the content image, rather than a blending of style and content. Given that the algorithm is still able to extract and separate style and content representations, we believe the distinction between a style *overlay* and a style *transfer* could be explored as an interaction between parameters.

3.8. Discussion

Over the course of the experimentation we conducted with various parameters of the original baseline algorithm, we discovered some notable independent impacts on output images. Most notably, we observed a trade-off between style and content weights that can be leveraged as a hyperparameter to improve transfer quality. In giving more importance to either aspect, the algorithm could be tuned to specific images based on user-defined requirements. A common theme across multiple experiments was that lower total loss did not necessarily result in the best stylized output. This challenge is unique to the use-case for NST, as subjective evaluations of the output cannot be ignored when the goal is to apply a loosely-defined style to an equally ambiguous content. Additionally, our experimental design and observations beg the consideration of parameter interactions. Our experimentation assessed the parameters in isolation, preventing us from taking into account any further positive (or negative) implications from considering them simultaneously. We saw most promise of this concept when considering $I_C - I_S$ contrast and different I_i initializations. It is well-established in machine learning theory that random, unbiased input leads to more reliable results, so to witness the exact opposite in an isolated setting should not be enough to rule out its effectiveness when incorporated with other parameters.

4. Revising the Algorithm

After experimenting with multiple parameters and analyzing the relationships between them, we propose an improved NST algorithm based on some key findings. Our final algorithm will include a random style shifting before each transfer iteration. Since Adam, Adagrad, and SGD all generated similar results, we opted for an Adam optimizer as it is common practice in CV research. Content weights, style weights, and I_i initialization will be treated as hyperparameters and will be tuned with every combination of

330 style and content image. We opted for feature extractions
 331 via the ReLU layers, because despite producing visually
 332 similar results as convolutional layers, the nature of the
 333 operation would make running the algorithm less computa-
 334 tionally intensive. Table 4 describes the fixed parameters of
 335 our proposed revision.
 336

Parameters	
Number of Iterations	600
Content Error Layer	ReLU_5
Style Error Layers	ReLU_1
	ReLU_2
	ReLU_3
	ReLU_4
	ReLU_5
Optimizer	Adam (lr: 0.05)
Loss Function	L1

Table 4. Revised Parameter Settings

4.1. Results

After feeding and tuning our algorithm for various combinations of content and style images, we achieved the following output images (Fig. 19). We also report \mathcal{L}_C and \mathcal{L}_S for the baseline and revised algorithms for a quantitative comparison (Table 5).

Row	Baseline \mathcal{L}_S	Baseline \mathcal{L}_C	Our \mathcal{L}_S	Our \mathcal{L}_C
1	0.265	2.352	3.276	0.916
2	4.370	13.113	1.575	0.933
3	5.811	19.208	4.115	3.816
4	4.707	13.526	2.861	7.440
5	4.560	20.193	10.701	0.879

Table 5. Loss Comparison Across Results

4.2. Discussion

Overall, our revisions significantly reduced total loss when compared to the baseline and produced more visually appealing results. Revisiting our decision to continue investigating different I_i initializations proved beneficial in enhancing transfer quality. For instance, applying the baseline style image onto the dancer with a zero-pixel input image generated a depiction of the dancer built from the irregular polygons present in the style image, with blue streaks and orange highlights. We theorize this was likely due to the algorithm having to reconstruct *both* style and content onto an uninformative input image, encouraging it to better learn the style and incorporate the patterns effectively into pieces of content.

5. Conclusion

In this paper, we conducted a rigorous assessment of the individual impacts of common parameters on an original NST algorithm proposed by Gatys et al. and devised an improved algorithm from our findings^[7]. We discovered a trade-off between style and content weightages, efficient feature extraction via ReLU layers, and influence of input initialization.

In future research directions, we hope to consider experimental designs that take parameter interactions into account. This way, we may glean a better understanding of any potential interplay that may have been missed by our computational limitations. Our study only tested 512×512 images and transformed any nonconforming inputs to these dimensions. Exploring algorithmic adjustments to accommodate photos of higher fidelity may lead to different outputs.

385
386
387
388
389
390
391
392
393
394



395
396
397
398
399
400
401
402
403
404



405
406
407
408
409
410
411
412
413
414



415
416
417
418
419
420
421
422
423
424



425
426
427
428
429
430
431
432
433
434



435
436
437
438
439

. Style

. Content

. Gatys et al.

. Ours

Figure 19. Results Array

440 **Software and References**

441 **Software**

443 All of the code for this report was run in Jupyter notebooks
444 on UCSD's Datahub. The built-in VGG-19 network from
445 PyTorch was used for the NST algorithm implementation.
446 We also utilized Kornia, an open source Python package
447 specialized in computer vision.

448

449 **References**

450

451 [1] A. J. Champandard, "Semantic style transfer and
452 turning two-bit doodles into fine artworks," (2016)
453 <https://arxiv.org/abs/1603.01768>

454

455 [2] A. Krizhevsky, I. Sutskever, and G.E. Hinton, "ImageNet
456 classification with deep convolutional neural networks," (2017) <https://doi.org/10.1145/3065386>

457

458 [3] C. Olah, A. Mordvintsev, and L. Schubert, "Feature
459 Visualization", (2017) <https://distill.pub/2017/feature-visualization/>

460

461 [4] J. E. Kyprianidis, J. Collobosse, T. Wang, and T. Isenberg. "State of the "Art": A Taxonomy of Artistic Stylization Techniques for Images and Video", (2013) doi:
462 10.1109/TVC.2012.160

463

464 [5] J. Johnson, A. Alahi, L. Fei-Fei. "Perceptual Losses
465 for Real-Time Style Transfer and Super-Resolution," (2016)
466 <https://arxiv.org/abs/1603.08155>.

467

468 [6] K. Simonyan, A. Vedaldi, and A. Zisserman. "Deep
469 Inside Convolutional Networks: Visualising Image Classifica-
470 tion Models and Saliency Maps", (2013), arXiv:1312.6034

471

472 [7] L. A. Gatys, A. S. Ecker and M. Bethge, "Image Style
473 Transfer Using Convolutional Neural Networks," (2015) doi:
474 10.1109/CVPR.2016.265.

475

476 [8] V. Dumoulin, J. Shlens, and M. Kudlur. "A
477 Learned Representation for Artistic Style", (2017)
478 <https://arxiv.org/abs/1610.07629>

479

480 [9] X. Huang, S. Belongie, "Arbitrary Style Transfer in
481 Real-time with Adaptive Instance Normalization", (2017)
482 <https://arxiv.org/abs/1703.06868>

483

484

485

486

487

488

489

490

491

492

493

494