



# **CLIP FOR IMAGE STYLE TRANSFER: EXPLORING TEXT-IMAGE CORRELATIONS**

Lappeenranta-Lahti University of Technology LUT

Master's Program in Computational Engineering, Master's Thesis

2024

Nadine Bisanukuli Cyizere

Examiners: Professor Zhisong Liu

(Use \SecondExaminer to replace this text.)

# ABSTRACT

Lappeenranta-Lahti University of Technology LUT  
School of Engineering Science  
Computational Engineering

Nadine Bisanukuli Cyizere

## **CLIP for Image Style Transfer: Exploring Text-Image Correlations**

Master's thesis

2024

x pages, y figures, z tables, w appendices

Examiners: Professor Zhisong Liu and (Use \SecondExaminer to replace this text.)

Keywords: computer vision, machine vision, image processing, pattern recognition

# ACKNOWLEDGEMENTS

The author can decide the contents of this page. Usually the place of work and related people (supervisors, collaborators, friends, relatives, etc.) are acknowledged.

I would like to thank my supervisors ...friends ... family ...

Lappeenranta, January 16, 2024

*Nadine Bisanukuli Cyizere*

Advice: All symbols and abbreviations are listed on this page in the alphabetical order. Remember to introduce the abbreviation when it is used in the text for the first time.

You may use the automated system, depending on your LaTeX environment:

```
\glsnogroupskiptrue
\setlength{\glsdescwidth}{1.0\hsize}
\printglossary[title=LIST OF ABBREVIATIONS,type=\acronymtype,
style=long, nonumberlist, nopostdot]
```

# CONTENTS

<b>1</b>	<b>INTRODUCTION</b>	<b>5</b>
1.1	Background . . . . .	5
1.2	Objectives and delimitations . . . . .	6
1.3	Structure of the thesis . . . . .	7
<b>2</b>	<b>RELATED WORK</b>	<b>8</b>
<b>3</b>	<b>PROPOSED METHODS</b>	<b>11</b>
3.1	Results . . . . .	11
<b>4</b>	<b>DISCUSSION</b>	<b>12</b>
4.1	Future work . . . . .	12
<b>5</b>	<b>CONCLUSION</b>	<b>13</b>
	<b>REFERENCES</b>	<b>14</b>

# 1 INTRODUCTION

## 1.1 Background

In the field of Computer Vision and Natural Language Processing, the combination of textual description and visual elements has given rise to innovation in many real-world applications with image style transfer as the crucial and captivating one. Image style transfer is an accommodating tool that facilitates the development of artistic works and enhances the quality of visual aesthetics. This technique is useful in many fields, such as graphic design, virtual reality, photo editing, film production, and even social media optimization. The field of image style transfer has had different advancements over the years with the involvement of different machine learning techniques mostly the use of deep learning. The main goal of image style transfer is to apply the style of an image usually referred to as **style reference**, to other images while preserving the original content. This thesis will dive deep into the intersection of Contrastive Language-Image Pre-training (CLIP) [1,2] of the original target image and image style transfer to leverage the text-image correlations for more advancements in the domain. OpenAI's CLIP model has demonstrated new and better possibilities in the field of image style transfer. Unlike some traditional methods that require reference style images, CLIP can do style transfer by only using a text description of the desired style. This is very useful in case a user doesn't have reference style images but is interested in transferring styles based on their imagination.

Several methods have been proposed to better the performance of image style transfer. Despite the advancements, challenges persist in the field of image style transfer. One of the challenges is the need for methods that can adapt to styles without loss of content. Traditional methods depend on reference style image which limits their applicability in scenarios where users might not have specific reference images but wish to transfer styles. Additionally achieving balance between style transfer and preservation of content remains an issue. This lies in more development of algorithms that can incorporate the desired style while preserving content of the image. Another challenge lies in the domain specificity since many existing methods lie to specific domain such as portraits, landscapes and industrial settings. Hence generalising these models to handle arbitrary styles is an ongoing challenge but have tried to solve. Lastly computational efficiency is a concern since as image style transfer models become more and more complex, there is a growing need of development of less complex models that can deliver results without much consumption of computational resources. There have been multiple solutions to

the challenges such development of PCA-based knowledge Distillation method to distill lightweight models demonstrating its usability with different architectures hence improvement of efficiency and suitability for real time applications. Also, advancements in the Industrial Style Transfer method have shown promising results in creating new visual products with a nice appearance for industrial designers' reference.

The exploration of CLIP for image style transfer presents an exciting edge in the field of computer vision. The ability to manipulate and transfer styles using text descriptions could revolutionize how we interact with digital imagery. This thesis aims to contribute to this field by addressing the challenges of style transfer without explicit style references, achieving a balance between style and content.

## 1.2 Objectives and delimitations

The primary objective of this thesis is to explore the potential of CLIP for image style transfer with a particular focus on leveraging text-image correlations. The specific research questions that this work aims to address are:

1. **Can CLIP be effectively used for image style transfer?** This involves the development and evaluation of a method that uses CLIP to transfer the style of a text description to an image. The performance of this method will be evaluated based on the quality of the style transfer and the preservation of the original content.
2. **How can text-image correlations be leveraged to improve image style transfer?** This involves investigating how the correlations between text descriptions and their corresponding images, as learned by CLIP, can be used to enhance the style transfer process. This could involve modifying the style transfer method based on the specific text-image correlations identified.
3. **Analysis and experiments on artistic text-driven style transfer:** This objective entails a comprehensive examination of text-driven style transfer with a specific emphasis on artistic styles. The goal is to conduct in-depth analyses and experiments to understand the nuances and effectiveness of applying text-driven style transfer in the realm of artistic expressions. This involves exploring various artistic styles, such as those of specific artists or art movements, and evaluating the results to provide insights into the potential and challenges of this approach.

The scope of this thesis is limited to the following delimitations:

- The study will focus on the use of OpenAI’s CLIP model for image style transfer. Other models or methods for image style transfer or text-image correlation learning will not be considered.
- The performance of the proposed method will be evaluated using available datasets. The collection of new data is beyond the scope of this work.
- While the aim is to develop a method capable of high-quality style transfer, the computational efficiency of the method will not be a primary focus of this work.
- The thesis will not explore the use of CLIP for other tasks beyond image style transfer, such as image generation or text-to-image synthesis.

By addressing these research questions, this thesis aims to contribute to the ongoing efforts to leverage text-image correlations for image style transfer. However, it is important to note that the proposed methods are subject to the inherent limitations and uncertainties of machine learning techniques. Future work may be needed to refine the methods and address any limitations identified in this study.

### **1.3 Structure of the thesis**

## 2 RELATED WORK

In the realm of Computer vision, style transfer has been an important yet interesting field. style transfer has witnessed significant advancements with an exploration of different techniques [3,4]. Style transfer involves producing a content image in the style of another image, which allows for adaption to arbitrary new styles via a feed-forward neural network. It entails matching the mean and variance of the content features to those of the style features to perform real-time style transfer without being limited to a predefined range of styles. Style transfer has been under study for a long time with it originating from unrealistic photo output [5] also it is closely related to texture synthesis and transfer [6–8]. There have been multiple approaches that were put into practice to help solve the issue of unrealistic photo output, such approaches include Non-parametric sampling [7,9] and histogram matching on linear filter responses [10]. These approaches are low-level statical methods and tend to fail to capture the underlying structures. This was followed by Gatys et al. [11] who demonstrated interesting style transfer methods using convolutional layers by matching features of a DNN. As time went by, more improvements were made to [11]. Li and Wand [12] came up with an approach based on the Markov field(MRF) in the deep space to focus on local patterns. Gatys et al. [13] later proposed ways to preserve color, spatial location, and the style of style transfer. This was followed by Ruder et al. [14] where they improved the quality of video style transfer by imposing constraints.

Going deeply on Gatys et al. [11], it is based on a slow optimization process that updates the image while minimizing loss of content. It takes minutes to converge with modern GPUs. A common solution is to replace the optimization process with a feed-forward neural network trained to minimize the objective function [13,15,16]. These approaches are three orders of magnitude faster than the optimization alternative since they use a feed-forward style transfer approach. More was done on the feed-forwards style transfer by Wang et al [17]where they enhanced the method by making a multi-resolution architecture. Later on, Ulyanov et al [18] proposed an improved way to improve the quality and diversity of generated samples. Nonetheless, the feed-forward methods were limited since each network was limited to a fixed style. The problem was addressed by Dumoulin et al. [19] who introduced a single network able to encode 32 styles and interpolations. More was done by Chen and Schmidt who introduced a feed-forward method that can transfer arbitrary styles by the use of a style swap layer. Many more problems were addressed such as the loss function to be used, and more effective loss functions were introduced such as MRF [12], Adversarial loss [15], histogram loss [20], CORAL loss [21], MMD loss [22]. Style transfer methods can be divided into three main types:



text-driven style transfer, image-driven text transfer, and attention-based style transfer.

**Image-Driven Style Transfer.** It is a task in image processing that involves displaying a picture's semantic content in several styles. Convolutional Neural Networks facilitate the creation and alteration of high-quality images by separating and combining the image's content and style. Transferring styles from one image to another is an issue of texture transfer. In texture transfer, the goal is to produce texture from a source image while of course preserving the content of the expected resulting image. There exists a wide range of powerful non-parametric algorithms that can produce realistic photos by resampling pixels of a source texture [7, 8, 23, 24]. There have been multiple approaches for solving image transformation tasks, such approaches include training a feed-forward convolutional neural network which is done in a supervised manner by use of a per-pixel loss function to keep track of the differences between output and the true/expected images. The approach was used by Dong et al. [25] to increase the resolution of many types of images and create better artistic-looking images from photographs while having several orders of magnitude faster than many state of art techniques. Additionally, the approach was used by Cheng et al [26, 27] for colorisation also by Long et al. for segmentation and last but not least Eigen et al for d

Furthermore Johnson et al . introduced the use of the VGG network for style transfer emphasizing the effectiveness of feature extraction from pre-trained networks in capturing and applying artistic styles across images [28]. More research was done by Huang et al. who showed real-time abilities of arbitrary style transfer through adaptive instance normalization in the AdaIN framework [29]. The two papers showed the foundation of Image-driven style transfer across multiple domains. Later on, Image-driven style transfer extended to more complex scenarios such as indoor 3D scene constructions. Hollein et al. introduced styleMesh, emphasizing the ability to complex three-dimensional environments [30]. Additionally, the issue of unpaired cartoon images was addressed by introducing gated cycle mapping which showed the potential for image-driven style transfer without explicitly paired training data [31]. This also gave rise to another style transfer that's known as Attention-based style transfer whose architectures emerged as a significant theme in image-driven style transfer with the introduction to StyTr2 that incorporated transformers to capture complex style patterns [32]. Additionally, efforts were made to handle the computational efficiency and quality of the styles issue as seen in [33]. Moreover, there was work done for industrial style transfer addressing industrial contexts for large-scale geometric warping and content preservation for preserving content while applying complex styles [34]. More efforts were made to arbitrary style transfer and domain

generalization with exact feature distribution matching provided insights into handling diverse styles and applicability of models across different domains [35].

**Text-Driven Style Transfer.** It emphasizes using textual descriptions to guide the transformation of images. It has been a significant area of research in style transfer over the past years marked by different discoveries. Such discoveries include the discovery of a novel framework. This work introduces a novel neural style transfer framework addressing practical scenarios where users lack reference style images but desire style transfer through text descriptions. Leveraging CLIP’s pre-trained text-image embedding model, the proposed method achieves style transfer without a style image, relying solely on a single text condition [3]. This was followed by another crucial discovery that introduced text-driven style transfer (TxST) as a flexible alternative to traditional image style transfer, leveraging advanced image-text encoders like CLIP [36]. The proposed method employs a contrastive training strategy and a novel cross-attention module, enabling arbitrary artist-aware style transfer with superior performance, demonstrating potential for future advancements in mimicking various artistic styles.

**Attention-Based Style Transfer.** It uses the power of attention mechanisms through the use of advanced architectures like transformers. This has become a focus recently due to its crucial application in image-driven style transfer since it employs sophisticated architectures, specifically transformers, to enhance the quality and diversify the style transfer results. A very early contribution that dives deep into the integration of transformers using attention mechanism to capture complex style patterns and elevate expression abilities of style transfer models [32]. The adaptability of attention-based architectures showcased in StyTr2 signifies a move from traditional methods while offering more possibilities to create visually exciting images. Although the paper [35] doesn’t explicitly mention attention-based mechanisms, the use of transformers suggests the use of attention-based features. The paper hints at the importance of aligning feature distributions for arbitrary style transfer and domain generalization that showcase the ongoing development of methodologies in attention-based styles. Furthermore, the paper [37] talks about the importance of the attention mechanism in style transfer. This work introduces an attention-based approach allowing the model to create realistic styled images with multiple strokes, this ensures focus on specific regions enhancing the stylization and adding good characteristics to the output image. These papers show how attention-based mechanisms are transforming artistic style transfer.

## **3 PROPOSED METHODS**

### **3.1 Results**

## **4 DISCUSSION**

### **4.1 Future work**

## **5 CONCLUSION**

## REFERENCES

- [1] Suchen Wang, Yueqi Duan, Henghui Ding, Yap-Peng Tan, Kim-Hui Yap, and Jun-song Yuan. Learning transferable human-object interaction detector with natural language supervision. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 929–938, 2022.
- [2] Jimmy Lei Ba, Kevin Swersky, Sanja Fidler, and Ruslan Salakhutdinov. Predicting deep zero-shot convolutional neural networks using textual descriptions. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 4247–4255, 2015.
- [3] Gihyun Kwon and Jong Chul Ye. Clipstyler: Image style transfer with a single text condition. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18041–18050, 2022.
- [4] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- [5] Jan Eric Kyprianidis, John Collomosse, Tinghuai Wang, and Tobias Isenberg. State of the "art": A taxonomy of artistic stylization techniques for images and video. *IEEE Transactions on Visualization and Computer Graphics*, 19(5):866–885, 2013.
- [6] Michael Elad and Peyman Milanfar. Style transfer via texture synthesis. *IEEE Transactions on Image Processing*, 26(5):2338–2351, 2017.
- [7] Alexei A. Efros and William T. Freeman. Image quilting for texture synthesis and transfer. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '01*, page 341–346, New York, NY, USA, 2001. Association for Computing Machinery.
- [8] A.A. Efros and T.K. Leung. Texture synthesis by non-parametric sampling. In *Proceedings of the Seventh IEEE International Conference on Computer Vision*, volume 2, pages 1033–1038 vol.2, 1999.
- [9] Oriel Frigo, Neus Sabater, Julie Delon, and Pierre Hellier. Split and match: Example-based adaptive patch sampling for unsupervised style transfer. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 553–561, 2016.

- [10] D.J. Heeger and J.R. Bergen. Pyramid-based texture analysis/synthesis. In *Proceedings., International Conference on Image Processing*, volume 3, pages 648–651 vol.3, 1995.
- [11] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2414–2423, 2016.
- [12] Chuan Li and Michael Wand. Combining markov random fields and convolutional neural networks for image synthesis. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2479–2486, 2016.
- [13] Leon A. Gatys, Alexander S. Ecker, Matthias Bethge, Aaron Hertzmann, and Eli Shechtman. Controlling perceptual factors in neural style transfer. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3730–3738, 2017.
- [14] Sun Fengxue, Sun Yanguo, Lan Zhenping, Wang Yanqi, Zhang Nianchao, Wang Yuru, and Li Ping. Image and video style transfer based on transformer. *IEEE Access*, 11:56400–56407, 2023.
- [15] Chuan Li and Michael Wand. Precomputed real-time texture synthesis with markovian generative adversarial networks. volume 9907, pages 702–716, 10 2016.
- [16] Dmitry Ulyanov, Vadim Lebedev, Andrea Vedaldi, and Victor Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. *ICML’16*, page 1349–1357. JMLR.org, 2016.
- [17] Xin Wang, Geoffrey Oxholm, Da Zhang, and Yuan-Fang Wang. Multimodal transfer: A hierarchical deep convolutional neural network for fast artistic style transfer. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7178–7186, 2017.
- [18] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4105–4113, 2017.
- [19] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. page 9, 04 2017.
- [20] Pierre Wilmot, Eric Risser, and Connelly Barnes. Stable and controllable neural texture synthesis and style transfer using histogram losses. 01 2017.

- [21] Xingchao Peng and Kate Saenko. Synthetic to real adaptation with generative correlation alignment networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1982–1991, 2018.
- [22] Yanghao Li, Naiyan Wang, Jiaying Liu, and Xiaodi Hou. Demystifying neural style transfer. pages 2230–2236, 08 2017.
- [23] Li-Yi Wei and Marc Levoy. Fast texture synthesis using tree-structured vector quantization. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '00, page 479–488, USA, 2000. ACM Press/Addison-Wesley Publishing Co.
- [24] Vivek Kwatra, Arno Schödl, Irfan Essa, Greg Turk, and Aaron Bobick. *Graphcut Textures: Image and Video Synthesis Using Graph Cuts*. Association for Computing Machinery, New York, NY, USA, 1 edition, 2023.
- [25] Xiancai Ji, Yao Lu, and Li Guo. Image super-resolution with deep convolutional neural network. In *2016 IEEE First International Conference on Data Science in Cyberspace (DSC)*, pages 626–630, 2016.
- [26] Zezhou Cheng, Qingxiong Yang, and Bin Sheng. Deep colorization. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 415–423, 2015.
- [27] Richard Zhang, Phillip Isola, and Alexei A. Efros. Colorful image colorization. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 649–666, Cham, 2016. Springer International Publishing.
- [28] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 694–711, Cham, 2016. Springer International Publishing.
- [29] ManMan Peng and Zhongrui Zhu. Enhanced style transfer in real-time with histogram-matched instance normalization. In *2019 IEEE 21st International Conference on High Performance Computing and Communications; IEEE 17th International Conference on Smart City; IEEE 5th International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*, pages 2001–2006, 2019.
- [30] Lukas Höllein, Justin Johnson, and Matthias Nießner. Stylemesh: Style transfer for indoor 3d scene reconstructions. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6188–6198, 2022.



- [31] Yifang Men, Yuan Yao, Miaomiao Cui, Zhouhui Lian, Xuansong Xie, and Xian-Sheng Hua. Unpaired cartoon image synthesis via gated cycle mapping. pages 3491–3500, 2022.
- [32] Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. Stytr2: Image style transfer with transformers. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11316–11326, 2022.
- [33] Tai-Yin Chiu and Danna Gurari. Pca-based knowledge distillation towards lightweight and content-style balanced photorealistic style transfer models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7834–7843, 2022.
- [34] Jinchao Yang, Fei Guo, Shuo Chen, Jun Li, and Jian Yang. Industrial style transfer with large-scale geometric warping and content preservation. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7824–7833, 2022.
- [35] Yabin Zhang, Minghan Li, Ruihuang Li, Kui Jia, and Lei Zhang. Exact feature distribution matching for arbitrary style transfer and domain generalization. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8025–8035, 2022.
- [36] Zhi-Song Liu, Li-Wen Wang, Wan-Chi Siu, and Vicky Kalogeiton. Name your style: text-guided artistic style transfer. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 3530–3534, 2023.
- [37] Yuan Yao, Jianqiang Ren, Xuansong Xie, Weidong Liu, Yong-Jin Liu, and Jun Wang. Attention-aware multi-stroke style transfer. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1467–1475, 2019.