

Extensions to Neural Style Transfer

Team 4D Tensor

Praveen Iyer, Swarnita Venkatraman, Sneha Bandi, Neha Chawla

Abstract

Given a content and style image, the process of applying styles to generate an output image that retains the core elements of the content image but appears to be in the style of the reference image is known as Neural Style Transfer (NST) (Gatys et al., 2015a). In recent times, NST has achieved good performance, but studies highlighting improvements in stability, quality, flexibility, and evaluation are still in its early stages, specifically using user feedback. To build upon these challenges, we propose an extension to the original NST technique to enhance the final style of the output image, using the novel pipeline of Cascade Style Transfer (CST). Our key idea focuses on using CST to combine multiple style images, alongside novel quantitative and qualitative evaluation methods to assess the stylistic quality of generated image. Our cascade framework involves using multiple architectures in serial based on traditional Gram Matrices, finally evaluating the stylized output using three quantifiable factors - Content Fidelity (CF), Global Effects (GE) and Local Patterns (LP) and qualitative evaluation.

1 Introduction

Multiple Style Transfer (MST) in Neural Style Transfer (NST) enables styles of multiple images to be transferred to a single content image by modification of style losses. Multi-Objective network (MO) uses multiple style images to optimize multiple objectives, for content preservation and style transfer, simultaneously while Cascade style transfer (CST) is a neural style transfer method that iteratively applies style to a content image to preserve both style and original content, sequentially.

Existing challenges with multiple style transfer include increased training time due to a separate network tied to one style. Adding multiple styles in MO can lessen the impact of each individual style, especially when the styles are somewhat different (Makow and Hernandez, 2017). Higher style

loss exists when blending two images compared to having just one style image (Makow and Hernandez, 2017). On the other hand, CST is shown to improve the stylistic quality and flexibility of the generated image. (Wang et al., 2021b)(Wang et al., 2020) Style from the style image is known to be preserved more in CST, while in traditional MOs, content from the content image is known to be preserved more. (Wang et al., 2021b) Exploring possibilities of performing multiple style transfer using CST to exploit its benefit of preserving style gives us room to incorporate multiple styles. Besides, existing literature on using CST for Multiple Style Transfer has not been found yet. Two style images might have contradicting styles that can lead to unsatisfactory output. The ordering of style images in the CST will also matter leading to vastly different outputs for different ordering of input. There is no ground truth for NST.

Based on the complementary strengths that CST has when compared to the challenges of Multiple Style Transfer using MO, the **hypothesis** of the project was formulated. Hypothesis of the project is, “Using CST for Multiple Style Transfer enables generation of images that incorporate multiple styles with better quality and style preservation”

Evaluating and improving the stylization quality remain two important open challenges, standard quantifiable metric implementation to measure NST performance does not exist. This means that for evaluation we rely on qualitative methods which can be challenging. This can lead to subjective opinions mattering more rather than having a clear objective which is being optimized. Designing more fine grained metrics might also be needed to improve the quality of results.

We evaluate the hypothesis ¹ to understand if the style is known to be preserved more in CST than

¹<https://github.com/praveen-iyer/CSCI566-Project>

081 in traditional MO to enable generation of images
082 with better quality. For testing the above hypothesis
083 we use evaluation metrics Content Fidelity, Global
084 Effects, and Local Patterns along with a qualitative
085 user study.

086 2 Related Work

087 In recent times, though NST has achieved good
088 performance, studies highlighting improving its
089 stability, quality, flexibility, and evaluation are still
090 in its early stages. There is a dearth of literature
091 on quantitative metrics for evaluating the quality
092 of generated images through NST, especially when
093 using CST with multiple styles.

094 Traditional Convolutional Neural Networks
095 (CNN) hierarchically stack to form content rep-
096 resentations. It is possible to form style represen-
097 tations by capturing the correlations of different
098 filter responses within each layer. In (Gatys et al.,
099 2015a) produced the first work investigating the
100 extraction style representations from images. Ad-
101 ditionally, much work has been done on exploring
102 the perceptual possibilities of neural style, where
103 novel extensions to the original neural style algo-
104 rithm have greatly improved the perceptual quality
105 of generated images, as well as introduced a num-
106 ber of new features. These features include mul-
107 tiple style transfer (Dumoulin et al., 2017), color-
108 preserving style transfer (Gatys et al., 2016), as
109 well as content-aware style transfer (Gatys et al.,
110 2017).

111 To allow for the styles of multiple images to be
112 transferred to a single content image, simple mod-
113 ifications to the style loss function have been pro-
114 posed in this paper (Makow and Hernandez, 2017).
115 The paper explains two perceptual loss functions
116 of interest: feature reconstruction loss and the style
117 reconstruction loss. Both make use of a pre-trained
118 loss network trained for image classification on
119 the ImageNet dataset. For Feature Reconstruction
120 Loss, the pixels of the output image have similar
121 feature representations as computed by the loss net-
122 work, and in Style Reconstruction Loss, the style
123 reconstruction loss penalizes differences in style,
124 such as differences in colors and textures (Gatys
125 et al., 2015b), common patterns, etc. This loss
126 formulation selects the style layers and weights
127 independently for each style image enabling the
128 model to weigh overall effects of each style image
129 on the generated image.

130 The paper (Wang et al., 2020) focuses on lever-

131 aging CST by combining existing approaches in
132 Serial Style Transfer (SST) architecture to improve
133 quality and Parallel Style Transfer (PST) architec-
134 ture to improve quality and flexibility. It is one
135 of the first papers to focus on domain independent
136 style transfer across three domains namely artis-
137 tic, semantic, and photo realistic using one content
138 and one style image. For SST architecture, it pro-
139 poses a linear combination of n methods chosen
140 thoughtfully based on each of their strengths with
141 regards to factors like content fidelity, global color,
142 and local patterns with performance levels of High,
143 Medium, and Low. The paper suggests reasonably
144 concluded serial schemes for each of the domains
145 based on their end goal. In PST architecture, a
146 single backbone is shared by all methods and the
147 total loss is optimized in parallel. The proposed
148 ParallelNet architecture uses weighted losses and
149 by increasing the weights of respective losses both
150 global color and local textures of the style image
151 of two different methods can be achieved. A user
152 study was conducted by randomizing various com-
153 binations of style and content images where users
154 were shown results of 10 compared methods and
155 SST/PST outperformed the state of the art models
156 for all three domains. Besides, the paper also men-
157 tions more effective and efficient scheme designs.

158 In this paper (Wang et al., 2021b), three new
159 quantitative factors (Content Fidelity, Global Ef-
160 fects, and Local Patterns) for evaluating NST are
161 introduced. Using these factors ideas have been
162 touched upon which involve the usage of CST and
163 user feedback. It brings up the point where CST
164 as an architecture is biased more to style transfer
165 than content preservation. This can be an important
166 feature in the case of multi style transfer, which we
167 aim to do. Additionally, the paper also discussed
168 how the user feedback can be incorporated using
169 the three quantitative evaluation factors by modify-
170 ing the loss function to include the user preferences.
171 This paper also provides a thorough comparison of
172 traditional multi-objective networks and CST archi-
173 tectures with average stylization ranking metrics
174 clearly highlighting the superiority of CST over
175 MO in this aspect.

176 "Interactive Artistic Multi-style Transfer" (Wang
177 et al., 2021a) is a research paper that proposes a
178 novel approach for interactive multi-style transfer
179 in artistic images. A new system is introduced that
180 allows users to manipulate the style and content of
181 an input image by selecting different style images

and adjusting various parameters. Their method combines feature normalization, multi-layer transformation, and iterative optimization to achieve high-quality style transfer while maintaining the content of the original image. The proposed approach enables users to create customized artistic images with multiple styles, which can be used for a variety of applications such as graphic design, digital art, and visual communication.

3 Implementation Details

For evaluating the hypothesis, three different model architectures were combined in SST. Each generated result was then compared to a baseline MO model to check if using CST improved our results as hypothesized. For evaluating the performance, three quantitative evaluation metrics were implemented from scratch. User-based qualitative evaluation was also designed to support it.

3.1 Dataset

Since NST does not fall under supervised learning or unsupervised learning, there is no fixed dataset that we are using. For testing and development purposes we have been using the same set of images as provided in the GitHub repository ([Gordić, 2020](#)).

3.2 Baseline Model

As a starting point, we decided to modify the traditional Gram Matrix network as a MO which will be able to handle our multi style image setting. This was giving satisfactory results visually and was used as a baseline model, shown in Fig 1.

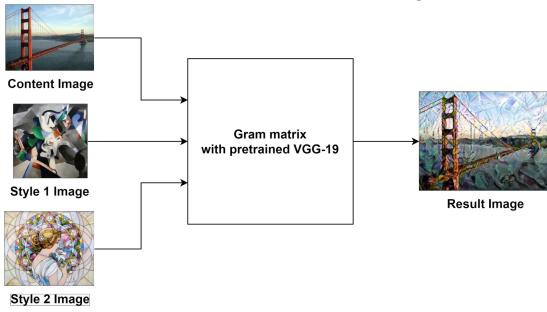


Figure 1: Baseline for Multiple Style Transfer using MO model

As a part of the baseline model ([Gatys et al., 2015a](#)), style transfer is performed, using gradient descent on random noise to minimize the deviation of target image from content. The deviation from the style representation of the style image, is defined as a set of Gram correlation matrices G_l with the entries $G_{ij}^l = (F_i^l, F_j^l)$, where l is the network layer, i and j are two filters of the layer l and F_k^l is the

array (indexed by spatial coordinates) of neuron responses of filter k in layer l. We have therefore used the VGG-19 network as backbone. It performs better at style transfer, as its architecture captures more information compared to other networks like AlexNet, GoogLeNet which strongly compress the input at first convolutional layer, thus losing a lot of fine details. L-BFGS is preferred over Adam due to its success as a second-order optimization algorithm that can handle larger batch sizes.

3.3 Models Explored

SST is a technique under CST that involves applying a sequence of style transfer algorithms or models to an image, with each subsequent model building on the output of the previous one. SST was implemented with two different architectures connected serially one after the other.

Many models were considered for exploration based on the paper ([Wang et al., 2020](#)). Three different models were used as architecture 1 while the architecture 2 was kept constant to be the Gram Matrix model as shown in Fig 2. A summary of each model used is given below. The content image and style 1 image are given to architecture 1. The output of architecture 1 becomes the content image for architecture 2 alongside the style 2 image and the output images are generated.

3.3.1 Model 1: Gram Matrix with pretrained VGG-19

The traditional Gram Matrix model with pretrained VGG-19 backbone (as seen in the baseline) was connected serially two times, as both architecture 1 and 2. The stylistic characteristics of style 1 were replaced by the stylistic characteristics of style 2 on the content image. Results obtained were not promising.

3.3.2 Model 2: Variational AutoEncoder (VAE)

A VAE model was used as architecture 1 along with the gram matrix model as architecture 2. ST-VAE (Style Transfer VAE), is a VAE for latent space-based style transfer which performs multiple style transfer by projecting nonlinear styles to a linear latent space, enabling to merge styles via linear interpolation before transferring the new style to the content image ([Liu et al., 2021](#)). According to the paper this method is faster and flexible for multiple style transfer compared to baseline methods. The ST-VAE method was implemented using exist-

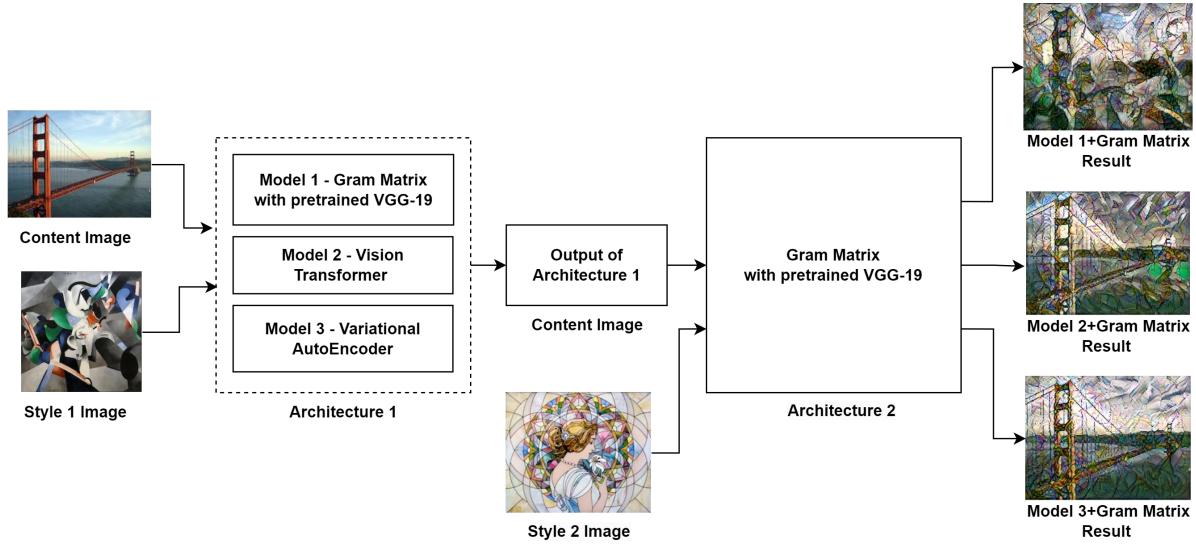


Figure 2: Three different models explored as Architecture 1 in SST

ing code for both single and multiple style transfer. For multiple style transfer the existing code base focuses on using four style images which was modified according to the project goals which aims at combining two style images. It showed promising results when implemented separately and hence was chosen to be further investigated in SST. The results obtained looked promising when compared to the gram matrix model connected serially twice.

3.3.3 Model 3: STyle TRansformer (STTR)

A Vision Transformer model was used as architecture 1 along with the gram matrix model as architecture 2. Existing approach adopts a global feature transformation to transfer style patterns into content images. Such a design usually destroys the spatial information of the input images and fails to transfer fine-grained style patterns into style transfer results. To solve this problem, we use the STTR network proposed in (Wang et al., 2022) which breaks both content and style images into visual tokens to achieve a fine-grained style transformation based on attention mechanism (Deng et al., 2021). When implemented for single style transfer, it showed good results. Hence it was used with gram matrix in serial to investigate its properties for combining multiple style images, in SST.

4 Evaluation

4.1 Quantitative Evaluation

We implemented three metrics (Wang et al., 2021b) from scratch to incorporate them into the existing code base for result evaluation. Creating an

efficient implementation for local patterns was particularly challenging as it involved generating and comparing all pairs of patches in two images.

4.1.1 Metric 1: Content Fidelity (CF)

It is a way to measure the extent to which the content of the original image is preserved in the stylized image. Since NST involves transforming the content of the original image to match the style of a style image, it is important to evaluate the quality of the output with respect to how much it retains the original content. In the equation below, \vec{c} represents the content image while \vec{x} represents the stylized result. N represents the total number of layers and $f_l(\cdot)$ represents the deep feature activations that are extracted as we go through each layer from $l=1$ to N . This results in the content fidelity factor.

$$CF(\vec{x}, \vec{c}) = \frac{1}{N} \cdot \sum_{l=1}^N \frac{f_l(\vec{x}) \cdot f_l(\vec{c})}{\|f_l(\vec{x})\| \cdot \|f_l(\vec{c})\|}$$

4.1.2 Metric 2: Global Effects (GE)

It is a measure of overall visual quality in NST that evaluates the similarity between the transferred style and the style image and it reflects the initial impression that the image leaves on human visual perception. It encompasses two aspects: Global Colors (GC) and Holistic Textures (HT), which were identified as important factors in user studies. GC compares the similarity of color histograms, while HT uses the Gram matrix to capture the HT across multiple layers. The GE factor is the average of both GC and HT. In the equation below,

331 \vec{s} represents the style image and \vec{x} represents
 332 the stylized result. c represents the channel and
 333 $hist_c(\cdot)$ represents the color histogram vector for
 334 the c^{th} channel. This results in the global colors
 335 factor.

$$336 \quad GC(\vec{x}, \vec{s}) = \frac{1}{3} \cdot \sum_{c=1}^3 \frac{hist_c(\vec{x}) \cdot hist_c(\vec{s})}{\|hist_c(\vec{x})\| \cdot \|hist_c(\vec{s})\|}$$

337 Then, in the next equation below, $G(\cdot)$ represents
 338 the gram matrix calculation and this equation re-
 339 sults in the HT.

$$340 \quad HT(\vec{x}, \vec{s}) = \frac{1}{N} \cdot \sum_{l=1}^N \frac{G(f_l(\vec{x})) \cdot G(f_l(\vec{s}))}{\|G(f_l(\vec{x}))\| \cdot \|G(f_l(\vec{s}))\|}$$

341 Then, in the final equation below, both the GC and
 342 HT factors are taken into account to compute GE.

$$343 \quad GE(\vec{x}, \vec{s}) = \frac{1}{2} \cdot (GC(\vec{x}, \vec{s}) + HT(\vec{x}, \vec{s}))$$

344 **4.1.3 Metric 3: Local patterns (LP)**

345 It is an evaluation metric to measure the extent
 346 to which patterns like brush strokes and exquisite
 347 motifs from the style image are transferred to the
 348 stylized image. This is done by using two com-
 349 ponents: the first assesses the similarity between
 350 the local patterns in the original and stylized im-
 351 ages, while the second compares the variety of the
 352 pattern categories identified in both images.

353 To calculate this we first extract 3x3 patches from
 354 $f_l(\vec{x})$ and $f_l(\vec{s})$ denoted by $\{\phi_i^l(\vec{x})\}_{i \in n_x}$ and
 355 $\{\phi_j^l(\vec{s})\}_{j \in n_s}$ where n_x and n_s are the number of
 356 patches respectively.

$$357 \quad CM(i) := argmax_{j=1 \dots n_s} \frac{\phi_i^l(\vec{x}) \cdot \phi_j^l(\vec{s})}{\|\phi_i^l(\vec{x})\| \cdot \|\phi_j^l(\vec{s})\|}$$

$$358 \quad LP_1(\vec{x}, \vec{s}) = \sum_{l=1}^N \sum_{i=1}^{n_x} \frac{\phi_i^l(\vec{x}) \cdot \phi_{CM(i)}^l(\vec{s})}{\|\phi_i^l(\vec{x})\| \cdot \|\phi_{CM(i)}^l(\vec{s})\|}$$

$$360 \quad LP_2(\vec{x}, \vec{s}) = \frac{1}{N} \sum_{l=1}^N \frac{t'_{cm}}{t'_s}$$

361 Here t'_{cm} and t'_s denote the number of unique
 362 patches in $\{\phi_{CM(i)}^l(\vec{s})\}$ and $\{\phi_j^l(\vec{s})\}$ respec-
 363 tively.

$$364 \quad LP(\vec{x}, \vec{s}) = \frac{1}{2} \cdot (LP_1(\vec{x}, \vec{s}) + LP_2(\vec{x}, \vec{s}))$$

4.2 Qualitative User Evaluation

To conduct a qualitative analysis, a Google form with 4 questions was created based on TA feedback. The form consisted of four questions, with one for each metric and one for the overall quality. We received over 30 responses for this exercise from technical individuals working in the field of Computer Science. The questions asked to the user were:

1. Which of the resulting images above is able to preserve the original content image? (CF) 374
375
2. Which result image has a similar color palette and textures compared to the two style images shown above? (GE) 376
377
378
3. Which resulting image has more lines/edges/patterns from the two style images shown above? (LP) 379
380
381
4. Which image do you feel is overall able to contain both the styles and visually pleasing? (Overall Quality) 382
383
384

We permitted users to choose multiple responses for the first question related to Content Fidelity, as our initial assumption was that the variables were independent in both architectures. For the other questions, we opted for a single option as our aim was to determine the relative ordering of output images for each metric.

4.3 Observations

The output image was generated by combining content image with two style images using different models. Fig 3 shows the results from different model combinations.

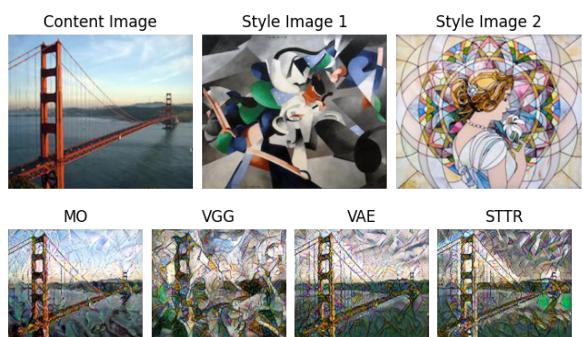


Figure 3: Content image: ‘Golden Gate’, style 1 image: ‘Udnie’, style 2 image: ‘Mosaic’

Based on results in Fig 3 different observations were made after running evaluation metrics, listed in table 1 and 2.

Q Survey Questions		MO(%)	VGG(%)	VAE(%)	STTR(%)
1. Which of the resulting images above is able to preserve the original Content image?	100.00	0.0	100.0	40.00	
2. Which result image has a similar color palette and textures compared to the two style images?	26.67	46.67	6.67	20.00	
3. Which resulting image has more lines/edges/patterns from the two style images?	0.00	46.67	13.33	40.00	
4. Which image do you feel is overall able to contain both the styles and visually pleasing?	33.33	0.00	46.67	20.00	

Table 1: Qualitative observations of MO vs CST

A parallel user study was performed against each quantitative metric (CF, GE, LP), to evaluate quality of output image. The above table denotes observations of the percentage distribution of users for each question asked in survey.

Metrics	MO	VGG	VAE	STTR
CF	0.38	0.26	0.82	0.51
GE	style 1	0.73	0.78	0.76
	style 2	0.81	0.84	0.81
LP	style 1	0.55	0.51	0.49
	style 2	0.53	0.59	0.43

Table 2: Quantitative metric observations MO vs CST. The table shows, different values calculated for the three quantitative metrics (CF, GE, LP). GE and LP is calculated for each style image.

Table 1 shows observation about the human study. For each question asked, larger percentages were picked towards Cascaded Networks - VGG, VAE, STTR. For question 1, since users were enabled multiple selections, they picked MO and CST outputs equivalently. For the second question, higher or equivalent percentages were picked by the user for CST. Intriguingly in the third question for LP, nearly no one picked MO and favored CST majorly. Finally the overall winner observed was CST, following our hypothesis over MO.

Further **Table 2** depicts inference about the three metrics that CF increased, GE slightly moved in both the directions, such that it decreased for style 1 and increased for style 2 and LP decreased for both the styles.

5 Results and Discussion

Table 1 and Table 2 depict the results, and based on them, we discuss our initial hypothesis and the findings in the Table shown in **Fig 4**.

In Fig 4 the Blue area denotes MO and rest is a part of different Cascaded architectures. The pie charts denote the user evaluation (outer circle) and quantitative metrics (inner circle) projected against each other. We can clearly see non-blue areas are

higher on the outer circle as well as cogent to the inner circle. One of the significant results was that none of the users picked MO output in question 3 (Table 2) that is for local patterns. For GE there was even distribution, when comparing MO and CST. Overall the users picked CST outputs more than MO.

Originally MST architecture was implemented by jointly minimizing the feature loss and the style loss formulated as the difference of Gram matrices. As a part of the alternate architectures, VAE transfers textures and styles efficiently and also successfully preserves the content information during style transition (Liu et al., 2021), which we could see in our results as well. Transformer network (STTR) preserves spatial information of the content structure and style patterns as it matches the style tokens onto content tokens in a fine-grained manner, which is coherent with our results.

6 Conclusion

Based on our results, we summarized a few conclusions. Content was equivalently preserved in CST as in MO. Style qualities like color palette and texture were significantly preserved in CST. Style patterns like lines, edges and patterns were more preserved in CST, true to our hypothesis. Based on the different architecture combination results we can conclude that CST results seem promising and can be considered to be a viable alternative to MO for better quality and style preservation in the generated image. Additionally we also show that quantitative metrics can be a feasible and reproducible option that can be used for evaluating NST performance in line with the prevalent qualitative evaluation techniques.

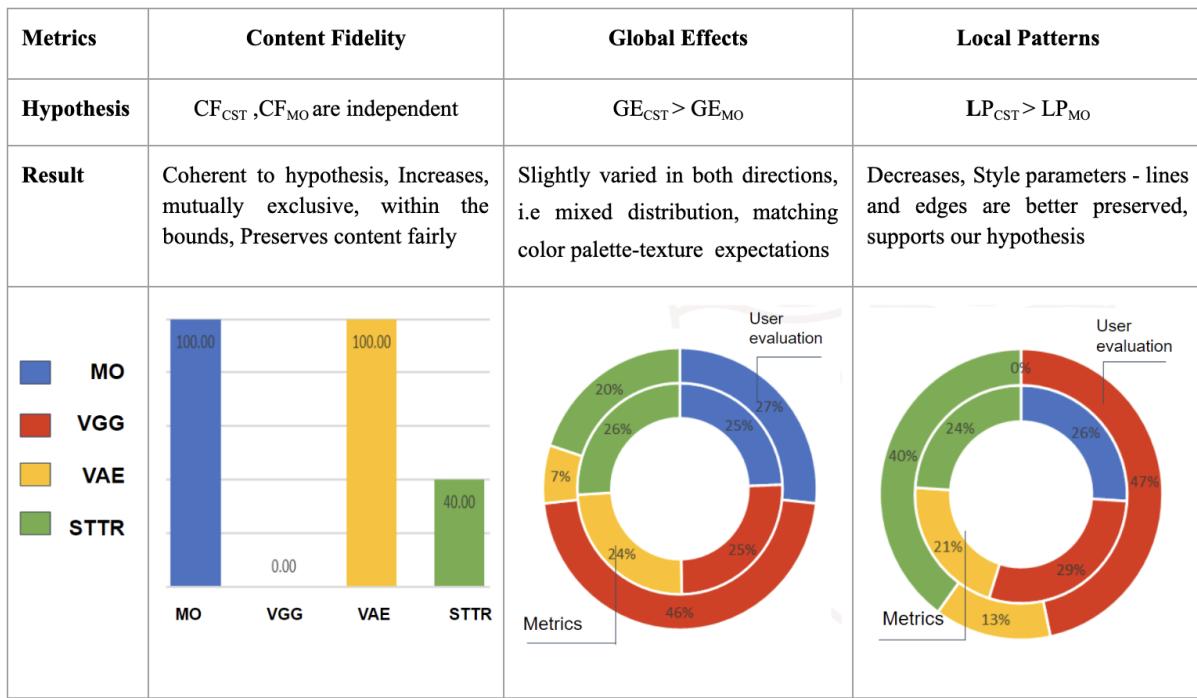


Figure 4: Evaluation Hypothesis and Results

Table shown above denotes the comparison between our initial hypothesis and our deductions against the same, based on our observations for each metric, both quantitatively and qualitatively. It also depicts the parallel conclusion between both kinds of evaluations (quantitative and qualitative)

463 7 Future Scope

464 The future scope of our work involves further im-
 465 provements to our proposed CST method. One
 466 extension consists of examining alternate architec-
 467 ture approaches, such as PST to extend our CST
 468 technique. Additionally, the original loss function
 469 can be investigated to extend its capability to trans-
 470 fer multiple style images optimally. Furthermore,
 471 enhancement based on user feedback can be ex-
 472 plored by streamlining it to existing neural style
 473 models. User evaluations can be mapped to current
 474 metrics for adjusting different parameters of an im-
 475 age. Lastly, the incorporation of visual parameters
 476 such as brightness, contrast, hue, etc. can be con-
 477 sidered for real-time rendering to further improve
 478 the style transfer process.

479 Acknowledgements

480 We thank **Bingjie Tang** for providing valuable feed-
 481 back and support throughout our project. We thank
 482 Professor **Jesse Thomason** for providing the op-
 483 portunity to pursue our research. Lastly, we thank
 484 Google Cloud Platform which provided us with
 485 vital computational resources that helped us ac-
 486 complish our project goals.

487 References

- Yingying Deng, Fan Tang, Xingjia Pan, Weiming Dong, Chongyang Ma, and Changsheng Xu. 2021. Stytr2: Unbiased image style transfer with transformers. *ArXiv*, abs/2105.14576.
- Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. 2017. A learned representation for artistic style. In *International Conference on Learning Representations*.
- Leon A. Gatys, Matthias Bethge, Aaron Hertzmann, and Eli Shechtman. 2016. Preserving color in neural artistic style transfer. *CorR*, abs/1606.05897.
- Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. 2015a. A neural algorithm of artistic style. *Journal of Vision*, abs/1508.06576.
- Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. 2015b. Texture synthesis using convolutional neural networks. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*, NIPS’15, page 262–270, Cambridge, MA, USA. MIT Press.
- Leon A. Gatys, Alexander S. Ecker, Matthias Bethge, Aaron Hertzmann, and Eli Shechtman. 2017. Controlling perceptual factors in neural style transfer. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3730–3738.
- Aleksa Gordić. 2020. Reference GitHub repository. <https://github.com/gordicaleksa/pytorch-neural-style-transfer>.

516
517
518
519
Zhi-Song Liu, Vicky Kalogeiton, and Marie-Paule Cani.
2021. [Multiple style transfer via variational autoencoder](#). In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 2413–2417.

520
521
Noah Makow and Pablo Hernandez. 2017. [Exploring style transfer: Extensions to neural style transfer](#). *Stanford University Report*.

523
524
525
526
Jianbo Wang, Huan Yang, Jianlong Fu, T. Yamasaki,
and Baining Guo. 2022. Fine-grained image style transfer
with visual transformers. In *Asian Conference on Computer Vision*.

527
528
529
530
Xiaohui Wang, Yiran Lyu, Junfeng Huang, Ziying
Wang, and Jingyan Qin. 2021a. [Interactive artistic multi-style transfer](#). *International Journal of Computational Intelligence Systems*, 14.

531
532
533
534
Zhizhong Wang, Lei Zhao, Haibo Chen, Zhiwen Zuo,
Ailin Li, Wei Xing, and Dongming Lu. 2021b. [Evaluate and improve the quality of neural style transfer](#). *Computer Vision and Image Understanding*, 207:103203.

535
536
537
538
Zhizhong Wang, Lei Zhao, Qihang Mo, Sihuan Lin, Zhi-
wen Zuo, Wei Xing, and Dongming Lu. 2020. [Cascade style transfer](#). *Under review as a conference paper at ICLR 2020*.

539 A Division of Labor

540
541
542
543
544
545
546
Each team member played a pivotal role in ensuring
the successful completion of the project. The
division of work was based on individual strengths
and interests, with regular discussions to ensure
equal distribution of responsibilities. The team col-
laborated smoothly using GitHub as the version
control technology.

547
548
The **project proposal** was meticulously planned,
with:

- 549
550
551
552
553
554
555
556
 - **Sneha** focusing on the introduction and outlining the objective and goals,
 - **Swarnita** identifying existing work and key ideas for novelty,
 - **Neha** highlighting the potential impact of the project; and
 - **Praveen** identifying potential challenges and solutions.

557
558
The team collectively estimated the computation costs, divided the work, and established a timeline.

559
560
561
During the **project pitch**, all four members collab-
orated to prepare the slides and content, with Sneha
delivering the pitch.

562
563
564
In the **survey report**, each team member con-
tributed to the Related Work section by reading
and summarizing one research paper.

- 565
566
567
568
569
570
571
 - **Praveen** worked on quantitative evaluation related work,
 - **Swarnita** worked on cascade style transfer architecture research,
 - **Sneha** worked on multiple style techniques; and
 - **Neha** worked on exploration of interactive styling.

572
573
574
575
576
577
Each team member also contributed two challenges to the ‘Challenges’ section that were identified from the respective papers and earlier from the proposal. All team members worked on contributing a set of references to the survey report and the final addition to Overleaf Latex format.

578 Before the midterm report:

- 579
580
581
582
583
584
585
586
 - Implementation of the three evaluation metrics was a collaborative effort of **Swarnita, Neha and Praveen**,
 - **Praveen** worked on the Gram matrix-based method,
 - **Sneha** implemented the Vision Transformer-based method; and
 - **Swarnita** handled the VAE-based method.

587 Before the final report:

- 588
589
590
591
592
593
594
595
596
597
 - **Praveen** handled the final efficient implementation and verification of the three evaluation metrics,
 - **Neha** designed the qualitative study,
 - **Swarnita** integrated the VAE-based method in a SST setting; and
 - **Sneha** integrated the Transformer-based architecture in a SST setting and interpreted observations from the different architecture combinations and qualitative study.

598
599
600
601
All team members collaborated in gathering responses for the study. Finally, for the final presentation and report, all team members worked together and contributed equally.

602 B Appendix

603
604
605
606
607
608
Fig 5, shows the transition of output image alongside style addition, simultaneously from both the input styles. We can see different levels or percentage effect of styles, for MO. Similarly Fig 6, shows the transition of output image alongside style 1 addition, and Fig 7 shows style 2 addition transitions.



Figure 5: Multi-objective Network Transitions at various iterations for only style 1 addition



Figure 6: Gram Matrix + Gram Matrix SST transitions at various iterations for style 1

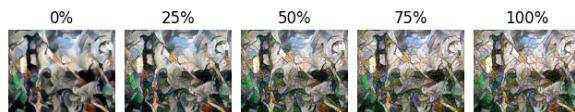


Figure 7: Gram Matrix + Gram Matrix SST transitions at various iterations for style 2

Based on the inputs in Fig 8, left side of Fig 9 depicts visually satisfactory results for MO and the right side shows style 2 dominating over style 1. This results from using the same architectures for SST.

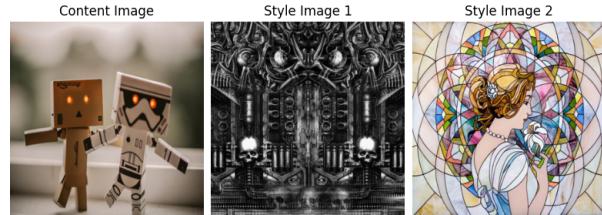


Figure 8: Content image: 'Figures', style 1 image: 'Giger Crop', style 2 image: 'Mosaic'

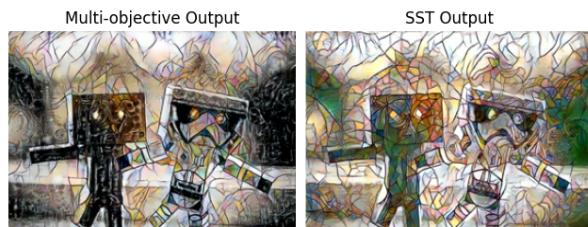


Figure 9: MO Output and Gram Matrix + Gram Matrix SST Output

For further exploration with different models following were the results upon which we built the cascaded networks i.e. combinations of different architectures. Implementation of just the VAE model separately for MST is shown in Fig 10. The obtained results on combining Transformer Network serially two times for MST are shown in Fig 11. For additional exploration we also changed the order of the style incorporated in the final image to study its effect, as shown in Fig 12.



Figure 10: VAE Output with Style Image 'Ben Giles' then 'Starry Night'



Figure 11: STTR Output with Style Image 'Udnie' then 'Mosaic'



Figure 12: STTR Output with Style Image 'Mosaic' then 'Udnie'

Fig 13, 14, 15, 16 and 17 are the screenshots showing the Google form created for the user study.

We perform the following combination of the images :

(Content Image: "Bridge") + (Style 1 : "Udnie") + (Style 2 : "Mosaic")

Following results were generated :

(A) (B) (C) (D)

Figure 13: Qualitative User Study Form Image Options

Q.1) Which of the resulting images above is able to preserve the original Content * image (Bridge) ?
(You can select multiple options if they are equal for you)

- A
- B
- C
- D

Figure 14: Qualitative User Study Form Question 1

Q.2) Which result image has a similar color palette and textures compared to the * two style images (udnie & mosaic paintings) shown above ?

- A
- B
- C
- D

630

631

632

633

634

635

Figure 15: Qualitative User Study Form Question 2

Q3) Which resulting image has more lines/edges/patterns from the two style images (udnie & mosaic paintings) show above? *

- A
- B
- C
- D

637

Figure 16: Qualitative User Study Form Question 3

Q.4) Which image do you feel is overall able to contain both the styles and visually * pleasing?

- A
- B
- C
- D

638

Figure 17: Qualitative User Study Form Question 4