

# **Artificially Painting Pictures & Artworks**



Universität Heidelberg

**Fakultät für Physik**

# **Artificially Painting Pictures & Artworks**

**Masters Thesis**

Arthur Willy Heimbrecht

June 2nd 2020

supervised by Professors Björn Ommer and Tilman Plehn

# **Abstract**

This is where the abstract is going to be.

# Contents

<b>Abstract</b>	<b>iv</b>
<b>Contents</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
<b>THEORETICAL BACKGROUND &amp; RELATED WORK</b>	<b>5</b>
<b>2 Machine Learning &amp; Computer Vision</b>	<b>7</b>
2.1 Artificial NNs . . . . .	7
2.2 Fully Connected NNs . . . . .	15
2.3 Convolutional NNs . . . . .	17
2.4 Recurrent Neural Networks . . . . .	21
2.5 Types of Learning . . . . .	21
<b>3 Optimization &amp; Gradient Descent</b>	<b>23</b>
3.1 Optimization Problem . . . . .	23
3.2 Gradient Descent . . . . .	25
3.3 Backpropagation . . . . .	29
3.4 Vanishing Gradients . . . . .	31
<b>4 Models</b>	<b>35</b>
4.1 Discriminators & Classifiers . . . . .	35
4.2 Generators & Decoders . . . . .	36
4.3 Auto-Encoders . . . . .	37
4.4 Generative Adversarial Networks . . . . .	39
<b>5 Artistic Computer Vision</b>	<b>43</b>
5.1 Brushstroke Extraction . . . . .	43
5.2 Style Transfer . . . . .	45
5.3 Painterly Rendering . . . . .	56
5.4 Conclusion . . . . .	64

<b>CONTRIBUTION AND EXPERIMENTS</b>	<b>65</b>
<b>6 Approach</b>	<b>67</b>
6.1 Motivation . . . . .	67
6.2 Neural Renderer . . . . .	69
6.3 Stroke Approximation . . . . .	80
<b>7 Evaluations &amp; Ablation Experiments</b>	<b>107</b>
7.1 Neural Renderer . . . . .	107
7.2 Optimization Procedure . . . . .	107
<b>CONCLUSION</b>	<b>109</b>
<b>8 Discussion</b>	<b>111</b>
<b>9 Outlook</b>	<b>113</b>
<b>APPENDIX</b>	<b>115</b>
<b>A Appendix</b>	<b>117</b>
<b>Bibliography</b>	<b>119</b>

# 1

## Introduction

The American Cambridge dictionary defines the word "picture" as

**picture:** a drawing, painting, photograph, etc.

(As submitted to the Cambridge American Dictionary)

this need at least some pictures in it

add phonetic signs

which right away separates between different media that can make up a picture. If one were to look up the definitions of "drawing", 'painting" and "photograph" as well, this is what distinguishes the three words from one another.

Paintings are only pictures that are made with paint.

A drawing confines oneself to pictures drawn with a pencil or pen.

And finally, a photograph is a picture taken with a camera.

All of these three methods could be used to picture the same image, but will always give different results, as each medium comes with its own limitations.

The painting will almost always deviate in its details from the original image, as the paint tends to mix and flow. The drawing will present large evenly colored areas often in a grainy way due to the structure of the paper which it has been drawn on. The photograph will limit an artist in showing anything more than the raw visual input.

Yet, all of these techniques also have their own characteristics that other techniques do not possess.

When comparing paintings to photographs this becomes especially clear. The layered texture of paintings which varies

vastly in thickness within a single painting is something that photographs will never be able to capture and display, as photographs are just the projection of a 3D world onto a 2D plane. This will not only be confirmed by art historians or the like, but can be seen when looking at the industry of painting replicas of masterpieces with oil such as the Starry Night by van Gogh or most famously the Mona Lisa by DaVinci. There can only be a market for such replicas can only if there is enough of a difference between an actual painting and a very high resolution photo of a painting.

This does not hold as well for drawings on the other hand. As drawings are also inherently a 2D representation with just an infinitesimal height to them, they CAN be captured by photos quite well. Also, as technology has become more and more advanced, it is possible to capture and print a drawing in such a high quality and fidelity to the original drawing, that it is hard to tell the original drawing and the print apart.

This shows also in the popularity of drawing applications on computers and tablets that have been useful tools to digital artists for over 2 decades now. And further development it has now become a feasible challenge to mimic the feel of pen on paper with styluses on glass tablets.

All of this makes it clear, that the art of drawing has been object to the digital revolution as much as the art of taking photographs.

But what about paintings? Are there not also applications that try to imitate painting techniques, or 3D printed replicas of famous paintings? Yes there are, but it is clear that this comes with a lot of limitations or huge effort, as 3D scans and prints and real-time fluid simulations are cutting-edge technology. And even then the majority of digital content is consumed through 2D screens which are not able to display what makes a painting unique.

So why even bother then, if all of these limitations are in place?

As it already has been hinted, there is existing technology that is capable of taking paintings into the digital realm. At the same time Augmented Reality (AR) and Virtual Reality (VR) gain more traction every year, which will only strengthen this development.

Thus, the question arises: Is it possible to digitalize paintings? One way to achieve this are the mentioned 3D scans and Gigapixel photographs of paintings. Another way are painting applications which simulate brush strokes virtually. But could these two approaches be combined such that a painting can be made up of digital brush strokes that replicate the original brush strokes in a painting?

Such a combined approach would help to promote the current development in this area and bring interesting applications in other fields. Two exemplary applications would be the conservation of images and the study of paintings.

As far as conservation goes, an image that is disassembled into its individual brush strokes will store information about the original painting in a more efficient manner than Gigapixel images or 3D scans. Both of these archiving methods generate huge amounts of data, which is why only a few hundred images in the world are archived in this way [[googleartproject](#)]. Having extracted brush stroke information at hand could store this information more efficiently and improve reconstruction as simulation techniques improve as well.

Considering the study of paintings, there are already computer assisted methods to distinguish forgeries from real artworks. A way of extracting brush strokes from paintings would open up new possibilities to study the style and techniques of an artist that goes beyond visual inspection.

At last there are other fields that deal with artworks and could benefit from such an approach such as artistic style transfer in computer vision. Artistic style transfer relies exclusively on images of paintings when it comes to learning the style of an artist. As it has been laid out images do not

capture everything that makes up a painting and it would be interesting if style transfer could benefit from more advanced information about a painting.

This thesis will perform an experimental evaluation whether it is possible to extract individual brush strokes from a photograph of a painting. Furthermore, it shall make a first attempt at performing style transfer using brush strokes.

As there is a vast realm of different painting techniques and materials, this work shall only be evaluated on oil painting brush strokes, as there is little data available on other techniques and it is quite possibly the easiest to replicate as well. Even though there are many artists which created oil paintings in the past, this work will mainly focus on later works by van Gogh. These works are well known and thus there are many high quality photographs of his works. Also, most of these paintings show very clear brush strokes which is decisive of van Gogh's style.

introduce the approach  
broadly

reformulate

**Structure of this work** First, the theoretical background for such an attempt along with related work will be presented. Then the approach itself will be inferred in two parts that describe the training of a differentiable renderer and the subsequent retrieval of brush strokes from an image. This is followed by experiments as well as an ablation study of the approach. Lastly, there will be a discussion about the results and the future direction of research.

**Allocation of individual contributions** The results, which are presented in this thesis are the result of the combined efforts by Dmitryo Kotovenko, Matthias Wright and Arthur Heimbrecht.

# **THEORETICAL BACKGROUND & RELATED WORK**



# Machine Learning & Computer Vision

# 2

In this chapter the theoretical backgrounds and motivation of neural networks ought to be explored.

First, a physical/biological motivation will set the foundation for the theoretical background accompanied by a more mathematical motivation later on. Thereafter, more advanced lines of thought will be introduced that finally lead to convolutional neural networks as the current driving force in computer vision.

## 2.1 Artificial Neural Networks

### Neurobiological Inspiration

Many human achievements have at least been partially inspired by studying nature. A very popular example is that of airplanes and birds. By studying how birds can fly, people found out that the shape of their wings is essential. Inventors and engineers have taken this as inspiration and slowly but steadily came up with planes and the like. However, planes are not birds, as they are not flapping their wings (yet?), but in combination with other inventions like jet engines, they obtain the same (or better) capability of flying like birds.

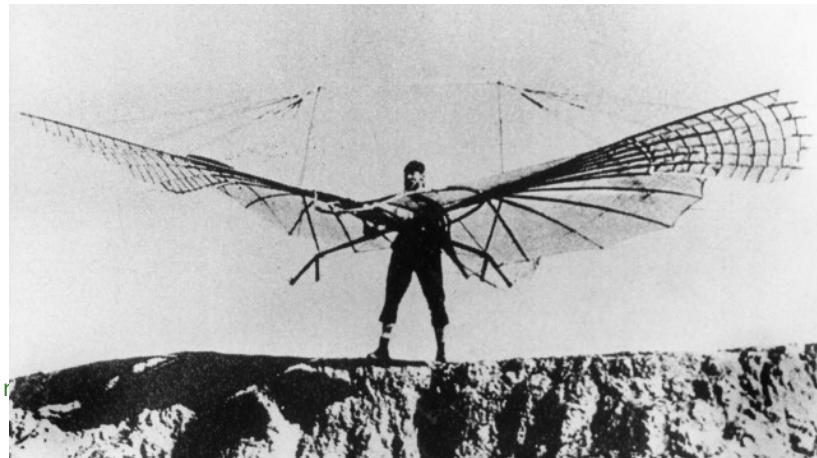
Similarly, the brain features a lot of insights, how intelligence or something that seems like it can be modeled. In the same way, that humans studied birds to understand flying, researchers are now studying the brain to create better artificial intelligence.

In the earliest stages of this research, they would try to imitate the brain, as people have tried imitating birds at first; and they failed similarly.

2.1 Artificial NNs . . . . .	7
Inspiration . . . . .	7
Regularization . . . . .	14
2.2 Fully Connected NNs .	15
Multi-Layer Perceptron	15
2.3 Convolutional NNs . .	17
Convolutions . . . . .	18
Pooling . . . . .	20
2.4 Recurrent Neural Networks . . . . .	21
2.5 Types of Learning . . . .	21
Supervised Learning .	22
Reinforcement Learning	22
Unsupervised Learning	22

There is current research on replicating the brains structure to the neuron level on hardware with more success [brainscales].

**Figure 2.1:** Otto Lilienthal with his flying apparatus. One example of how people failed in trying to imitate nature too closely. [https://www.deutschlandfunkkultur.de/geschichte-der-fliegerei-wie-die-976.de.html?dram:article\\_id=308043](https://www.deutschlandfunkkultur.de/geschichte-der-fliegerei-wie-die-976.de.html?dram:article_id=308043)



Notably McCulloch and Pitts (1943) even preceded Hodgkin and Huxley's (1952) Nobel price winning description of a neuron.

Warren McCulloch and Walter Pitts proposed one of the earliest models of an artificial neuron. They aimed at simplifying models of a biological neuron at the time.

### Biological Neuron

For this reason, the functioning of a neuron shall be described shortly.

Human cells that make up the brain are called **neurons**. They connect through dendrites, synapses, and axons. These connections allow neurons to exchange signals with other neurons.

A neuron receives signals through its dendrites, which all lead to the cell body (soma). The cell body accumulates these signals. If the summed value of the signals' potentials reaches a certain threshold, an action potential (spike) is generated. The spike then travels as a signal along the axon and its branches towards other neurons. Axons end in synapses that connect to other neurons' dendrites. Synaptic transmissions are usually mediated by chemicals and not by electrical signals. The chemical nature of the synapse allows it to forward either an excitatory or an inhibitory signal. Excitatory signals will bring the cell potential closer to the threshold, while inhibitory do the opposite [coloratlas].

Cell potentials are decreased by excitatory signals, as the action threshold sits below the resting potential.

What makes the brain so powerful, though, is not the neuron itself with its arguably simple structure but the vast network of billions of these neurons. Each neuron is connected to thousands of other neurons with which it communicates. How a signal is transported between neurons depends on the interplay of synaptic weights, neural connections, and the threshold of each neuron.

## Artificial Neuron

McCulloch and Pitts saw that powerful things could be achieved when connecting lots and lots of simple structures. Thus, they proposed an even simpler model of a neuron: They restricted their neuron to a binary state (on or off). Each neuron gathers signals from other neurons, which are either positive or negative. A neuron only becomes active if the number of incoming positive signals minus the number of negative signals exceeds the neuron's threshold. McCulloch and Pitts then also changed the highly parallel and complex nature of biological neural networks to a single layer feed-forward network architecture.

add equation with activation function, when neuron fires

In a feed-forward network, neurons are grouped into layers and operate in parallel within a layer. They do not interact within a layer. Each neuron in a layer is fed the same input signal (often described as an input layer). The neuron's state is then computed according to an equation such as equation ??.

The activation value of each neuron then defines an output.

As this model deviates from nature quite a bit, these structures are better referred to as **units** instead of neurons.

In a single layer architecture, a layer often consists of a single unit (see Figure ??). Basically, input signals come from one side, and output signals go out the other side, which can be expressed in a simple equation like equation ??.

This is not only done for practical reasons but also inspired by the observation of layered neuron structures in the brain.

The McCulloch-Pitts model is capable of emulating simple logical relations (AND, OR, NOT) but nor XOR which will be explained later.

## Perceptron

This McCulloch's and Pitts' is very much simplified yet comes with some weaknesses. Thus, the **perceptron** has been introduced by Frank Rosenblatt in 1957. Instead of having a binary unit, he came up with a linear threshold unit (LTU). The LTU allows for numeric instead of binary signals, which can be weighted with independent factors. Also, a unit's threshold/bias is parameterized as another weight with a constant input using the 'bias trick'. The resulting

equation for LTU with bias trick

equation for each LTU then reads: Equation equation ?? can then be formulated in a vectorized form such that:

vectorized LTU

In this form, the calculations for each unit become mathematically and computationally relatively easy. Each layer can be expressed as a vector of activations  $\mathbf{x}$ . Multiplying this vector with the weight matrix  $\mathbf{W}$  and applying the element-wise activation function returns the activations for the subsequent layer.

This capable yet straightforward description of a unit built the base for the first surge of interest in neural networks (connectionism).

The word perceptron describes the whole function  $f$  in equation equation ?? which can consist of many LTUs

## Mathematical Interpretation

The simplistic mathematical formulation of the perceptron suggests that there might be a mathematical meaning besides the biological analogy. Indeed, the perceptron is equal to the definition of a **binary linear classifier**.

A binary linear classifier categorizes inputs into two classes, hence binary. It does so by drawing a virtual hyperplane in input space and predicts the class for each input. The

decision depends on whether a point lies above or below this hyperplane.

As the hyperplane (or decision boundary) is quantified by a linear function, the classifier is described as linear.

A classic problem would be classifying email as spam. Given two data inputs (*i.e.* frequency of the words ‘weight loss’ and ‘invest’), the classifier has to make a decision. For this reason, the binary classifier defines a **decision boundary**. Any data point that lies above this decision boundary is classified as ‘spam’, any point beneath is classified as ‘not spam’. Since regular emails rarely use the two words, they do not classify as spam.

Other people *i.e.* nutritionists, on the other hand, will have the word weight loss come up more frequently. This means the decision boundary must differ for different users and their mail.

To compute where a point lies relative to the decision boundary, a data point’s value along each axis  $x_1, x_2$  is weighted individually  $w_1, w_2$ , and summed with a bias  $b$ . The result is checked whether it is above or below a threshold  $t$ .

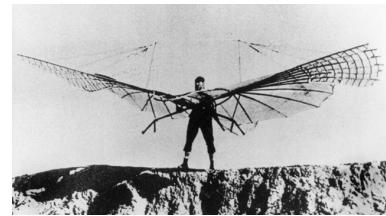
$$z = w_1x_1 + w_2x_2 + b \quad (2.1)$$

With the classification ‘spam’ if  $z > r$  and ‘not-spam’  $z \leq r$  (*in dubio pro reo*). Equally  $z - r > 0$  holds for spam as well such that the threshold can be brought into the equation and  $z$  is checked against 0.  $r$  can then be absorbed into the bias.

$$\rightarrow z = w_1x_1 + \dots + w_Dx_D + b - r = w_1x_1 + \dots + w_Dx_D + b' \quad (2.2)$$

For arbitrary dimensions  $D$  this becomes

$$z = w_1x_1 + \dots + w_Dx_D + b \quad (2.3)$$



**Figure 2.2:** As a simple example, this binary classifier has data on how often the words ‘weight loss’ and ‘invest’ appear in an email. Any time these two words appear too often, the data point is above the decision boundary. An email is then classified as ‘spam’

where  $\mathbf{x}$  and  $\mathbf{w}$  can be defined by vectors

$$z = w_1x_1 + \dots + w_Dx_D + b = \mathbf{w}^T\mathbf{x} + b \quad (2.4)$$

The decision boundary can be easily derived from this, since  $\mathbf{w}$  is the orthogonal vector to the hyperplane and  $\frac{b}{\|\mathbf{w}\|}$  is the displacement of the plane along  $\mathbf{w}$ .

For a simpler notation, one can define an additional 'virtual' input which has a constant value of 1 as  $x_0$ , the bias  $b$  can then be elegantly included into  $\mathbf{w}$  (same 'bias trick' as before)

$$z = w_0b + \mathbf{w}^T\mathbf{x} = w_0b + w_1x_1 + \dots + w_Dx_D = \hat{\mathbf{w}}^T\mathbf{x} \quad (2.5)$$

with  $x_0 = 1$  and  $w_0 = b$ . The output  $z$  can also be seen as a function of  $x$ . Then,  $w$  can describe any function linear in  $f$ .

This description also holds for perceptrons with more than one unit. In that case, the input vector  $\mathbf{x}$  and the weights  $\mathbf{w}$  become matrices with multiple column vectors.

$$\mathbf{x} \rightarrow (\mathbf{x}_1, \mathbf{x}_2, \dots) = \mathbf{X} \mathbf{w} \rightarrow (\mathbf{w}_1, \mathbf{w}_2, \dots) = \mathbf{W} \rightarrow \mathbf{z} = \hat{\mathbf{W}}^T \mathbf{X} \quad (2.6)$$

## Loss Function

With a mathematical definition at hand, the next step is to quantify the output. In order to train a classifier, an objective has to be formulated through a **loss function**. Usually, there is already a data set for the network to train on.

In the given example this would be mails which were read beforehand and then declared either 'spam'  $\tilde{y}_i = 1$  or 'not-spam'  $\tilde{y}_i = -1$ .  $\tilde{y}_i$  is called the **label** for a sample  $x_i$  with index  $i$ .

With the binary linear classifier, the decision boundary has been introduced (check  $z_i > 0$ ) to predict a label for any given sample. This is sufficient to predict a class, but much information is lost this way. For training and evaluation, the

information available in  $z$  should be used.

e.g., a large  $z_i$  implies that the data point  $x_i$  sits far from the decision boundary. Thus, the classifier is very sure of this classification. For  $z_i \approx 0$ , the classifier is not that sure, and for  $z_i = 0$  the classifier is indecisive. Ultimately,  $z_i$  can be seen as a score that is calculated for each data input.

The question then becomes how to quantify how well the classifier performs on given data. Hence, a loss function is defined, which measures the classifier's performance on the data. A popular choice is **least squares**, where the score's distance to the label is measured.

$$\mathcal{L}_{\text{LS}} = \sum_i (y_i - z_i)^2 \quad (2.7)$$

The value of the loss function becomes minimal for  $y_i = z_i, i = 1, \dots, N$ . Yet, this score function is especially susceptible to outliers which will cause the decision plane to skew towards outliers with  $z_i > 1$  or  $z_i < -1$

For this reason **support vector machines (SVM)** employ a **maximum margin** classifier.

A maximum margin classifier seeks to find a decision boundary which is as far from the closest representatives of each class as possible (see Figure ??). The maximum margin is defined as

$$\text{margin} = d_+ + d_- \quad (2.8)$$

with  $d_+$  the distance to the nearest training sample with class  $+1$  and  $d_-$  to closest training sample with class  $-1$ . Noticeably, this requires the data to be linearly separable, which means that a hyperplane must exist that perfectly separates the data according to its class. The margin becomes ideal for  $d_+ = d_-$ . Since  $w$  is orthogonal to the hyperplane,  $w$  can always be rescaled such that

$$d_+ = d_- = \frac{1}{\|w\|} \quad (2.9)$$

Additionally,  $\mathbf{w}$  can be chosen such that  $z = \mathbf{w}_i^T \mathbf{x}_i + b_i \geq +1$  for  $\tilde{y}_i = +1$  and vice versa for  $\tilde{y}_i = -1$ . Thus,

$$\tilde{y}_i z_i \geq 1 \quad (2.10)$$

will hold, for all inputs  $x_i$  with equality for points on the margin, as there is always at least one point of each class on the margin.

Thus,

$$d_- = d_+ = \frac{1}{\|\mathbf{w}\|} \quad (2.11)$$

and the margin

$$d_- + d_+ = \frac{2}{\|\mathbf{w}\|} \quad (2.12)$$

is maximized when  $\|\mathbf{w}\|$  is minimized.

Subsequently, the classification can be expressed as a relatively simple optimization problem.

$$\arg \min_{w,b} \frac{1}{2} \|\mathbf{w}\|^2 \quad (2.13)$$

under the constraints

$$\tilde{y}_i z_i = \tilde{y}_i (\mathbf{w}_i^T \mathbf{x} + b) \geq 1 \forall i \quad (2.14)$$

How to solve this optimization problem shall be explained in Subsection ??.

## Regularization

Regularizations, combined with loss functions, penalize possible solutions deemed wrong for any particular reason. A very common regularization is L2 regularization on the weights. It is expected that networks with very large weights generalize poorly. Hence, all weights are summed up and

multiplied with a regularization constant  $\lambda$  [grosse].

$$\mathcal{R} = \frac{\lambda}{2} \sum_{i=1}^D w_i^2 \quad (2.15)$$

## 2.2 Fully Connected Neural Networks

### Multi-Layer Perceptron

With the perceptron, which is capable of classifying any linear separable data, at hand, the question becomes: What are the limitations to this? Minsky and Papert found the limitations in 1969 with their book 'Perceptrons'. They outlined the limitations of perceptrons with the XOR problem [perceptron]. The problem becomes obvious when looking at the XOR problem in a 2D plane (see Figure 2.3)

As it has been explained in the previous section, the perceptron is equal to a binary linear classifier. As such, the perceptron can only classify linear separable data perfectly. Since the XOR problem can obviously not be solved with a straight line separating the two classes, the perceptron is also not able to compute such an operation.

This realization led to the first decline in interest in artificial neural networks.

Since then, there have been ways of solving this problem for SVMs by projecting the data into a higher dimensional space with various kernel functions [ommmer]. Another approach keeps the logic but goes from the shallow network approach to a deep neural network (DNN). **Deep Neural Networks (DNN)** are ANNs which consist of more than one hidden layer. A single perceptron may not be capable of computing XOR but it is capable of calculating AND, OR and their negated forms. By using one perceptron with two units to compute AND and OR, a second layer perceptron can in fact compute



**Figure 2.3:** OR and XOR operations visualized. The XOR problem cannot be solved by drawing a single line.

XOR

$$XOR(x, y) = AND(OR(x, y), NOT(AND(x, y))) \quad (2.16)$$

Thus, DNNs solve the XOR problem.

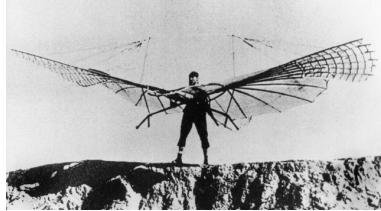
The ability to stack *and train* multiple layers in ANNs stems from the discovery of backpropagation which shall be explained in Sub-section ??.

Perceptrons, until then, could approximate only linear functions. Multiple layers of perceptrons now promise to approximate any higher degree function just as well. Thus, **Multi-Layer Perceptrons (MLP)** sparked new interest in the field of artificial neural networks.

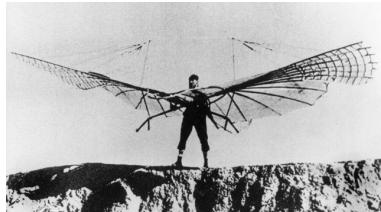
This interest also originated in the similarly layered structure that has been found in the brain.

Ultimately, MLPs really start to show the connected structure in a network that is typically expected.

MLPs are also called **fully connected networks** since each unit is connected to all unit in the previous layers as well as all units in the next layer.



**Figure 2.4:** The brains structure under a microscope



**Figure 2.5:** Layers of an MLP

## Activation Functions

These newfound capabilities for MLPs are not only restricted to binary operations but will translate into continuous space. In this case the hidden-layer perceptrons get stripped of their activation function. The activation function of a perceptron has been used, up until now, to make a class prediction  $y_i$  from a score  $z_i$ , which is also called pre-activation [GrosseNotes].

Replacing the step function with a linear activation function (*i.e.* identity function), each hidden-layer's perceptron would initially seem to increase the capabilities of the MLP. Unfortunately, this is not the case as any subsequent perceptrons with linear activations can be reduced to a single preceptron.

$$\mathbf{z}_2 = \mathbf{W}_2^T \mathbf{z}_1(\mathbf{x}) = \mathbf{W}_2^T \mathbf{W}_1^T \mathbf{x} = \mathbf{W}' \mathbf{x} \quad (2.17)$$

Consequently, the step-function that was used in the XOR problem played an important role. The reason for this is the non-linear nature of the step-function in contrast to any linear activation function. It can easily be shown that equation 2.17 does not hold if a non-linear activation function is used.

Thus, the question becomes which other activation functions there are that go beyond binary classification. An early popular choice was sigmoid functions (logistic function or tanh). Especially the logistic function is popular due to its similarity to the step-function amongst other things.

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (2.18)$$

Another popular choice are rectified linear units (ReLU) which are identical to a linear activation for  $z > 0$  but mimic the step function for  $z < 0$ . Basically, a ReLU suppresses the signal of a perceptron until it reaches the threshold of 0 and then forwards the signal unaltered.

Another activation function 'leaky ReLU' attenuates the signal below the threshold with a factor  $\alpha$  instead of fully suppressing it.

## 2.3 Convolutional Neural Networks

With fully connected networks at hand, it seems as if any complex function could be solved by just stacking enough hidden layers.

While this seems compelling at first, the issue of the computational burden arises very quickly. Especially for broad layers with many units the number of connections and weights becomes problematic. For two exemplary layers with 1,000 units each, there would be 1,000,000 connections as well as 1,000 biases. This causes two kinds of problems.



**Figure 2.6:** A sigmoid function. It saturates to 1 for very large inputs and 0 for very small inputs, similar to the step function.



**Figure 2.7:** ReLu activation function



**Figure 2.8:** leaky ReLu activation function with  $\alpha = 0.2$

1. Each connection represents a multiplication and summation to the activation value with an associated computational cost.
2. All connections have separate weights which must be trained, which requires a vast number of training data.

Another problem is invariance for spatially or temporally distributed data. If a pixel is to be shifted by a single pixel, or an audio track delayed by a second, the content does not change, and as such the result should not either. An FCN would likely over-fit the data and be susceptible to such variations.

Computer vision as a field is particularly affected by this issue. Since the input are often images with upwards of  $64 \times 64 = 4096$  pixels. Hence, a different solution is needed.

## Convolutions

A solution to the previously stated problem are convolutional layers.

A **convolution** is a mathematical operation which produces for any given point of a function  $f$  the weighted average of its surroundings. The way the surroundings are weighted is through a kernel function  $g$ . This principle can be translated into image space where a filter  $g$  is applied to each pixel and its surroundings.

A convolution is commutative  $f * g = g * f$ , associative  $f * (g * h) = (f * g) * h$  and distributive  $f * (g + h) = f * g + f * h$ .

Specifically, many filters of size  $k \times k$  are put on top of the image similar to tiles on a roof. Tiles may overlap but are evenly spaced over the area they cover. The spacing between the centers of each tile is called a stride  $s$  and will be assumed the same along each dimension. The return value of each filter then defines a new grid of values much like the original input. Typically, the output of such a convolution is smaller than the input, since the outermost filter must still fit fully into the image. Yet, it is possible to avoid this problem by padding an image with values such that the output has the same size.

These filters may seem very simple but they are able to capture very interesting properties in an image and transform them as well.

The easiest example would be a  $3 \times 3$  filter which blurs the image it is applied to. The filter given in Figure 2.10 results in such blurring since the signal of the pixel at the center is dominant but mixed with the signal of its surrounding pixels.

Another filter like Figure 2.11 will sharpen edges in an image since the signal of the surrounding pixels is subtracted from at the center pixel signal.

One can easily imagine that these kernels can also perform blurring only along one axis or sharpening gradients in one direction. As many such filters are imaginable, there are often many filters applied per layer called filter-banks. The output of each filter then defines an input channel for the next layer much like the three color channels in an image. With many channels convolutions can become even more complex as each filter takes all previous channels as input and combines them into a new channel.

Interestingly, the weights of each kernel may not be hand-crafted but can be learned similar to fully connected layers.

As convolutions qualify as linear functions, there exists an equal fully connected layer to each convolutional layer. The weights of such an FC layer would be sparse, though, which means there are many values equal to zero. This property makes it clear that convolutions are more efficient than fully connected layers by imposing some simple restrictions.

Similar behavior also occurs when comparing convolutional layers with different kernel sizes. By using a larger kernel, the 'field of view' (the area that each pixel in the result layer 'sees') increases as well. At the same time, a larger kernel increases the number of parameters significantly. The same field of view can also be obtained, though, when using two convolutions in series.

e.g. two layers with kernel size  $3 \times 3$  have the same field



**Figure 2.9:** A convolution in 1D space



**Figure 2.10:** 3x3 filter for blurring.



**Figure 2.11:** 3x3 filter for sharpening edges.

A convolutional kernel with the size of the input would be equivalent to a fully connected layer.

of view like a single layer with kernel size  $5 \times 5$  but fewer weights.

## Pooling

A sub-type of convolutional layers are pooling layers. Pooling layers have the purpose of downsampling an image or layer to a lower resolution. Basically, downsampling can already be achieved by choosing a convolutional layer with stride  $s > 1$ . Another way is to downsample an image without any learnable weights is by using one of two specific hand-crafted kernels.

The first option is average-pooling. By taking the mean of four neighboring pixels (each weight is  $\frac{1}{4}$ ) with a stride of  $s = 2$  in each channel, the input is scaled down by a factor of 2 along each axis. Figuratively, four pixels will be combined into a single pixel by taking their mean.

On the other hand, max-pooling does not take the average but only takes the maximum value of the four inputs. The result will then propagate only the most dominant signals.

Either way, pooling will result in loss of information but is often necessary to reduce the computational load. Especially when the number of channels increases for deeper layers, it is preferable to reduce memory usage and computational load along with the spatial size.

Typical CNN architectures will often stack several convolutional layers then apply a single pooling layer. Then another stack of convolutional layers is applied. This way, the input's spatial size is gradually reduced, while the number of channels is increased at the same time.

## 2.4 Recurrent Neural Networks

All previously presented networks have in common that they generate simple outputs  $y_i$  such as a class from simple inputs  $x_i$  such as the frequency of a certain word. In fact, images and even videos can act as inputs and outputs. Nevertheless, the larger the output, the harder it is to train. Recurrent neural networks are specialized for sequences of data as inputs and/or outputs.

Again, the email-spam problem can be manipulated such that a recurrent network fits the task. The classifier should check whether an email classifies as spam by analyzing whether an email is well written or just a bunch of buzzwords. Such an approach could be realized by taking the whole email as an input of an FC layer. However, similar to how convolutional layers are more efficient than FC layers for images, are recurrent networks more efficient for sequences.

Instead of feeding a whole sentence at once, RNNs take only one word at a time as input. Then the activation for a hidden layer is calculated and stored. During the next part of the input sequence, the previously stored activations provide an additional input besides the sequence snippet. Then new activations are generated, which are again stored for the next step of the sequence. This process is repeated until the sequence is complete, and then an output is generated.

Recurrent neural networks are prevalent in natural language processing as language can easily be visualized as a sequence.

## 2.5 Types of Learning

Neural networks are often categorized into different groups. One popular way of grouping is by categorizing the training according to the way the data is structured [grosse].

## Supervised Learning

If data is available, which already shows the network's desired behavior, training falls into the supervised learning category. Referencing the previous example, this would mean that there are emails available that are already labeled as 'spam' or 'not-spam'. The network can then imitate this behavior from the data.

Most problems fall into the supervised category as they are also the easiest to solve [grosse].

## Reinforcement Learning

In reinforcement learning, there is no such paired data available. However, the result can be evaluated how good it solves the problem. The resulting score is called a reward.

Most AIs trained to play games such as Google Alpha Go [alphago] are based on reinforcement learning. They also overcome the problem that calculating the reward is often based on undifferentiable outputs such as the score in a video game [grosse].

## Unsupervised Learning

Unsupervised learning has to deal with data that has no labels. Usually, this means finding patterns in the data on its own [grosse]. Considering the email-spam problem: A network is given all the email data without labels spam or not-spam. Ideally, the network is then able to find patterns in emails which allow grouping these emails. These groups are then evaluated, and hopefully, one group coincides with the 'spam' label or any other desired label.

# Optimization & Gradient Descent

# 3

## 3.1 Optimization Problem

With the foundation of neural networks in computer vision laid out, it seems as if one could easily start off with these right away. Yet, one very important part is still missing, which is the optimizer. Every neural network as it has been presented is based on weights which are meant to be trained. Without training, neural networks are of little use and would basically just generate pseudo-random output.

This is where optimizers come in. Optimizers have played a very important role for the development of neural networks and their comeback after the so-called 'AI winter'. The goal of an optimizer, is to find weights such that a neural network fulfills its intended task. Mathematically this equal to finding a minimum of a loss function which has been introduced in Section ??.

The optimizer is defined as an algorithms which searches for weights, which minimize said loss function. Hence, the network is optimized according to its loss function.

Very early implementations of networks like the perceptron or the McCulloch-Pitts model, had very straight forward algorithms for finding these weights. As these networks were inspired by nature and so was the optimizer. In accordance to observations made on neurons, weights were trained according to Hebb's rule [ommer]. Simply speaking it says "What fires together, wires together"[hebb]. Say, two neurons are in subsequent layers of a network and given some input, they are both activated. Then Hebb's rule states that the

3.1 Optimization Problem	23
3.2 Gradient Descent . . . . .	25
Stochastic Gradient Descent . . . . .	27
Momentum . . . . .	28
AdaM & AdaGrad . . . . .	28
3.3 Backpropagation . . . . .	29
3.4 Vanishing Gradients . . . . .	31
Activation Functions . . . . .	32
Normalization . . . . .	33
Residual Networks . . . . .	33

weight that connects them is increased like

$$\Delta w_{ij} = \eta \cdot a_i \cdot b_j \quad (3.1)$$

$$w_{ij} \rightarrow w_{ij} + \Delta w_{ij} \quad (3.2)$$

for  $\Delta w_{ij}$  the weight update,  $\mathbf{a}, \mathbf{b}$  the activation vectors for the respective layers and  $\eta$  the learning rate. Since the McCulloch-Pitts only knows the states active  $a_i = 1$  and inactive  $a_i = 0$ , this rule will increase the weights between synchronously firing neurons.

The **learning rate**  $\eta$  is an important parameter in many optimizers that decides how large an update step will be. A learning rate too high will often result in uncontrolled behaviour as the optimizer changes the weights too much and overshoots the optimal point. A too-low learning rate results in slow convergence and thus wasted time. Also, an optimizer with a low learning rate can more easily get stuck in local minima (more on that later).

Even though the conformity of Hebb's rule with nature has been shown [Lomo], it is rather obvious that it does not factor in the loss function. Other approaches have tried to improve Hebb's learning rule and also address problems like instability [ojas\_rule]. Still, Hebb's rule is clearly not well suited for finding suitable weights.

A slightly more intuitive, yet not inspired by nature, optimization algorithm is the perceptron algorithm. As the name states, it is used to update a perceptron.

Section ?? already stated, that the perceptron is equal to a linear classifier. The perceptron algorithm leverages this analogy and implements a very straight forward optimization algorithm:

If the perceptron correctly classifies a point, then do not change anything. If a point is wrongly classified, move the decision boundary in direction of the wrongly classified point.

This update rule will ideally move the decision boundary

repeatedly until the boundary crosses all wrongly classified points.

$$\Delta w_{ij} = \eta \cdot \frac{1}{2}(\tilde{y}_i^k - y_i^k) \cdot x_j^k \quad (3.3)$$

$$w_{ij} \rightarrow w_{ij} + \Delta w_{ij} \quad (3.4)$$

with  $\tilde{\mathbf{y}}$  the desired activations/classification,  $\mathbf{y}$  the observed classification,  $i, j$  the indices for units in each layer.  $\mathbf{x}^k$  is the input for a sample with index  $k$  from the previous layer.  $(\tilde{y}_i^k - y_i^k)$  is non-zero only if  $\mathbf{x}^k$  is misclassified and the factor then scales this term to 1. The sign of this term will indicate in which direction the decision boundary should move. The weight update also scales with  $x_j^k$  since the amplitude and the sign of the input are important to pick the correct weight.

If the optimization algorithm comes to a stop, since all classifications are correct, then the classification error has been minimized as well.

## 3.2 Gradient Descent

The idea of the perceptron algorithms can now be generalized further such that it uses the previously defined loss or cost functions (see Section ??). Instead of updating each neuron at a time, update everything!

Thus, the term  $\frac{1}{2}(\tilde{y}_i^k - y_i^k)$  in equation 3.3 must change. The purpose of this term is that it indicates in which direction the boundary should move. Now the question is how to determine this direction for an arbitrary loss function  $L^k$  that return a scalar score for a given sample  $\mathbf{x}^k$ .

By looking at a single weight at a time, the scalar output of the loss function depends on the scalar value of each single weight. Thus,  $L^k$  can be visualized in a 2D graph.

Looking at such a depiction (see Figure 3.1) the weight should change such that the loss function is minimized. This can be achieved by following the slope at position  $w$ , as it points



**Figure 3.1:** Example depiction of loss function  $L^k$  against  $w$ . At point  $w'$  the slope point into the other direction but still towards the minimum.

towards the minimum. The slope in this graph is equivalent to the negative gradient of the function at point  $w$ , which can be easily calculated, if the dependence between  $L(w)$  and  $w$  is known. By taking a small step along the direction the negative gradient points in (see previous section) the loss function becomes smaller but still has likely not reached its minimum. Thus, the gradient again calculated again for a new  $w$  and –again–  $w$  is updated in direction of the gradient.

Applying this procedure iteratively will reduce the loss function in a step-wise manner by descending along the gradient, thus the name **gradient descent**.

Besides giving information about the direction, the gradient also includes information about the steepness of the slope. For a continuous loss function  $L^k$  the gradient will become smaller the closer the  $w$  is to the minimum as  $\tilde{w}$ . Consequently, a larger gradient will indicate that a larger update step is possible.

The equation for calculating  $\Delta w_{ij}$  then becomes:

$$\Delta w_{ij} = -\eta \cdot \frac{\partial L^k}{\partial w_{ij}} \cdot x_j^k \quad (3.5)$$

$$w_{ij} \rightarrow w_{ij} + \Delta w_{ij} \quad (3.6)$$

with  $\frac{\partial L^k}{\partial w_{ij}}$  the partial derivative of  $L^k$  with respect to  $w_{ij}$ .

The same principle can be applied to many weights at the same time, where each gradient is calculated with respect to the respective weight.

Also, the gradient can easily be evaluated over the whole data set now  $k = 1, 2, \dots$  Thus, can be rewritten in vectorized notation:

$$\Delta w = -\eta \cdot \frac{1}{N} \sum_{k=1}^N \frac{\partial L^k}{\partial w} \cdot x^k \quad (3.7)$$

$$w \rightarrow w + \Delta w \quad (3.8)$$



Figure 3.2: 3D plot of loss function for two weights  $w_1$  and  $w_2$ .

Gradient descent allows to find minima for any loss function

as long as the output is differentiable in any weight  $w$ . This means that the same principle also applies to deep neural networks. However, there are two problems left.

First, how to obtain the gradient easily (especially for many layers), as analytical differentiation is hard to do for every weight, yet numerical differentiation [**numerical\_diff**] is very costly. The answer to this is automatic differentiation with back-propagation, which will be explained in the next section. The second problem is, how to avoid local minima. One could argue, that local minima would also solve the problem sufficiently and especially for many weights it is unlikely to find the global minimum. Yet, shallow local minima can still be problematic if the optimizer gets stuck in these.

## Stochastic Gradient Descent

Stochastic gradient descent (SGD) is one possible solution for this problem alongside another problem that comes with gradient descent. For large  $N$  it takes very long calculate all gradients and then take the mean. In comes SGD which does not take the mean over all samples but applies the gradient for each sample.

This introduces some randomness into the training as two samples can very quite much (*e.g.* photographs). As two samples vary in their loss functions, so does the loss hyperplane, and local minima change constantly. Yet this approach comes with a trade-off since the hyperplane changes so much. It could prohibit convergence as the changes become too large and there is no common gradient vector along which the optimizer can move.

The solution to this problem is **mini-batch gradient descent** as an intermediate solution. Instead of taking the mean over the whole data set, it takes the mean over a mini-batch, which consisting of a smaller number of samples. Typical choices are in the order of 32 samples per mini-batch plus-minus an order of magnitude.

## Momentum

Another solution is adding momentum to SGD [**momentum**]. By saving the direction of the previous SGD update and including it into its current step, the 'momentum' from the previous step is carried over [**ommer**].

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \Delta\mathbf{w}_{t-1} \Delta\mathbf{w}_t = -\eta \cdot \frac{1}{N} \frac{\partial L^k}{\partial \mathbf{w}_t} \cdot \mathbf{x}^k + \gamma \Delta\mathbf{w}_{t-1} \quad (3.9)$$

(3.10)

There are notable improvements to this method such as the Nesterov momentum update [**ommer**]

## AdaM & AdaGrad

**AdaGrad** takes a different approach in viewing the learning rate  $\eta$  as a hyper parameter which is tuned for each weight. Basically, weights which are updated frequently (often large  $\Delta w_{ij}$ ) should receive smaller updates. On the other hand rarely updated weights should receive larger updates. The result is "elementwise scaling of the gradient based on the historical sum of squares in each dim" [**ommer**]

$$\Delta\mathbf{w}_t = -\frac{\eta}{\sqrt{\epsilon I + \text{diag}(\mathbf{G}_t)}} \cdot \frac{1}{N} \frac{\partial L^k}{\partial \mathbf{w}_t} \cdot \mathbf{x}^k \quad (3.11)$$

$$\mathbf{w}_{t+1} \rightarrow \mathbf{w}_t + \Delta\mathbf{w}_t \text{ with } \mathbf{G} = \sum_{\tau=1}^t g_{\tau} g_{\tau}^T \text{ and } g_{\tau} = \Delta\mathbf{w} \quad (3.12)$$

AdaM is a more advanced version of AdaGrad and adds a version of momentum. The momentum and the sum of squares are replaced by exponentially decaying average  $m_t$  and  $v_t$  which serve the same purpose. Both  $m_t$  and  $v_t$  are corrected ( $\hat{m}_t, \hat{v}_t$ ) for biases at early iteration due to the

exponential average

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t v_t = \beta_1 v_{t-1} + (1 - \beta_2) g_t g_t^T \hat{m}_t = \frac{m_t}{1 - \beta_1^t} \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \Delta \mathbf{w}_t = -\frac{\eta}{\sqrt{\epsilon I + \text{diag}(v_t)}} \cdot \frac{1}{N} \quad (3.13)$$

$$\mathbf{w}_{t+1} \rightarrow \mathbf{w}_t + \Delta \mathbf{w}_t \quad (3.14)$$

AdaM along with SGD has in recent years become the de facto standard for training neural network in computer vision. Other optimization algorithms such as L-BFGS (Quasi-Newtonian second order optimization method) are rarely used nowadays.

### 3.3 Backpropagation

Backpropagation (or just backprop) was the major invention by **backprop** in the 1980's which ended the first AI winter caused by the **xOR** problem.

It allows to efficiently and accurately calculate gradients and thus train networks with many layers, which in turn 'solved' the **xOR** problem.

The way it works, backprop splits a large equation with many operations into its individual parts such that a tree structure emerges.

Then the chain rule is employed which states that

$$(f \circ g(x))' = f' \circ g(x) \cdot g'(x) \quad (3.15)$$

or

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} \quad (3.17)$$



**Figure 3.3:** Tree graph for function  $f$  in equation ??.

By leveraging equation 3.17 a function like

$$f(x_1, x_2) = x_1^2 \cdot (x_1 + x_2) \quad (3.18)$$

can be split into smaller pieces

$$\begin{aligned} f(x_1, x_2) &= z_0(x_1, x_2) & z_3(x_1, x_2) &= x_1 \\ z_0(x_1, x_2) &= z_1(x_1, x_2) \cdot z_2(x_1, x_2) & z_4(x_1, x_2) &= x_1 \\ z_1(x_1, x_2) &= (z_3(x_1, x_2))^2 & z_5(x_1, x_2) &= x_2 \\ z_2(x_1, x_2) &= z_4(x_1, x_2) + z_5(x_1, x_2) \end{aligned}$$

Each function  $z_i$  now involves only a single operation with either one or two inputs. This makes differentiating  $z_i$  with respect to its child node(s) very easy. e.g.

$$\frac{\partial z_1}{\partial z_3} = 2z_3 \quad (3.19)$$

(3.20)

The term child node is used since each node  $z_i$  (rectangular shape) is made up of lower level child nodes which are connected via an operation (circular shape).

The result of  $f$  can now be calculated by inserting some values for the end nodes  $x_1$  and  $x_2$  and then moving through the tree, from the bottom to the top. This is called a **forward pass** as this is also how such equations are solved normally. At the end, each node is now assigned a value as the forward comes through.

Now, how can the gradient of  $f$  with respect to the end nodes be obtained?

By using the chain rule and walking the graph backwards (hence *backpropagation*), the gradients of  $f$  with respect to each node can be calculated. The very first node to start with would be the top-most node  $z_0$  which is identical to  $f$ . Thus,  $\frac{\partial f}{\partial z_0} = 1$ , can be easily obtained. Calculating the gradient for  $z_0$ 's child nodes is almost equally easy as the chain rule is applied.

$$\frac{\partial f}{\partial z_1} = \frac{\partial f}{\partial z_0} \cdot \frac{\partial z_0}{\partial z_1} = 1 \cdot z_2 = z_2 \quad (3.21)$$

$$\frac{\partial f}{\partial z_2} = \frac{\partial f}{\partial z_0} \cdot \frac{\partial z_0}{\partial z_2} = 1 \cdot z_1 = z_1 \quad (3.22)$$

(3.23)

So now the gradient for nodes  $z_1$  and  $z_2$  is known as well. Notably, the gradients of  $z_1$  depends on  $z_2$  and vice versa. This is not a problem, as the gradient is not calculated in

the variable node, but in the operator node, which knows all three values  $z_0, z_1, z_2$  from the forward-pass. These calculated gradients are then forwarded to their respective nodes and the process start again.

This way gradients for all nodes a propagated (hence *backpropagation*) until they reach the end nodes. As these gradients  $\frac{\partial f}{\partial x_1}$  and  $\frac{\partial f}{\partial x_2}$  are those which were sought originally, backprop is complete.

$$\frac{\partial f}{\partial x_1} = \frac{\partial f}{\partial z_3} \cdot \frac{\partial z_3}{\partial x_1} + \frac{\partial f}{\partial z_4} \cdot \frac{\partial z_4}{\partial x_1} = 2z_2z_3 \cdot 1 + z_1 \cdot 1 = 2z_2z_3 + z_1 \quad (3.24)$$

$$\frac{\partial f}{\partial x_2} = \frac{\partial f}{\partial z_5} \cdot \frac{\partial z_5}{\partial x_2} = z_1 \cdot 1 = z_1 \quad (3.25)$$

(3.26)

In case of neural networks the end nodes are weights, parameters and inputs. Most often just the weight gradients are of interest, as the data set or *fixed* parameters should not change.

Interestingly, such a structure assigns each operation a second meaning besides its impact on the forward-pass result. Each operator now alters the gradient during the backwards-pass as well. The plus operator as an example just distributes the gradient unaltered to its child nodes. On the other hand, the gradient is scaled when it passes a multiplication operation. One important fact must be stretched, though: If the gradient becomes zero, all gradient that come after that node will also be zero!

## 3.4 Vanishing Gradients

As the previous sectioned ended, shall this section evaluate a little bit further.

Vanishing gradients are a problem that especially deep networks with many weights and many operations have to deal with. If any weight in a multiplication operation were to

become zero kills the gradient for subsequent operations along the backpropagation. Subsequently, the very first layers, close to the input, are likely to be subject to smaller gradients than the last few layers, which are very close to the loss function. Yet, as this is almost unavoidable, other causes for vanishing gradients are.

## Activation Functions

Activation functions are often a cause of vanishing gradients. The first examples of the McCulloch-Pitts neuron and the perceptron used a step-function as activation/classification function. Looking at the **step function**, the gradient is zero everywhere except at  $x = 0$  where the gradient is infinity. Such an activation function is obviously of no use and thus not found in any neural network.

Another favorite activation function for nature inspired neural networks are **sigmoid functions** (logistic function and hyperbolic tangent function). These seem to not suffer from the same problems, yet in their limit  $|x| \rightarrow \infty$  these functions become constant and thus the gradient becomes zero.

**Rectified Linear Units (ReLU)** were already presented in the previous chapter and their advantage is that they are very easy to compute. Also, for  $x > 0$  ReLU behaves like linear function and just passes the gradient on. For  $x \leq 0$ , though, the gradient is sadly zero.

**Leaky rectifies linear units (leaky ReLU)** do not suffer from this problem. They do not compare the linear value against 0 but against the same attenuated linear function.

$$\text{leakyReLU}(x) = \max(\alpha x, x) \text{ with } 0 < \alpha < 1 \quad (3.27)$$

Leaky ReLU is still easy to compute and thus a popular choice in neural networks.

## Normalization

Another way of avoiding vanishing gradients is **normalization** often just called norm. Normalization shifts and scales activations such that they represent a normal distribution with  $\mu$  and  $\sigma$ .  $\mu$  and  $\sigma$  are both trainable parameters, also called **affine parameters**. By doing so, the activations typically are very similar distributed across layers and there are fewer activations with very large amplitudes. Having well distributed amplitudes helps to avoid gradients exploding or vanishing gradients.

The most common types of normalization are **batch norm** and **instance norm**. Instance norm takes activations for each sample in each channel in a layer and applies normalization over this set [IN].

Batch norm normalizes over the channel *and* the mini-batch to obtain broader statistics [BN].

## Residual Networks

Residual networks are another take on tackling the vanishing gradient problem in very deep neural networks. They introduce skip connections, which forward the signal of a layer add it to the activations of a different layer down the line. This way the gradients get also forwarded during backprop and maintain their amplitude [res\_net].



# 4 Models

After the previous two chapters focused on basic ideas and mathematics behind neural networks, this chapter will introduce basic model architectures and building blocks of many neural networks. Not every new publication reinvents the wheel and often the repeated use of merited building blocks can be observed.

It is particularly interesting that most networks make a single choice for a basic building block throughout the network. The choice affects the properties of convolution, normalization and activation layers as a sequence of these makes up such a block. Mainly, parameters like kernel size, normalization function and activation function are picked beforehand and remain the same for the entire network. Arguably, one could say, that it shows some similarities to the way network motifs make up networks [[uri\\_alon](#)].

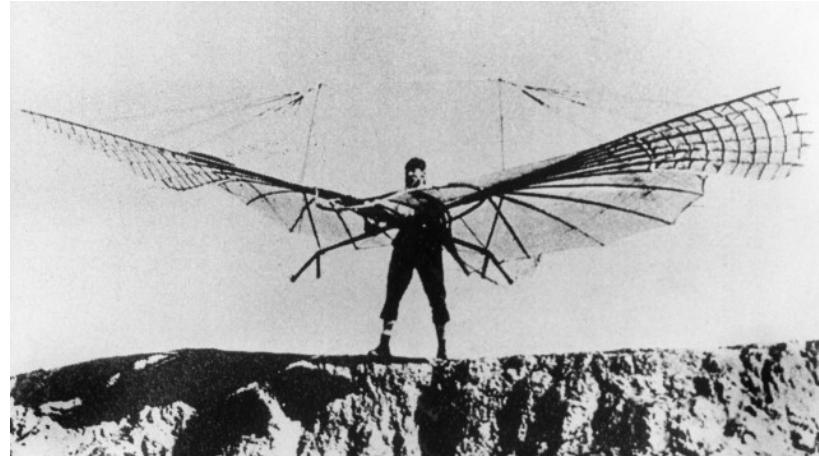
However, the focus of this chapter is the introduction of advanced model architectures like adversarial models or autoencoder. To aid this cause, this chapter starts with simpler architectures.

## 4.1 Discriminators & Classifiers

Discriminators and Classifiers are prime examples in many books which deal with neural networks. Both architectures are designed to take high-dimensional inputs and produce low-dimensional outputs from these. These outputs then represent a class prediction.

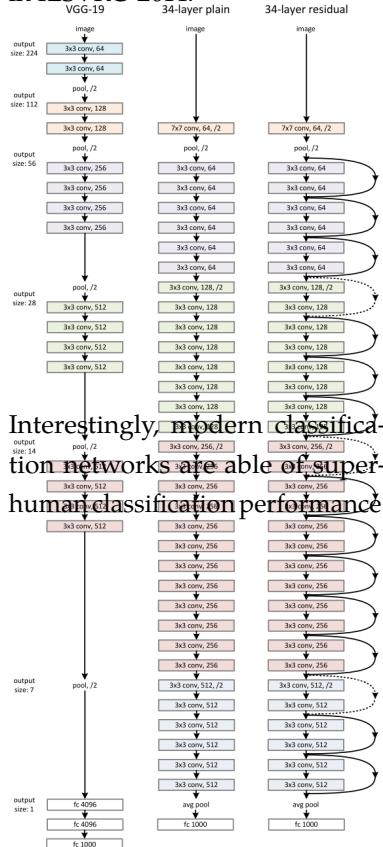
For discriminators these class predictions only include two classes (often ‘real’ or ‘fake’). Thus, both often display the same architecture.

4.1 Discriminators & Classifiers . . . . .	35
4.2 Generators & Decoders . . . . .	36
4.3 Auto-Encoders . . . . .	37
4.4 Generative Adversarial Networks . . . . .	39
Flavors of GANs . . . . .	40



**Figure 4.1:** Architecture of a VGG16 network.

ImageNet consists of 15 million high-resolution images in 22,000 categories. A subset of 1.2 million  $256 \times 256$  pixel images in 1,000 categories is used for the ILSVR challenge. VGG16 was runner-up in ILSVRC-2014.



Interestingly, modern classification networks are able of super-human classification performance [ILSVRC2014].

Probably the most common architecture is the VGG-architecture. This architecture was developed for the ImageNet classification challenge [[imagenet](#)]. The VGG network exists in different variants, which feature differing numbers of layers. VGG16, for instance, features 16 layers with weights.

VGG16 is exemplary for many classifier structures, as the spatial resolution is slowly reduced while the number of channels increases at the same time. This does not happen continuously but in a step-wise manner. After every two or three blocks (convolutional layer + activation layer) there is a pooling layer which reduces the resolution by  $\frac{1}{2}$  along each axis [[VGG](#)]. Every time the resolution is reduced by  $\frac{1}{4}$  the number of kernels for the following convolutions is doubled. In the end there are a couple fully connected layers which interpret the final features and generate a prediction. Figure 4.1 shows the networks architecture.

Another famous classification network is ResNet [[resnet](#)]. Presented only a year after VGG16, ResNet features many more by using residual connections.

## 4.2 Generators & Decoders

Decoders and generators are meant to generate images or high-dimensional representations from low dimensional input. For decoders the relation between input and output is often known, upscaling an image in a specific way. If there

**Figure 4.2:** Comparison of VGG19 architecture and ResNet architecture.

are no specific relations between input and output (input is noise), decoders are rather called generators.

Often the architecture for both is inverse to that of classifiers. The input is first processed by fully connected layers and then rearranged such that it can serve as an input to a convolutional layer. In contrast to classifier architectures like VGG16, feature decoders upsampling layers. Upsampling can be achieved through established upscaling methods like nearest neighbor upscaling, bilinear upscaling or bicubic upscaling. Sometimes transposed convolutions or subpixel upscaling are employed, although they add additional weights to be trained.

Usually, generators mirror the relation of scaling, channel size and number of convolutional blocks of VGG networks but in an inverse direction.

Both decoders and generators are more often used in larger structures like auto-encoders or adversarial networks than as stand-alone networks.

### 4.3 Auto-Encoders

Auto-encoders combine decoders with encoders. Encoders are similar to classifiers but instead of returning a class, they return a broader lower resolution representation of its high-resolution input. This is the inverse to a decoder's job. Auto-encoders combine an encoder with a decoder such that the low dimensional output of the encoder acts as input to the decoder. The decoder then tries to reconstruct the original input as good as possible. For this reason the encoder must compress all the information it gets as an input, such that the decoder can reconstruct the same information from the compressed representation.

The main goal of auto-encoders is to learn such a compressed/more efficient representation that still holds all relevant information. A very common use-case for auto-encoders are images of faces. Usually these faces are pre-processed

such that they are always aligned similarly in every image. Then a lower dimensional representation only needs to represent *e.g.* the eye-color since the position of the eyes should be the same in all images.

This example also points out another goal of auto-encoders, which is to make the compressed representation interpretable and/or controllable. This would mean, changing one parameter in the **feature space** *i.e.* the lower dimensional representation will equate to changing the color of the hair.

The now introduced features space allows for lots of research such as metric learning.

A popular variant of auto-encoders are **variational auto-encoders (VAE)**. VAEs add an additional sampling step between the output of the encoder and the decoder. Thus, the encoder predicts means and standard deviations each dimension in feature space. These parameters are then used to adjust a multi-dimensional Gaussian distribution from which the inputs of the decoder are sampled.

Such an additional layer of stochasticity is meant to solve a huge short-coming of auto-encoders. The features space of auto-encoders is often highly specific to the data. This means that there are many places in feature space which are not explored during training and it is not easy to find the meaningful areas in feature space. Taking such an undefined feature representation decoder input would result in equally undefined output. Thus, it is hard to find valid feature representations without using the encoder.

By adding a sampling layer in VAEs, the decoder is suddenly forced to deal with stochastic inputs. This means that any point in feature space could be an input by pure chance. Thus, the decoder gets more robust against various feature inputs. To further reinforce such a behavior, and additional constraint is put on the sampling parameters, that the encoder predicts. The sampling parameters should be as close to a standard Gaussian ( $\mu = 0, \sigma = 1$ ) as possible. This is enforced by calculating the Kulback-Leibler Divergence between the

encoder output and a standard Gaussian. The resulting network then allows to generate samples with the decoder and explore the feature space more easily.

## 4.4 Generative Adversarial Networks

Generative adversarial networks (GAN) were presented by GAN in GAN[GAN]. Nowadays, GANs are widely acclaimed and build the foundations for many approaches in computer vision and machine learning. The researchers behind GANs were even lauded with the prestigious Turing Award [turingaward].

The basic idea behind GANs is to combine a generator with a discriminator. The generator's job is to generate data from noise. Ideally, the generated data should fit the data which is present in a data set. The discriminator ensures, that this is the case by distinguishing generated samples ('fake') from data set samples ('real'). During training, the generator and the discriminator are in constant competition. While the discriminator tries to better distinguish fake and real data, the generator tries to come up with samples that fool the discriminator. This can be described by the following equation for a generator  $G$ , a discriminator  $D$ , some data  $x \in X$ , and sampled noise  $z \in Z$ .

$$G = \arg \min_G \mathbb{E}_Z [\log D(G(z))] \quad (4.1)$$

$$D = \arg \min_D \mathbb{E}_X [\log D(x)] + \mathbb{E}_Z [\log(1 - D(G(z)))] \quad (4.2)$$

As this game continues the quality of the generated data becomes better until the generated data becomes indistinguishable from real data. This way, the generator can for instance create images of faces, which have never been seen before, but look just like any other face. Such a simple principle can be taken even further, which is the reason why GANs can be found pretty much in any field of computer vision nowadays.

The Kulback-Leibler divergence is a measure for two distributions how much they differ. For two Gaussian distributions, the Kulback-Leibler divergence takes a well-defined form entirely depending on  $\mu_1, \sigma_1$  and  $\mu_2, \sigma_2$  [KLdiv].

Other researchers claim to have come up with the idea of GANs earlier, which is kind of a hot topic in the machine learning community. For safety measures the supposedly previous claim shall be cited as well [schmidhuber].

## Flavors of GANs

Since GANs have had such an impact, there has naturally been more research to improve GANs in general.

Some approaches are able to combine GANs with auto-encoders [AEGAN] or VAEs [VAE-GAN] and gain performance improvements from this. Other approaches add relevant information to the noise input and get the generator to act more like a decoder [cGAN]. Yet other approaches concentrate on improving the loss function, such that GAN training becomes more stable [WGAN]. Again, other approaches improve the quality for high-resolution images by generating samples at different scales [MSG-GAN] or progressively increasing the image scale [ProGAN, StyleGAN]

## Relativistic GANs

One particular approach to improving the loss function finds use in this thesis and thus shall be presented.

Relativistic Standard GAN (RSGAN) and Relativistic average standard GAN (RaSGAN) are two implementations presented by RGAN. They claim that by modifying the loss function the stability of training as well as the quality of generated samples can be improved [RGAN].

The main idea behind this approach is, that in a standard GAN approach the discriminator is more focussed on classifying real samples correctly. If the discriminator is to classify all real samples correctly, then there is no more gradient for real samples as  $\mathbb{E}_X[\log D(x)] \rightarrow 0$ . Then the discriminator mostly learns from generated fake samples and probably over-fit these.

A relativistic loss function avoids this by measuring the distance between real predictions and fake predictions with a

Standard GAN means the GAN as presented in equation ??.

sigmoid function.

$$G = \arg \min_G -\mathbb{E}_{X,Z}[\log(\text{sigmoid}(D(G(z)) - D(x)))] \quad (4.3)$$

$$D = \arg \min_D -\mathbb{E}_{X,Z}[\log(\text{sigmoid}(D(x) - D(G(z))))] \quad (4.4)$$

(4.5)

This more symmetric formulation is meant to push the generator closer to generating samples that follow the data distribution.

RaSGAN advances on the idea of RSGAN by reducing the variability that comes from comparing to single samples  $x \in X$  and  $G(z), z \in Z$ . Instead, the mean of real samples and generated samples is calculated across the mini-batch and then compared to each sample.

$$G = \arg \min_G -\mathbb{E}_Z[\log(\text{sigmoid}(D(G(z)) - \mathbb{E}_X[D(x)]))] - \mathbb{E}_X[\log(1 - \text{sigmoid}(D(x) - \mathbb{E}_Z[D(G(z))]))] \quad (4.6)$$

$$D = \arg \min_D -\mathbb{E}_Z[\log(1 - \text{sigmoid}(D(G(z)) - \mathbb{E}_X[D(x)]))] - \mathbb{E}_X[\log(\text{sigmoid}(D(x) - \mathbb{E}_Z[D(G(z))]))] \quad (4.7)$$

(4.8)

This further improves the quality of generated samples in adversarial training.



# Artistic Computer Vision

Applying computer vision techniques to images with artistic content is an interesting and challenging task at the same time. Due to the often observed change in appearance in many artistic images (*e.g.* abstraction), existing CV pipelines for *i.e.* classification usually do not work as intended. This makes it all the more interesting to see, how these pipelines are affected by this domain gap and whether this can be overcome.

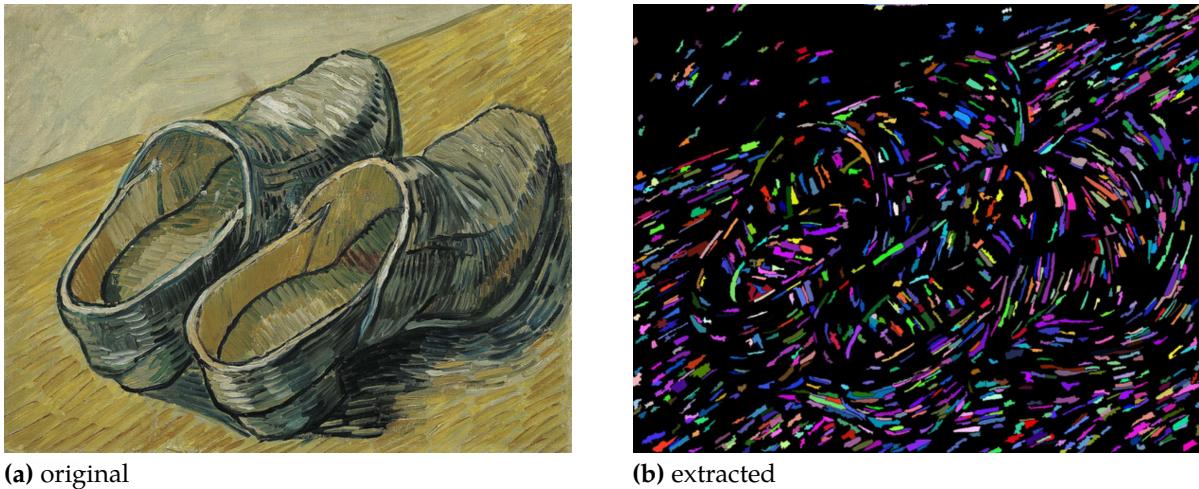
Yet, at the same time artistic images combined with computer vision often results in interesting and fascinating applications for users. A popular example for this is **style transfer**, where normal images are filtered such that they look like paintings. Another example is art created with the help of neural networks [**GANart**]. Research in the artistic field also incorporates techniques from computer vision and natural sciences to aid the detection of forgeries [**picassomatisefake**] [**guardianarticle**].

These examples show that computer vision can play an important role in the artistic domain and the other way around. Since this thesis is deeply anchored at the intersection of both fields, different problems and their proposed solutions shall be presented.

## 5.1 Brushstroke Extraction

Brushstroke extraction is one way of retrieving data for image analysis. Other approaches rely on patch-based analysis via neural networks [**patchbased**] or texton statistics [**textons**]. Yet, the most impressive results have been obtained when analyzing drawings on the pencil-stroke level

5.1	Brushstroke Extraction	43
5.2	Style Transfer . . . . .	45
	Early Approaches . . . . .	46
	Neural Style Transfer . . . . .	47
	State of the Art . . . . .	52
5.3	Painterly Rendering . . . . .	56
	Stroke-Based Rendering	57
	Drawing Networks . . . . .	60
	Genetic Algorithms . . . . .	63
5.4	Conclusion . . . . .	64

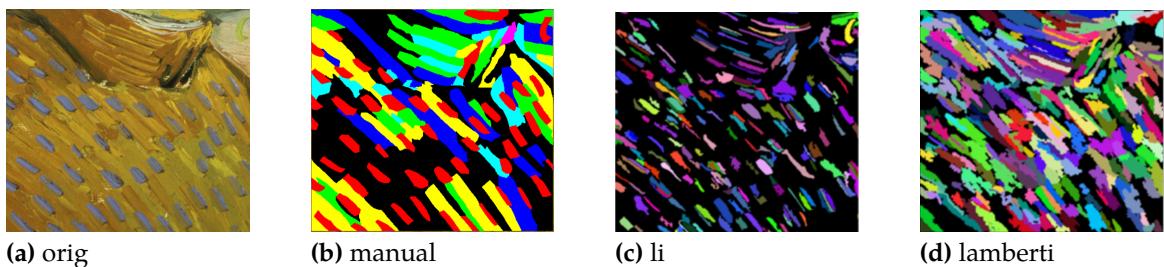


**Figure 5.1:** Original painting and extracted brushstrokes by Li *et al.*

[picassomatissefake]. Similar (if few) approaches exist for paintings and show that such approaches are feasible and open the door for one major application of brush stroke extraction.

Li *et al.* were the first to extract brush strokes from many paintings on a larger scale and combine it with forgery detection [1]. They introduced an algorithm that leverages edge detection as well as image segmentation to find connected components (regions enclosed by a continuous edge) in an image. These connected components then serve as stroke candidates. Each candidate is further processed to a single-line representation and then checked for various properties such as width-to-length ratio. As a result Li *et al.* claim that they detect 60% of brush strokes which were manually labeled on average. They further obtain features from their line representation and are able to achieve respectable results in forgery detection and classification of artistic periods.

Lamberti *et al.* later took an advanced approach and specifically modeled a brush stroke detection algorithm on the pixel level. They propose to randomly sample evenly spaced pixels over the image and build up connected regions by adding adjacent pixels if they are similar enough. This is repeated until the process comes to a stop or the brushstroke grows too large. If the region is too large, the similarity metric



**Figure 5.2:** Comparison of techniques on a cut out patch from van Gogh’s *The Little Arsiennne* (a). Manually labeled brushstrokes (b) [1], extracted brushstrokes by Li *et al.* (c) [1], extracted brushstrokes by Lamberti *et al.* (d) [2]

is adjusted and the growing process starts from scratch. If region are too small the whole seed is discarded. Candidate brushstrokes of the correct size are then approximated with an ellipsoid and shape constraints as well as alignment with the connected region are tested. Lamberti *et al.* specifically fine-tune parameters in their algorithm to improve compliance with manually labeled brushstrokes. In the end their results fare worse on metrics proposed by Li *et al.* such as detection rate but better in their own metrics. Notably, Lamberti *et al.* seems to be the only publication solely aimed at brushstroke extraction.

## 5.2 Style Transfer

Style transfer is probably the most popular and advanced field in artistic computer vision. The goal of style transfer is to generate an image representation which resembles the content of an input image but features the style of another image or image collection. With the focus lying on altering the image on the pixel-level such that the stylized image and the style-template seem to belong to the same group of images (*e.g.* images drawn by van Gogh)

This definition of style is hard to quantify, and thus relies on human judgement. Subsequently, many publications offer expert- or non-expert-based quality estimates, which are hard to compare and thus shall not be further mentioned. Then again, such an open definition of style allows the same

techniques to be applied to related problems such as texture transfer or image-to-image translation.

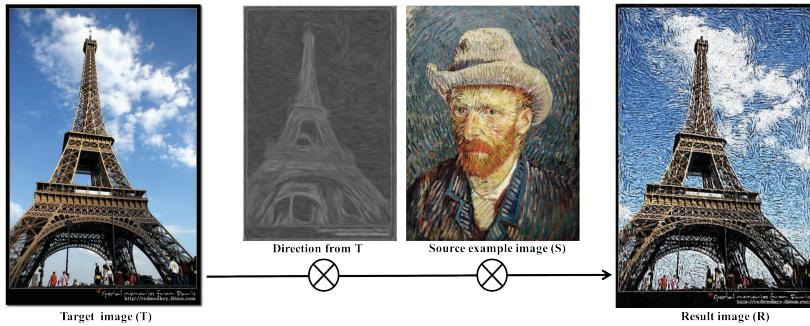
The high quality that recent style transfer has achieved over the years requires to compare any stylization approach with this field as well.

## Early Approaches

As one of the earliest approaches, Hertzmann *et al.* proposed a fast way of stylizing images by using trainable filters, while contemporary approaches relied on painterly rendering (see Section 5.3). They call their technique ‘image analogies’ since the transformation between two training images is analogously applied to a test image.

Basically, for creating the impression of brush strokes, a photograph  $A$  and a painting of that photograph  $A'$  are required. As this is rarely the case, **imageanalogies** showed that using anisotropic diffusion works also reasonably well to generate  $A$  from  $A'$ . Then an algorithm searches locally around a spatial position  $p$  for the best filter parameters  $F(p)$  that transform  $A(p)$  into  $A'(p)$ . By using another search algorithm to match similar regions  $B(p)$  and  $A(p)$ ,  $F(p)$  can then be used to transform  $B(p)$  into  $B'(p)$ . Thus  $B'$  – a stylized version of  $B$  – can be obtained by transforming  $B \rightarrow B'$  analogously to  $A \rightarrow A'$ .

Other works like Lee *et al.* [4] have built on this basic idea and advanced versions of it [**fasttexturetransfer**] [4]. Besides only matching local regions according to their pixel values, Lee *et al.* obtain a flow map (which they call ‘directions’), that is based on the content image’s gradient. This flow is then also matched against the style image. They are able to transfer texture (but only texture) from a given painting to a content image reasonably well.



**Figure 5.3:** Target image  $T$  is combined with a dictional map specially obtained from  $T$  and a style image  $S$ . The result maintains the direction's flow while presenting the texture from  $S$ .

## Neural Style Transfer

Despite some previous admirable results[4], style transfer really gained traction in 2015 with the publication of ‘A Neural Algorithm of Artistic Style’ by Gatys *et al.* It was the first approach to transfer the style as more than just texture of one image to another and at the same time maintaining a high contextual fidelity. In retrospective, this work really kicked started neural style transfer for the following years. Because of the weight this publication has and some similarities in the approach of this thesis this work is presented more extensively.

Gatys *et al.* themselves pinpoint the novelty of their approach as ‘manipulations in feature spaces’ as opposed to previous approaches that ‘directly manipulate the pixel representation of an image’[5]. They use existing neural architectures and extract information in two separate ways, such that content and style can be separated.

Previous works already used **perceptual loss** to gain information on content in an image [**percep\_loss**], or check whether two images have the same content [**other\_percep\_loss**]. Perceptual loss is based on the VGG-19 architecture [VGG] which is a deep CNN trained for object classification on ImageNet [**imagenet**]. Gatys *et al.* argue that the network’s layer activations increasingly respond to the content when following the networks’ hierarchy. So much even, that it is possible to reconstruct the content of an image by imitating the activations of one such layer (same content does not imply pixel-wise identity).

For reconstruction of an image's content, gradient descent is performed on a white noise image. The gradient descent aims to minimize the perceptual distance between the reconstruction and the target image. Perceptual distance is defined as the L2-distance between the activations of two images in a deep layer of the VGG-network.

An image vector is often used to represent a 2-dimensional image as a single column vector such that spatial information is neglected at first. Yet, the position inside the image vector still corresponds directly to a specific position in the original 2D representation

For an image vector  $\mathbf{x}$  with  $\mathbf{x} \in \mathbb{R}^{M_0}$ ,  $M_0 = H_x \dot{W}_x$ , a layer  $l$  of the network has  $N_l$  feature maps of size  $M_l$ . In this case  $M_l$  is equal to the height times the width of the feature map of the  $l$ -th layer  $M_l = H_l \dot{W}_l$ . The activations of the  $i$ -th filter ( $i \in N_l$ ) at position  $j$  ( $j \in M_l$ ) at layer  $l$  can then be represented by a matrix  $F(\mathbf{x})_{ij}^l \in \mathbb{R}^{N_l \times M_l}$ . The perceptual distance is then defined as

$$d_{\text{percep}}(\mathbf{x}, \mathbf{y}) = \sum_{i,j} (F(\mathbf{x})_{ij}^l - F(\mathbf{y})_{ij}^l)^2 = \left\| \mathbf{F}(\mathbf{x})^l - \mathbf{F}(\mathbf{y})^l \right\|_2^2 \quad (5.1)$$

, which allows to define the perceptual loss or content loss as

$$\mathcal{L}_{\text{content}} = \frac{1}{2} \sum_{i,j} (F(\mathbf{x})_{ij}^l - F(\mathbf{y})_{ij}^l)^2 = \frac{1}{2} \left\| \mathbf{F}(\mathbf{x})^l - \mathbf{F}(\mathbf{y})^l \right\|_2^2 \quad (5.2)$$

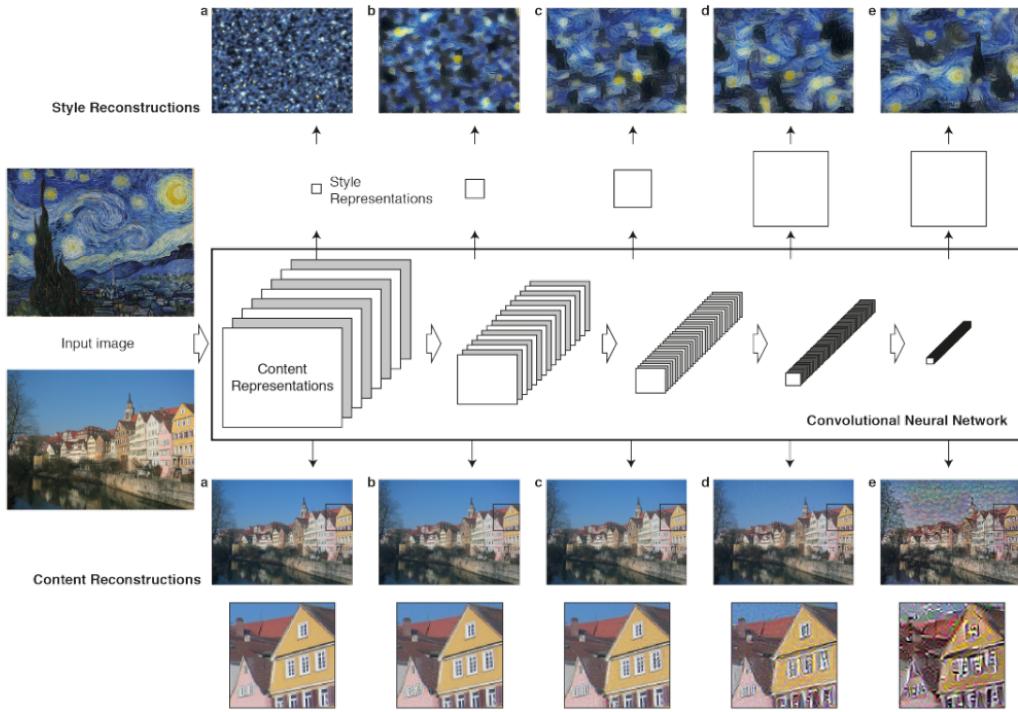
The factor  $\frac{1}{2}$  will cancel out when deriving  $\mathcal{L}_{\text{content}}$  with respect to  $F(\mathbf{x})_{ij}^l$ .

Minimizing the content loss between two images by using gradient descent restores the content of an image (see Figure 5.5).

As approximating the content of an image is possible, the question is whether the same is possible for style. Gatys *et al.* again turn to a pre-trained VGG network for this. They identify style as texture and thus search for a feature space that captures **texture** much like perceptual loss captures content. Subsequently, Gatys *et al.* propose the use of **Gram matrices** as they capture the correlations of feature-activations over their spatial extent.

The Gram matrix of a given matrix  $\mathbf{A}$  is the inner product of all column vectors in  $\mathbf{A}$ .

$$G = \langle a_i, a_j \rangle = \mathbf{A}^T \mathbf{A} \text{ if } a_1 \dots a_j \text{ are column vectors of } \mathbf{A} \quad (5.3)$$



**Figure 5.4:** Reconstructions of content(bottom) and style(top) using different layers. [5]

The resulting Gram matrix  $G$  now has the form  $j \times j$  and captures texture information but no longer the global content.

Arguably, style/texture is very complex and exists at various scales at the same time which Gatys *et al.* address by using many layers at the same time. As these layers sit at increasing depths their field of view increases as well and each layer captures information at a different scale. Early layers will tend to hold small scale information, later layers will hold larger scale information.

Gatys *et al.* first compute the Gram matrices of each layer  $l$  for both the target-style image and the current image. Then they use the L2 distance metric to measure the discrepancy

between them.

$$G(\mathbf{x})^l = \frac{1}{(2N_x^l M_x^l)^2} F(\mathbf{x})^{lT} F(\mathbf{x})^l \quad (5.4)$$

$$G(\mathbf{y})^l = \frac{1}{(2N_y^l M_y^l)^2} F(\mathbf{y})^{lT} F(\mathbf{y})^l \quad (5.5)$$

$$d_{\text{style}}^l(\mathbf{x}, \mathbf{y}) = \|G(\mathbf{x})^l - G(\mathbf{y})^l\|_2^2 \quad (5.6)$$

$$(5.7)$$

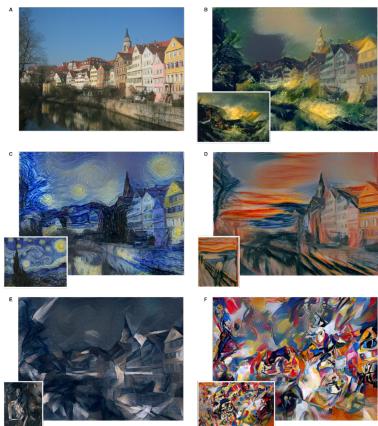
The denominator of  $\frac{1}{4N_x^l M_x^l}$  is squared since the Gram matrix is the product of a matrix with its transposed self.

The style distances at each layer are then weighted and summed up to make up the style loss:

$$\mathcal{L}_{\text{style}} = \sum_l w_l d_{\text{style}}^l(\mathbf{x}, \mathbf{y}) \quad (5.8)$$

This style loss can again be used together with gradient descent in order to check whether it is possible to reconstruct the texture of an image much like the content of an image. Figure 5.5 shows that it is in fact possible to reconstruct the texture of the image at various scales. Specifically the local consistency of each texture becomes larger, the deeper the layer sits.

Gatys *et al.* now combine the losses for a content image  $\mathbf{c}$  with a style image  $\mathbf{s}$  and optimize  $\mathbf{x}$  in the same way previous reconstructions have been obtained.



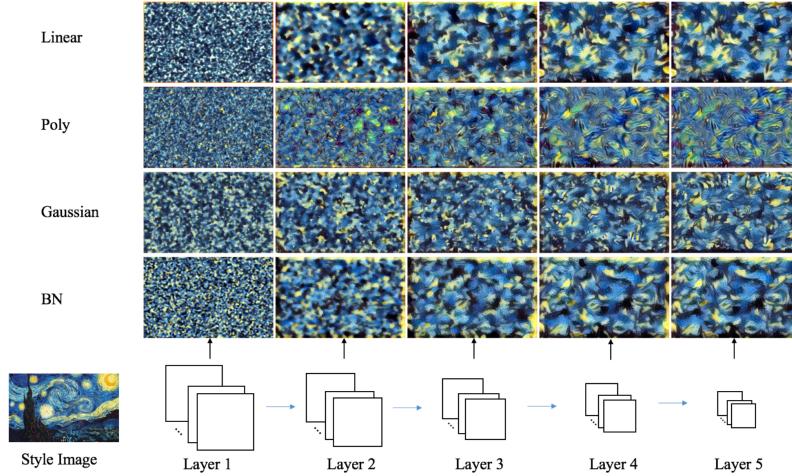
**Figure 5.5:** Style transfer examples by Gatys *et al.* [5]

$$\mathcal{L}_{\text{total}} = \lambda_{\text{content}} \mathcal{L}_{\text{content}}(\mathbf{x}, \mathbf{c}) + \lambda_{\text{style}} \mathcal{L}_{\text{style}}(\mathbf{x}, \mathbf{s}) \quad (5.9)$$

The result of this can be seen in Figure ??

## Follow-Up Research

There has been some follow-up research on Gatys *et al.*'s work which addresses mainly how the style loss works.



**Figure 5.6:** Reconstructed textures for Starry Night using different kernel functions  $k$  [6]

Li *et al.* have shown that the style loss is equivalent to calculating the **maximum mean discrepancy (MMD)** between the features of each layer [6]. MMD is a test-statistic for a null hypothesis  $p = q$  with data  $X = \{x_i\}_{i=1}^n$ , sampled from  $p$ , and  $Y = \{y_j\}_{j=1}^m$ , sampled from  $q$ , at hand. It can be used as a difference measure and vanishes only if  $p = q$ .

$$\text{MMD}^2[X, Y] = \frac{1}{n^2} \sum_{i=1}^n \sum_{i'=1}^n k(\mathbf{x}_i, \mathbf{x}_{i'}) \quad (5.10)$$

$$+ \frac{1}{m^2} \sum_{j=1}^m \sum_{j'=1}^m k(\mathbf{y}_j, \mathbf{y}_{j'}) \quad (5.11)$$

$$- \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m k(\mathbf{x}_i, \mathbf{y}_j) \quad (5.12)$$

MMD can be based on different kernel functions  $k$  and Li *et al.* have shown that the style loss is equivalent to the squared MMD with a polynomial kernel. Consequently they were able to show, that style transfer works with different kernel functions as well and even by explicitly matching the batch statistics (see Figure 5.6):

$$d_{\text{style}}^l(\mathbf{x}, \mathbf{y}) = \frac{1}{N_l} \sum_{i=1}^{N_l} \left( (\mu_{F(\mathbf{x})^l}^i - \mu_{F(\mathbf{y})^l}^i)^2 + (\sigma_{F(\mathbf{x})^l}^i - \sigma_{F(\mathbf{y})^l}^i)^2 \right) \quad (5.13)$$

**LenDu** tested whether pre-trained weights play an important role when performing style transfer. He was able to show

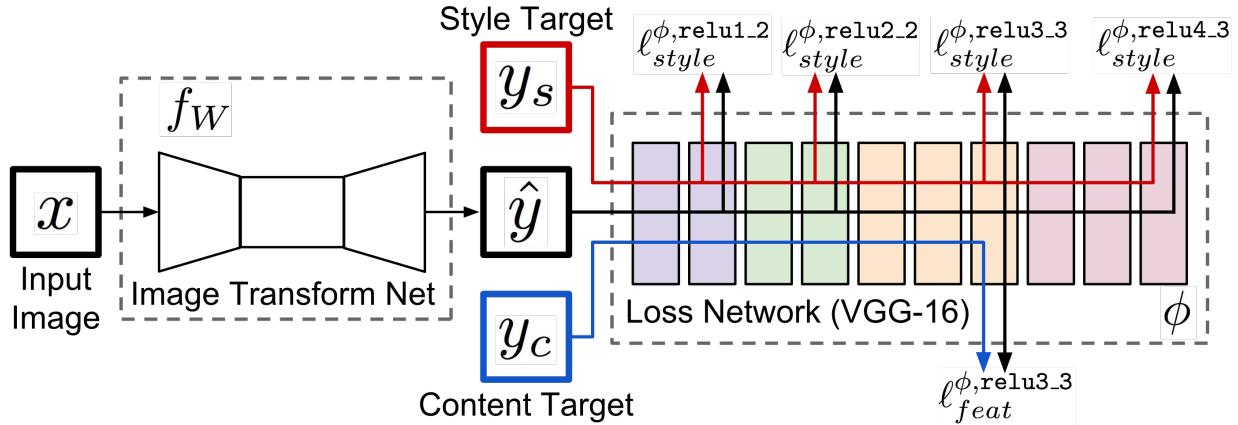


Figure 5.7: Training set-up by Johnson *et al.* [7]

basic style transfer even with random initialized networks, but results vary widely depending on parameters of the random initialization. Ultimately, it is possible to obtain some style transfer with this technique but the pre-trained weights seem to play an important role in stabilizing the reconstruction process.

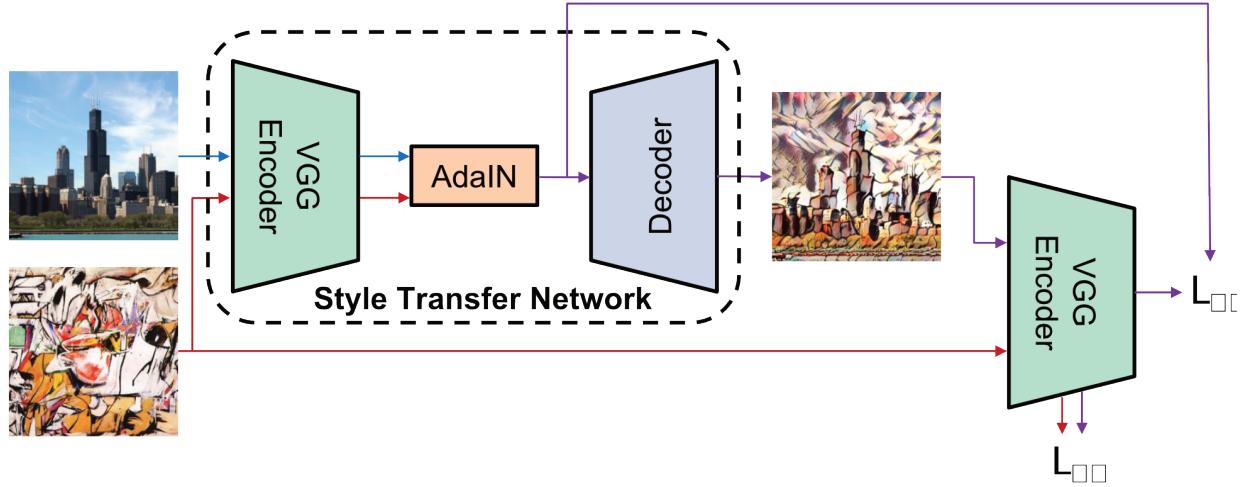
## State of the Art

### Real-Time Style Transfer

Following Gatys *et al.* seminal work, others have followed suit in trying to stylize images with neural networks. **Johnson** were the first to use the same losses but train a feed-forward architecture [7]. They were able to significantly speed up the stylization process like this as stylization is performed in a single feed-forward pass instead of a lengthy gradient descent optimization. Ultimately, this enabled them to generate stylized images in real-time from a given content image, using a deep residual convolutional neural network.

### Arbitrary & Universal Style Transfer

Huang *et al.* used a different feed-forward approach for arbitrary style. They first encode an arbitrary style image  $s$  as well as a content image  $c$  using a pre-trained VGG network.



**Figure 5.8:** Training set-up by Huang *et al.* [8]

This allows them to obtain the activations at a very deep layer of the network  $F^l(\mathbf{s})$  and  $F^l(\mathbf{c})$ . Then they compute the second order statistics for both  $\mu_F^l(\mathbf{s}), \sigma_F^l(\mathbf{s})$  and  $\mu_F^l(\mathbf{c}), \sigma_F^l(\mathbf{c})$ . Using adaptive instance normalization (AdaIN), they rescale the content activations such that they match the statistics of the style activations.

$$F'^l = \sigma_F^l(\mathbf{s}) \frac{F^l(\mathbf{c} - \mu_F^l(\mathbf{c}))}{\sigma_F^l(\mathbf{c})} + \mu_F^l(\mathbf{s}) \quad (5.14)$$

Finally, they train a decoder that minimizes the style and content loss, as they have been proposed by Gatys *et al.* They achieve comparable results to other style transfer approaches at a similar speed to Johnson *et al.* while allowing for any target style even though only training once.

A similar approach by Li *et al.* relies on matching the covariance and the mean of content and style activations. They do that through what they call 'whitening and coloring transform' [9]. First they whiten  $F^l(\mathbf{c})$  into  $\hat{F}^l$  such that  $\hat{F}^l \hat{F}^{lT} = I$

$$\hat{F}^l = E_c D_c^{-\frac{1}{2}} E_c^T (F^l(\mathbf{c}) - \mu_c) \quad (5.15)$$

Where  $D_c$  is the diagonal matrix of eigenvalues of the covariance matrix and  $E_c$  is the respective orthogonal matrix of

this is actually just PCA decomposition to a certain degree

eigenvectors. such that

$$F^l(\mathbf{c})F^l(\mathbf{c})^T = E_c D_c E_c^T \quad (5.16)$$

Then coloring is performed by rescaling the whitened representation  $\hat{F}^l$  into  $\tilde{F}^l$

$$\tilde{F}^l = E_s D_s^{\frac{1}{2}} E_s^T \hat{F}^l + \mu_s \quad (5.17)$$

$$F^l(\mathbf{s})F^l(\mathbf{s})^T = E_s D_s E_s^T \quad (5.18)$$

The resulting  $\tilde{F}^l$  is then decoded with a pre-trained decoder to render the final stylized result. Li *et al.* use pre-train the decoder solely on natural images and perceptual loss and reconstruction loss as objective. Additionally, they introduce a pipeline that performs style transfer on multiple scales sequentially. They achieve good results in real-time with just training the decoder once.

### Bidirectional Style Transfer

Zhu *et al.* went a different way on style transfer and rely on a generative adversarial objective to identify style in an image. Specifically, they transform images between any two domains  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ , not just photos and artworks. For this, they use two discriminators ( $D_X$  and  $D_Y$ , one for each domain) as well as two transformation networks ( $G : X \rightarrow Y$  and  $H : Y \rightarrow X$ ). Each translation network then transforms a given sample from one domain into the other and the discriminator assesses the result.

$$\mathcal{L}_{\text{adv}} = \log D_Y(y) + \log(1 - D_Y(G(x))) \quad (5.19)$$

Additionally, the transformed image is then transformed again and compared to the original input, in what they call **cycle loss**.

$$\mathcal{L}_{\text{cycle}} = \|F(G(x)) - x\|_2^2 \quad (5.20)$$

In the end, Zhu *et al.* are able to stylize and de-stylize images with their networks  $G$  and  $H$ . The main novelty here though is the good stylization quality they achieve without any of the previously introduced style losses. They have also shown one way, in which GANs are also capable of performing style transfer reasonably well. Other notable efforts were .

list GAN-based style transfer efforts

### Adversarial Style Transfer

Building on these GAN-based approaches, Sanakoyeu *et al.* improved the quality of adversarial style transfer and extended it to abstract styles as well. They argue that ImageNet-based approaches inherently favor photorealistic styles through the data set that ImageNet has been trained on [11]. Furthermore, approaches like cycleGAN suffer a similar fate as the back-transformation with cycle consistency opposes loss of detail in more abstract styles.

In order to retain content and global structure of an image, they introduce a fixed-point loss, which requires the stylized image to stay as-is when being re-stylized.

$$\mathcal{L}_{\text{content}} = \|E(G(E(x))) - E(x)\|_2^2 \quad (5.21)$$

To minimize this loss, the encoder must understand original content and stylized content. They also implement a transformed reconstruction loss for better visual quality of the stylized image

$$\mathcal{L}_{\text{transformed}} = \|T(x) - T(G(E(x)))\|_2^2 \quad (5.22)$$

The results show good visual quality, especially concerning the details and loss of details for abstract styles. Also this approach focusses on stylizing not only for a single image but the style of an artist in general.

**dima** take this further and focus on stylizing different content specifically. This means, a person is differently depicted than a tree, considering the level of detail, colorscheme etc.. , which holds with real-world experience. They achieve stylization that resembles this behavior by using the same fixed-point loss that Sanakoyeu *et al.* used but combine it with a second update step. In this second update step they require similar scenes to be placed closely in feature space and dissimilar scene to lie further apart. They add a transformation block between encoder and decoder shape the feature space accordingly.

### Others

There exist many other approaches that are capable of transferring style – often as a byproduct of texture transfer or image-to-image translation. Some focus very heavily on stylization of portraits using self-attention modules [**ugatit**]. Others choose an approach similar to cycleGAN but add a shared encoding space for content and separate attribute spaces where style is encoded [**unit**, **munit**, **drit**, **drit++**]. With the latter ones mainly focussing on separating shape and appearance of images and recombining them arbitrarily. One such example is taking the posture of a person in one image and combining it with the clothes and appearance of a person in another image.

## 5.3 Painterly Rendering

Painterly rendering is a field of computer vision that understand stylization of images as giving the impression that a certain medium was used. Most often this would be the looks of pencil drawings or paintings as their looks are very distinctive, especially when compared to regular photographs. Thus, painterly renderings rely on a brushstroke or a pencil line as its smallest unit from which the data is generated. Style transfer, in contrast, relies on pixels as its smallest unit.

The use of these larger units automatically comes with a level of abstraction which painterly rendering approaches often reinforce through adjusting the coarseness of said units. One could see this as a hard regularization. Style transfer (see Section 5.2), on the other hand, aims to learn this abstraction. Still, an impressionistic style in painterly rendering could be achieved by limiting the length and width of brush strokes in painterly renderings.

This thesis' approach follows similar regularization to painterly rendering and thus different key approaches shall be presented.

The earliest approaches to painterly rendering were stroke-based renderings, which artificially generate single brush-strokes. The challenge then becomes twofold 1) improve the quality of these strokes and 2) improve the algorithm which aligns them. Both problems will again come up in this thesis.

Early filter-based rendering approaches then appeared with superior computational efficiency and ultimately led to style transfer and texture transfer. This has already been explained in the previous section but shall be mentioned here to place it historically and contextually.

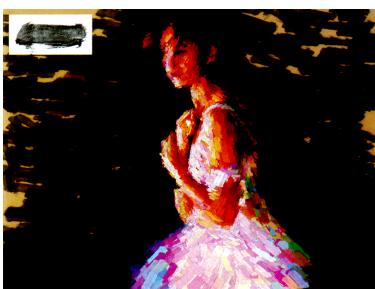
Lately, drawing networks revived stroke-based rendering with modern machine learning techniques. The focus of this field is really to remove any hand-tuning of parameters that stroke-based rendering required before and to a certain degree imitate the way humans would draw.

## Stroke-Based Rendering

Haeberli introduced stroke-based rendering with his seminal work 'Paint by numbers: abstract image representations' [12]. In this work he presents two methods which would allow for reconstructing images in an abstracted representation. His first approach is interactive and requires a user to click a certain points in an image to place a brush stroke there.



**Figure 5.9:** Interactively painted images using Haeberli's method with a hand-selected orientation (Figure ??) and a gradient-driven orientation(Figure ??).



**Figure 5.10:** Rendered image with a brush stroke texture

Then his software would automatically select a brush stroke color and size. He proposed several ideas on how to align the brush strokes on the canvas. Either users could do this on their own, or the orientation would be perpendicular to either the global gradient (uniform alignment) or local gradient (non-uniform alignment) of the image.

Haeberli even introduced the use of a scanned brush stroke texture in his algorithm.

Additionally to his interaction-based approach, Haeberli also introduced a relaxation based approach [12]. In this approach a given set of 100 brush strokes is iteratively perturbed. If the perturbation minimizes the energy function (L2-distance), the perturbation is kept. If not, another perturbation is applied.

This is described by Hertzmann as trial-and-error algorithm and very similar to genetic algorithms which shall be explained later [13].

Based on Haeberli's work, other authors automated the process of stroke positioning. At the same they improved various aspect of brush strokes which were well categorized in a review by Hegde *et al.* The relevant categories are position, path and orientation, length and width, ordering, and color [PRreview].

Litwinowicz proposed a straight-forward way of placing brush strokes evenly spaced over the entire image. Parameters are then inferred similar to Haeberli's approach. To give a more random feel to these brush strokes, the obtained parameters are randomly perturbed according to preset parameters. Also, a weakness of orienting brush strokes perpendicular to the gradient is dealt with. For large uniform areas with little to no gradient, the orientation could become arbitrary. Litwinowicz proposes to refine the gradient field by interpolating between the boundaries of large uniformly colored areas. Also, Litwinowicz introduced temporal coherence to brush strokes, which meant that brush strokes move with the optical flow between frame in a video.

**hertzmann** reformulated the problem as an energy minimization problem in two publications [**hertzmanreview**, **Hertzmann**].

$$E(R) = E_{\text{recon}}(R) + E_{\text{area}}(R) + E_{\text{nstr}}(R) + E_{\text{cov}}(R) \quad (5.23)$$

$$E_{\text{recon}}(R) = \sum_{x \in W, y \in H} w_{\text{recon}, x, y} \|I_{x, y}(R) - I_{x, y}\|_2^2 \quad (5.24)$$

$$(5.25)$$

where  $R$  is a brush stroke representation,  $I$  is the target image, and  $I(R)$  is the rendered representation.  $x$  and  $y$  are pixel positions. By adjusting the different weights, properties of the rendering can be altered.  $w_{\text{recon}}$  can vary spatially and dictate how well the reconstruction must fit the original image in certain areas.

Additionally, he added long strokes as B-splines with arbitrary control points. In contrast, Haeberli and Litwinowicz argued that short straight brush strokes would aid the perception of impressionistic style. Furthermore, **Hertzmann** added advanced rendering for brush strokes with synthetic textures in his work [**Hertzmann**].

**Hertzmann** combined all these aspects in his approach with advanced relaxation methods similar to Haeberli. Based on trial-and-error search, **Hertzmann** samples a local region along the many dimensions that represent a single brush stroke. The best set of parameters that minimizes the energy function  $E$  is then picked as new parameters.

In order to achieve better visual quality **Hertzmann** also employs a coarse-to-fine multi-layer rendering approach. Hereby, he blurred the image in the early iterations of his method and fixed the brush size at a large value. Blurring of the image would then be gradually reduced along with the brush size. The final implementation is further optimized to accelerate the relaxation algorithm and allow for more brush strokes than Haeberli's approach.



**Figure 5.11:** Image approximated by relaxation.

Ultimately, **Hertzmann** achieves respectable results with his approach and many ideas of this thesis can be found in his works as well.

### Stroke Rendering

Parallel to the advancements in arranging brush strokes in stroke-based rendering, others improved the rendering quality for many different styles. Most notably for this work, Baxter *et al.* were the first to implement a full 3D simulation of a brush and the process of placing paint on a canvas. They were even able to simulate mixing of color and different levels of dryness this way. Ultimately, they were able to showcase a drawing by a real artist based on the simulations in real-time [16].

**wetbrush** took this even further by including more accurate simulations of fluid dynamics and even single bristles in a brush stroke. They were able to generate a 3D model of paint on canvas from this which they could make subject to different lighting [wetbrush]. **adobe** tried to imitate this with the goal to achieve real-time simulation of the methods presented by **wetbrush**. Indeed, they were successful up to a point, where deviations to [wetbrush] became visible for too many stacked layers of paint.

### Drawing Networks

Drawing networks belong to the realm of painterly rendering approaches, yet are rarely compared with their predecessors as their architectures differ significantly. Drawing network often rely on recurrent architectures (see Section ??) or iterative approaches where the current state of the drawing is taken into account along the target-image.

Some of the first approaches used the mentioned recurrent neural networks to imitate predefined stroke sequences [17, 18]. Such approaches require paired data where the sequence

for a sketch or handwriting is already known. Obviously, this data is hard to acquire such that SPIRAL was the first publication to gain more attention in this field.

'Synthesizing Programs for Images using Reinforced Adversarial Learning' (SPIRAL) by Ganin *et al.* was the first approach to learn drawing sequences without such paired data [19]. They used a graphics engine and deep reinforcement learning, to train their network. Based on the current state of the canvas, SPIRAL predicts an action in form of brush stroke parameters. This process is iteratively applied to the canvas  $N$  times. Only at the end the result is evaluated in an adversarial fashion against the target image. Notably, the renderer is seen as a black-box and SPIRAL is agnostic to the given interface. Results show that this approach works well for simple stroke-based data such as handwriting data sets. Tests on images, such as the CelebA faces data set, show limitations as SPIRAL could only be trained for up to 20 strokes.

Other approaches followed suit, after Ganin *et al.*'s seminal work. They would often combine the idea behind SPIRAL with **worldmodel**'s idea to encode an environment through an autoencoder into latent space. A model is then trained on this fewer-dimensional latent space representation instead of the real-world data. **worldmodel** call their idea a 'world model'.

Zheng *et al.* employ differentiable rendering combined with two encoders to predict a single brush stroke in a feed-forward manner [20]. First they train a position encoder (which actually is a decoder) that transforms an input position to a coordinate in a 64x64 matrix. Then, they train the renderer with  $n$  control points decoded by the position encoder combined with additional information about brush size and pressure. The dataset is randomly generated with a fixed number of control points from a custom render engine. In the end, the agent is trained, which encodes the current canvas as well as the target image to predict a set of parameters for a single brush stroke.

THeir method works notably better for predicting single brush strokes on empty canvas, that dollwing brush stroks for more complex shapes. Zheng *et al.* also showed tests on images but got very poor results.

Nakano chose a different approach to differentiable rendering. He first trained a generator to decode action space input to brush stroke images in adversarial fashion. Then he used an LSTM-based agent to predict a sequence of brush stroke parameters which are rendered by the generator. As the whole pipeline is differentiable, Nakano was able to avoid the use of reinforcement learning. Results show better image reconstruction than SPIRAL when using the generator, but slightly worse results when combining the obtained action with the original renderer. Nakano also presents the capability of reconstructing content similar to Gatys, Ecker, and Bethge by minimizing the perceptual loss directly in action space with his network.

Long short-term memory (LSTM) is a basic building block in recurrent neural networks which allows to store information for many steps. It specifically learns which inputs to look at, when to forget something and which outputs to generate from the current state of the network.

Huang *et al.* also employ a pre-trained decoder and use more advanced reinforcement learning techniques [22]. This allows allows their network to predict strokes for an arbitrary number of steps, where SPIRAL could only predict up to 16 strokes. Their results show significantly more detail, especially as the number of strokes increases.

All previous approaches work only on very limited resolutions. This is why Jia *et al.* extend reinforcement-based approaches to higher resolutions by using more advanced reinforcemnt learning techniques [23, 24]. They do not use a differentiable renderer but a sliding-window approach which allows them to generalize to images of larger scales. Notably, they train each network specifically for an image, as the training does not generalize well for different images. They argue that this can be seen as a form of style transfer but such reconstructions mainly show a simple color shift.

Jia *et al.*'s approach is another approach close to this thesis, considering the goal of their approach. They are able to generate stroke-based renderings of high-resolution photographs

as well as images of painting just shy of 1 megapixel. Each network is specifically trained for a single target image, which takes at least 1h to accomplish [23].

## Genetic Algorithms

Genetic algorithms are typically not closely associated with painterly rendering or drawing networks, even though they represent just a different approach to algorithms for this problem.

Genetic algorithms already perform a similar task in order to approximate images by other geometric shapes or even smaller photos (also known as the popular photo mosaic effect). Starting with a random set of circles that are parameterized by their position, radius, and color, it then chooses the most successful samples and resamples in a region around these. This process is repeated until a certain level of convergence is reached.

As well as this does work, it is very much computationally expensive as most samples will not fit the image, thus searching for the small set of fitting shapes requires to evaluate all the wrong shapes as well. Considering artworks, brushstrokes have many more degrees of freedom, and artworks usually consist of upwards of a few thousand brushstrokes. Consequently, it would be considerably more challenging to apply to this problem until computational resources have become a few magnitudes more powerful.

**Figure 5.12:** Photo mosaic of Starry Night using only images by the Hubble Space Telescope. [http://www.astro.uvic.ca/~alexhp/new/figures/starrynight\\_HST.001.jpg](http://www.astro.uvic.ca/~alexhp/new/figures/starrynight_HST.001.jpg)



## 5.4 Conclusion

Works like these show that there is interest in improving the missing details in style transfer. Ultimately, it would be desirable that the fields of style transfer and painterly rendering converge further. Style transfer brings impressive perception of style. Painterly rendering brings realistic composition and/or low-level details.

This thesis aims to improve reconstruction quality for painterly renderings without reinforcement learning techniques while at the same time generating an interpretable brush stroke representation.

## **CONTRIBUTION AND EXPERIMENTS**



# 6

## Approach

### 6.1 Motivation

The basic approach of this work consists of two steps:

1. A differentiable renderer which can generate images of brushstrokes from a parameter representation.
2. An optimization procedure that iteratively approximates an image through brushstrokes representations.

Having two separate steps can be motivated by comparing the optimization procedure to the actual process of painting an image. An artist will most likely not pick single color particles and then place them on canvas. Instead, an artist uses a brush or other utilities (see Pollock or others) to place more paint with a single action.

Doing so – of course – limits the control over each drop of paint but maintains enough control to still create very delicate details in paintings. This trade-off depends on the brush's size, such that an artist must choose the brush size depending on the content.

An example would be the painting of a uniformly colored sky. Using a large brush size, the artist can cover a lot of canvas in relatively little time as well as keep the color well distributed over the canvas because the brush spreads the color more or less evenly within the brushstroke. On the other hand, if one were to draw a sky with the smallest brush available, not only would it take forever to paint, it would also be hard to keep the paint evenly distributed over multiple strokes.

Now, translating this onto the given problem of recreating/approximating an image through brushstrokes, it would

6.1 Motivation . . . . .	67
6.2 Neural Renderer . . . . .	69
Data Set . . . . .	70
Architecture . . . . .	78
Training . . . . .	79
Results . . . . .	80
6.3 Stroke Approximation . . . . .	80
Data Set . . . . .	80
Optimization Algorithm	81
Optimization Procedure	87
Placing & Blending . . .	94
Style Transfer . . . . .	104
Optimization Details .	105
Results . . . . .	106

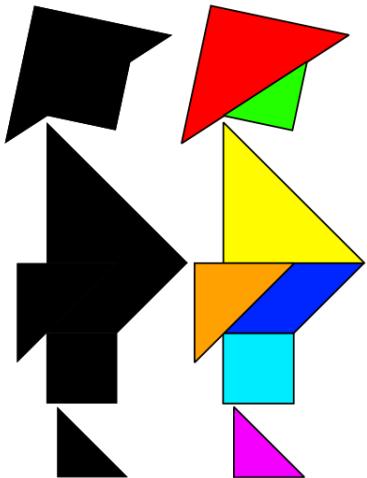


**Figure 6.1:** A typical set of brushes and spatulas used for oil paintings.

reformulate this

mean to limit the process to only use what we would describe as brushstrokes.

A comparable example is the game of Tangram.



**Figure 6.2:** An Example of Tangram.

Tangram is a Chinese puzzle game that has the objective of replicating a given silhouette only with a set of 7 unique shapes. The shapes may not overlap or be cut or anything. Quite similarly, the objective of an optimizer is to replicate an image by only using brushstrokes.

Genetic algorithms already perform a similar task in order to approximate images by other geometric shapes or even smaller photos (also known as the popular photo mosaic effect).

**Genetic algorithms** follow a random sampling approach that 'evolves' as genomes do. Starting with a random set of circles that are parameterized by their position, radius, and color, it then chooses the most successful samples and resamples in a region around these. This process is repeated until a certain level of convergence is reached.

As well as this does work, it is very much computationally expensive as most samples will not fit the image, thus searching for the small set of fitting shapes requires to evaluate all the wrong shapes as well. Considering artworks, brushstrokes have many more degrees of freedom, and artworks usually consist of upwards of a few thousand brushstrokes. Consequently, it would be considerably more challenging to apply to this problem until computational resources have become a few magnitudes more powerful.



**Figure 6.3:** Starry Night approximated by a genetic algorithm using only circles. <https://effyfan.com/2018/03/02/w6-van-gogh-flowfield/>

This premise can be overcome, though, by using a differentiable renderer. A **differentiable renderer** is capable of creating/rendering shapes in the pixel domain. In contrast to conventional renderers, it does so by solely using differentiable operations. Thus, the previously described task becomes feasible, as random sampling can be replaced by gradient descent. Ordinary renderers usually do not rely

on differentiable operations, as faster operations with more straightforward logic render images well enough already.

Nonetheless, it is theoretically possible to create a differentiable renderer [something].

When talking about rendering brushstrokes, it would be even harder to think of a differentiable pipeline to draw brushstrokes of reasonable quality from a set of parameters.

Neural networks turn this problem around. As neural networks are inherently differentiable, the question becomes rather how to make an existing neural pipeline render images from parameters. Previous works have shown that neural networks are capable of conditionally generating high-resolution and high-quality images. Conditioning the generator on brushstroke parameters as well as some noise should then output an image of the respective brushstroke with some variability to it.

The basic idea was proposed by [japaneseneuralrenderer](#) [japaneseneuralrenderer] as it facilitates training of reinforcement learning based networks.

Inspired by this, the approach becomes more apparent. First, a neural network is trained as a differentiable renderer. Then the same renderer is used by a gradient descent-based optimization procedure to approximate an artwork as a set of renderer input parameters.

Both steps require some tricks to avoid pitfalls like computational limitations, which are outlined in the following two sections.

## 6.2 Neural Renderer

The neural renderer in [japaneseneuralrenderer](#)'s [japaneseneuralrenderer] work improves an architecture based on SPIRAL [19]. Learning2Paint's efforts used a differentiable renderer to facilitate deep reinforcement learning for drawing images [Learning2Paint].

Move this into the related work section and shorten it further

As mentioned, the renderer is required to be differentiable and ideally require as few resources as possible.

## Data Set

Unfortunately, there is no data set available for this task, which means that a data set must be created explicitly for this approach.

There are several sources for virtual and real brushstrokes, to choose from. These resources will be evaluated in the following part. The main focus lies on three qualities for each data source:

1. Suitable data format: As brushstrokes usually overlap in a painting, the data should already provide information about opacity or transparency, favorably in the common RGBA format, which has a fourth alpha channel to hold opacity information.
2. Data set size and variability: Only with enough different data available, it will be possible to train a renderer reliably.
3. Image quality. Brushstrokes should be as close to real-world brushstrokes as possible to ensure high-quality renderings later on.

## Brushstroke Images

There are multiple sets of hand-drawn brushstrokes available online. Most notably, there is a set of various well-classified colors and brush styles created by 'zolee' on the platform [onlygfx.com](http://onlygfx.com). It consists of approximately 1000 brushstrokes that mostly follow rather straight horizontal paths. These brushstrokes are mostly grouped by color and painting technique (oil, acrylic, watercolor...). All images are in the PNG format with the background made transparent in a post-editing step.

reference this

This data set has the advantage that it consists of real-world brushstrokes that were painted under presumably reproducible conditions. On the other side, brushstrokes are of mostly the same width throughout the data set and also do not come with information which path the brush took or any other non-visual information. Also, the data is very sparse. Many color shades are not represented, which means that the generator would have to interpolate them or simply would not be capable of rendering any brushstrokes in this color.

It seems that this data set would be suitable to replicate single real-world brushstrokes as images. However, limitations to the data make it unlikely that a generator could learn a coherent representation from this.

## Painting Libraries

The mentioned work of SPIRAL [19] relies on synthetic data instead of real-world images. It used the painting library 'libmypaint' [[libmypaint](#)] to generate brushstrokes from parameters in real-time during training.

The apparent advantage of this and other painting libraries is the fact that one can fully control the output through parameters. As the whole space of input parameters for the renderer can be covered, it is much easier to avoid pitfalls like they were described in Section ??.

Still, this data set falls short regarding the authenticity of rendered strokes. Especially the inner area of the stroke shows a uniform color, which is far from what real brushstrokes would look like.

This data set is better suited for our task than the images from which were described in the previous section are. However, the synthetic data will tend to make all rendering look a bit 'cartoonish' or flat, which could, in turn, limit convergence during the latter optimization process.

## Fluid Simulation

Fluid Paint is a project by David Li [[fluidpaint](#)] that uses simple fluid dynamics to give artificial brushstrokes a more plastic look. It has been implemented in JavaScript and OpenGL.

There is a C++-version in the git-repository of SPIRAL along with Python bindings created by Yaroslav Ganin.

Using these Python bindings, it is possible to generate brushstrokes locally outside of a web browser environment.

The quality and controllability of FluidPaint fall right in the middle of the two previously mentioned data sets. The generated brushstrokes look distinctively better than those generated with libmypaint, but still lack the quality of the real-world images. Concerning controllability, FluidPaint allows controlling the path of the paintbrush's handle rather than the path of the brushstroke itself. As a result, there occur some offsets to a given path as opposed to libmypaint.

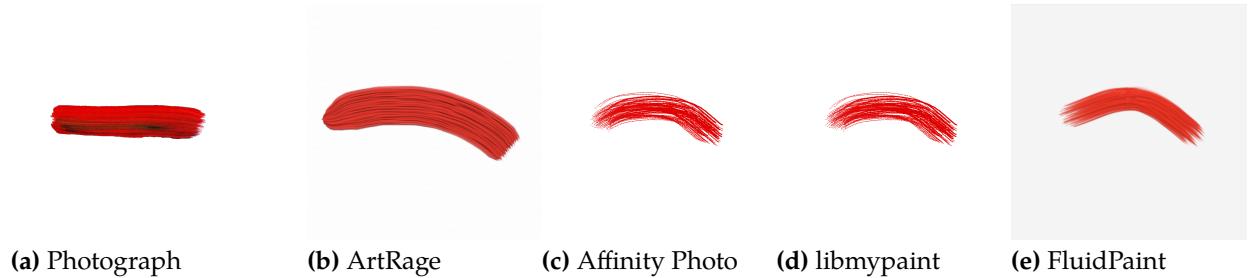
As a reasonable compromise between controllability, variability, and render quality, FluidPaint emerges as a sensible choice. Although real-time data generation is not possible with this library, data generation can be parallelized. Such a set-up still allows for the creation of large data sets in a reasonable time frame.

Even though this data set still has some weaknesses, it is probably the best choice for training a differentiable renderer because of the noted reasons.

these two weird software stuff things

Other 'honorable mentions' are painting programs such as , which allow for even more authentic brushstrokes but lack any well-documented interface in order to generate a vast number of brushstrokes.

talk about RGBA advantages



**Figure 6.4:** Comparison of similar brushstrokes in each data set

### Brushstroke Formalism

With the means of data set production seized , what is left is to formulate the parameters that define the brushstrokes. These parameters must quantify the following three properties of brushstrokes:

cut this joke

- ▶ color
- ▶ thickness
- ▶ path

The easiest of these three properties is quantifying the color. Naturally, computer vision relies on the RGB format, which defines color as a set of three 8-Bit integer values between 0 and 255. As for path and thickness, these two properties depend on the given coordinate system. FluidPaint represents the canvas as a 2D plane in the  $[0, 1]$  range. Thus, it makes sense to follow the same representation.

Thickness can be any value in  $[0, 1]$  for each brushstroke, where 0 is an infinitely small brushstroke, and 1 is a brushstroke as wide as the canvas. As both these edge cases do not make sense in this application, the range is constrained to  $[.03, .2]$ , which includes only brushstrokes that are visible and also do not cover the whole canvas.

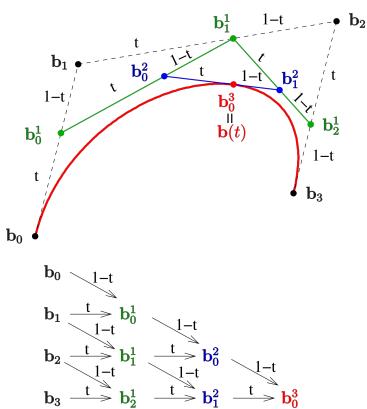
Quantifying the path now is a little more tricky. The fluid dynamics simulation that FluidPaint uses relies on internal time steps at which the equations are evaluated and subsequently rendered. At the same time, each step allows only a linear motion of the brush handle between positions  $a$  and  $b$ .

this is actually different as the steps depend on the length of the path.  
Look this up

Subsequently, any curved paths must be split into linear/straight segments that together should resemble a curved line. As more steps mean longer simulation times and fewer steps mean edgy movement, a value of 20 time-steps per stroke emerged as a good compromise. The same number of steps has been used in SPIRAL's implementation.

Another problem becomes how to express a curved path in numbers. The most straightforward representation would be a sequence of points that make up any curved path. Such an approach allows for the highest versatility, but at the same time introduces a noticeable amount of parameters as each point consists of 2 coordinates totaling to 40 values. These values are also not independent of one another but should follow a reasonable path as otherwise, the resulting brush-stroke would look somewhat like a random walk than an actual brushstroke. Since works such as SPIRAL or 'Learning to Paint' face a quite similar obstacle, their solution should be applicable in this case as well. Both used so-called 'Bezier curves' which parametrize curved paths by a limited set of numbers.

Bezier curves parametrize curved paths not as linear segments, but as analytical paths which can be used in computer models. They can be of different orders, which allows them to follow more complicated paths. In this case, the simplest form – the first order Bezier curve – is already sufficient. It defines a curve through its start- and endpoint, as well as a control point. The curve is then defined as the path of a point over the time interval  $[0, T]$ . First, one connects the start- and endpoints with the control point to get two lines in return. On these lines, lie two points that move linearly along their lines within a virtual time interval  $[0, T]$ . Then these two points are connected in the same way as before by a third line. Again, this line has a point moving along its path throughout  $[0, T]$ . As the first two points that define the third line will move, the line's orientation will change as well, thus translating the linear movement of the point into a complex curved path. The path that this point then takes in  $[0, T]$  defines the



**Figure 6.5:** Sample of a 3rd degree Bezier curve, using the De-Casteljau-algorithm, <https://de.wikipedia.org/wiki/B%C3%A9zierkurve#/media/Datei:Bezier-cast-3.svg>

bezier curves margin note

Bezier curve. A first-order Bezier curve will only bend into one direction or follow a straight path. For higher orders, the displayed process can be applied iteratively and allows for more complex curves. As brushstrokes usually follow a quite simple path and fewer parameters are preferred, Bezier curves of first-order are suitable as parametrization.

Ultimately, this gives ten values that are sufficient to parametrize brushstrokes with certain constraints:

- ▶ Three 2D coordinates that define the Bezier curves (6 values).
- ▶ One thickness parameter.
- ▶ Three values in RGB space.

## Data Constraints

Given the parameters listed in section 6.2, the data still needs further constraints to facilitate the generator's training even further.

Section 6.2 already hinted at the impracticality of online data generation. A rough estimation by timing the rendering of 100.000 FluidPaint brushstrokes reveals that a dedicated CPU server is capable of generating 300 strokes per second. A neural network with batch size 32 is limited to  $\approx 10$  iterations per second under these circumstances, which would mean a clear bottleneck. Thus, it seems advisable to generate data beforehand with enough samples to cover the data space sufficiently. It will allow for much faster access to data, as individual data samples are relatively small and can be stored in a binary data file such as HDF5.

experiment in appendix

Besides this constraint to the amount of data available, another set of constraints will be introduced to reduce the data space to 'valid' brushstrokes only. 'Valid' brushstrokes will be defined as brushstrokes that resemble real-world brushstrokes to a certain degree. This primarily concerns two relations within a brushstroke:

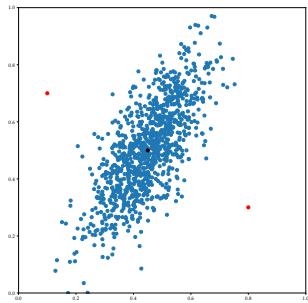
- ▶ Its width-to-length ratio.
- ▶ Its curvature.

The width-to-length ratio will be restricted to brushstrokes that are at least two times as long as they are wide.

$$\|\mathbf{s} - \mathbf{e}\| \stackrel{!}{\leq} 2 \times (\text{brush size}) \quad (6.1)$$

Due to the simulation background of FluidPaint shorter brushstrokes will show some artifacts due to the bristles' length in the simulation which depends on the width of the stroke. Another reason for this is the intended use-case, which will focus on van Gogh paintings. As van Gogh did not practice pointillism, most of his strokes have a length to them, which brings such a constraint in line with some characteristics of van Gogh's style.

The same argumentation applies to the curvature: Most brushstrokes (especially those by van Gogh) have a certain 'flow' or 'smoothness' to them, which can be described by using strokes with large curvature radii and without any corners in the strokes' path. Thus, the data set will also be restricted to strokes that follow these descriptions. In order to achieve this with random sampling in mind, a multivariate Gaussian distribution is placed between start ( $\mathbf{s}$ ) and end point ( $\mathbf{e}$ ). The two axes are rotated such that the short axis is in line with the vector  $\mathbf{a} = \mathbf{s} - \mathbf{e}$  while the other sits orthogonal. Then both axes are scaled with  $\|\mathbf{a}\|_2$  and also the handpicked values  $\frac{1}{200}$  and  $\frac{1}{25}$  for along  $\mathbf{a}$  and orthogonal to it, respectively. Figure ?? shows samples from this distribution for an exemplary brushstroke.



**Figure 6.6:** Exemplary scatter plot for given start and end point to visualize the covariance matrix

This distribution is intended to follow that of brushstrokes as they would appear in the real world. The majority of brushstrokes will be straight or just slightly bent due to the maximum of the PDF being at the center of  $s$  and  $e$ . Bent brushstrokes will mostly be symmetric as the long axis of the multivariate Gaussian is orthogonal to  $a$ . Still, there will be strokes that have their bend towards either end of the brushstroke as well as some strokes with a high curvature.

The area of interest, though, will be densely populated as intended.

$$p(\mathbf{c}|\mathbf{s}, \mathbf{e}) = \mathcal{N}(\mu, \Sigma) \quad (6.2)$$

$$\mu = \frac{\mathbf{s} + \mathbf{e}}{2} \quad (6.3)$$

$$\Sigma = \begin{pmatrix} a_x & 0 \\ 0 & a_y \end{pmatrix} \quad (6.4)$$

$$\mathbf{a} = \mathbf{s} - \mathbf{e} \quad (6.5)$$

The color of the brushstrokes is not constrained as the color distribution of the target data set is not known at this point.

why not van Gogh color distribution?

## Data Set Creation

The data set will be created with 100.000 samples that follow the constraints that were presented in section 6.2. As an underlying distribution, the uniform distribution is chosen as it allows a more evenly coverage of the data space.

First, a set of start and end points, as well as brush size, is drawn and checked against (6.1). If the constraint is not met, the set will be redrawn entirely. In case the constraint is satisfied, a checkpoint is sampled according to (6.2). If  $\mathbf{c}$  lies outside the render window, the checkpoint will be resampled. At last, an RGB set is sampled from a uniform distribution as a color.

The resulting tuple of start, end and control point, brush size, and RGB color is then added to the data set. Before rendering starts, the values of  $\mathbf{s}$ ,  $\mathbf{e}$  and  $\mathbf{c}$  are scaled with the handpicked factor of 0.7 to ensure the brushstrokes are rendered completely within the window and not cut by an edge of the render window. At last, the brushstrokes are rendered according to the data set and added as well.

The render canvas size was chosen to be 64x64 pixels for several reasons: First, even with such a small canvas size,

training for the renderer takes about one day. Secondly, the larger the render canvas size becomes, the deeper the renderer needs to be, which results in more computational overhead in the optimization routine as well as more layer through which the gradient has to be propagated. Lastly, as there will be upwards of a thousand brushstrokes in a single image, increasing the canvas size to 128x128 would require four times as much memory per rendered image. As 1000 brushstrokes would already account for  $1000 \times 64 \times 64 \times 4B \approx 16.4\text{GiB}$  a fourfold increase would be significant.

As the last step, the data set is renormalized to the range  $[-1, 1]$  for convenience and to facilitate training as well.

explain why zero-centered data is good

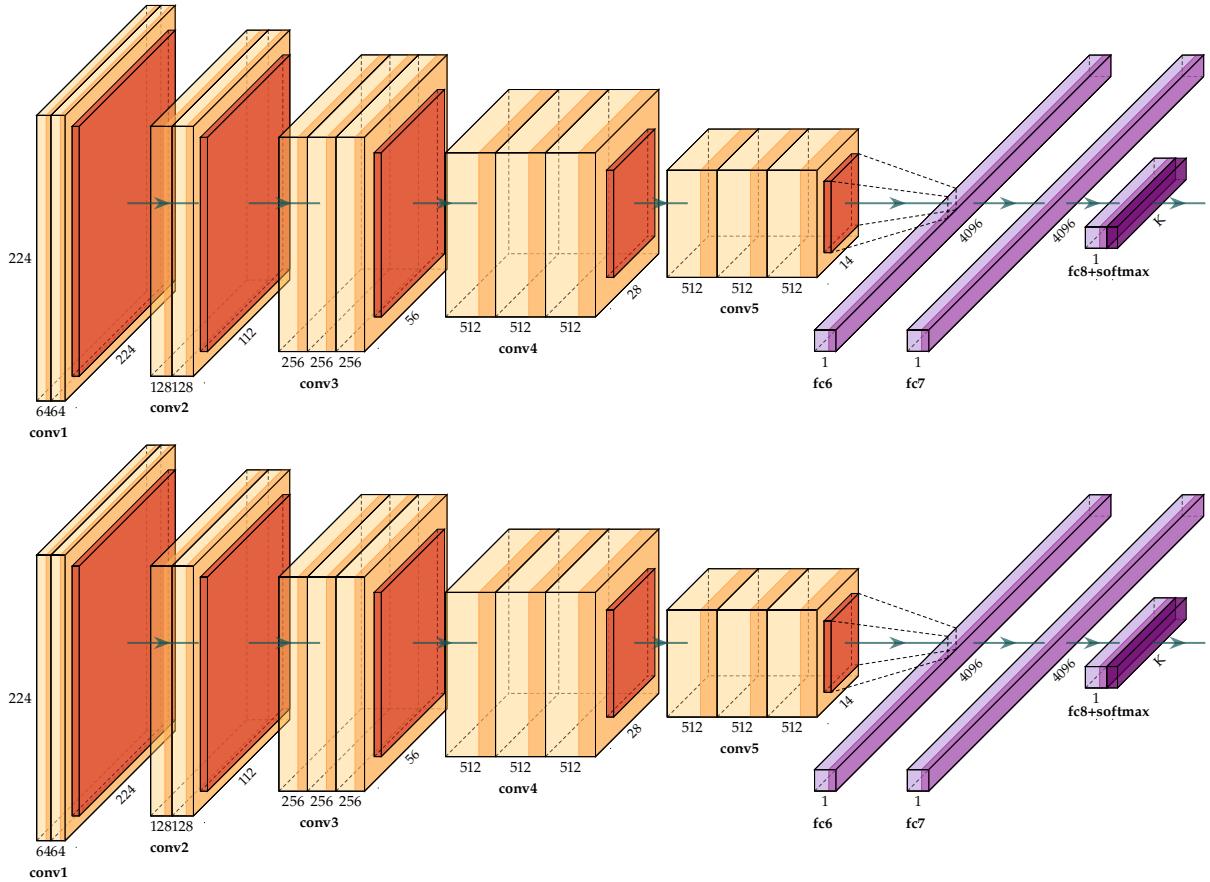
check the details

## Architecture

The architecture of the brushstroke generator follows that of an inverse VGG network. It is widely used and has shown in previous works that it should be capable of handling this task. The architecture consists of three dense layers at the beginning, followed by a two-times upsampling layer as well as three convolutional layers. The same pipeline with a two-times upsampling layer, and three convolutional layers is repeated until the target size is reached. After the last convolutional layer, a hyperbolic tangent function is applied to restrict the output to the  $[-1, 1]$  range. As part of the hyper-parameter search, different tweaks to the architecture have been tested:

- ▶ An additional noise input at every layer with a size equal to that of the existing signal.
- ▶ Additional information about the position in the pixel grid in every layer, so-called CoordConv [**coordconv**].
- ▶ Various combinations of activation and normalization functions.

The discriminator is designed after the same principles and resembles a VGG encoder network. First, three convolutional layers are applied, followed by a downsampling /pooling



**Figure 6.7:** Visualization of the generator and discriminator architectures

layer. This structure is repeated until a target resolution of 4x4 pixels is reached. Then a set of three dense layers is applied to give one final prediction per sample.

## Training

During training, the L2 distance and the FID score serve as evaluation metrics. The FID score becomes necessary as the visual comparison of the generated samples has proven difficult between different experimental runs. The L2 distance does not qualify as a sufficient metric for later training stages as the stochastic nature of the brushstrokes puts a lowerbound on the L2-distance.

A two-time-step update rule was implemented to stabilize training further.

write this down and  
look at the tricks that  
were used

## Results

pick some results and present whether they are any good

### 6.3 Stroke Approximation

#### Data Set

The data set will be presented in this section due to the data defining demands for later networks. Also, the data set for the optimization task should meet a few requirements.

In order to focus on the brushstrokes in an image, the data set should consist of relatively high-resolution images. At the same time, most brushstrokes should fit into a 64x64 window. Any larger strokes can not be approximated with a single rendering window and thus would falsely be approximated by multiple brushstrokes.

reference li and lamberti  
aaas they did the same

Thus, more information per image is required, than merely the resolution. The scale of the image, therefore, becomes inquired information in this task. Ideally, each image should be accompanied by its measurement as this allows for rescaling the image accordingly, given that one knows how large a typical brushstroke is.

At last, the painting's technique must be oil on canvas or similar techniques as this is what the renderer has been trained for.

All these requirements on the data are met by data directly available from the Van Gogh Museum in Amsterdam. Each high-resolution image is categorized by its technique as well as the period in which it was painted. This is accompanied by information on the dimensions of the image.

## Optimization Algorithm

The next stage of the approach takes the pre-trained neural renderer and uses it to approximate brushstrokes in images. In theory, the neural renderer should allow gathering meaningful gradients in parameter space even though the losses are calculated in image space, as explained in section 6.1. There, it was outlined already that the method of choice will be an optimization procedure that relies on the differentiability of the neural renderer.

Nevertheless, this is not the only possible approach to dissect images into sets of parameters. As explained in 5.3, there exist a variety of approaches for this kind of task. Some of these promise a fast inference of parameters from images, which makes the optimization approach in this work seem like a step back at first. The following section is dedicated to justifying the chosen approach by comparing the existing state-of-the-art approaches.

The targeted task is to approximate the representation of an image with roughly 1,000,000 pixels = 1MP by means of  $\approx 10.000$  brushstrokes (see section ??).

### Genetic Algorithms

Genetic algorithms are possibly the most straight-forward approach. As it was laid out in sections 6.1 and ??, genetic algorithms use random sampling and previous best results to approximate fitting solutions to the rendering problem, which was described in section ??.

put much of this into  
the related work section

Current state-of-the-art solutions are capable of finding a solution in  $\approx 1h$  when searching for an approximation of a 1MP image with simple geometric shapes like circles or triangles. This corresponds to roughly  $1 \times 10^9$  sampling steps. The time frame and the number of sampling steps depend on target shape density, target accuracy, sampling rate per

shape, and degrees of freedom (with the latter requiring more sampling).

Looking back at section ??, each brushstroke has 10 degrees of freedom with some inter-dependencies between them. Also it is known that the obtained neural renderer is capable of generating  $\approx 300 \frac{\text{images}}{\text{s}}$ . This means that rendering  $10^9$  sampled brushstrokes will take:

$$\frac{10^9 \text{samples}}{300 \frac{\text{images}}{\text{s}}} \approx 3.33 \times 10^6 \text{s} \approx 926 \text{h} \approx 38.6 \text{d}$$

Clearly, this is an absurd amount of time to spend **per image**. Furthermore, the sampling magnitude has not even been corrected for by the additional degrees of freedom that brushstrokes provide.

Ultimately, this means that genetic algorithms are not an option for this task, as they are too inefficient.

## Brushstroke Extraction

put much of this into  
the related work section

Next in the line are algorithm-based approaches that extract brushstrokes from an image using standard computer vision techniques.

Besides texton based image characterization [**textons**] and pure filter based approaches [**filters**] there have been approaches to extract brushstrokes or some of their characteristics [**brushstrokecharacteristics**] [**brushstrokeextraction**].

The latter would pose a valid option for the task at hand. Unfortunately, the best existing techniques fail to detect brushstrokes reliably over the whole image but only identify the most significant ones (see figure ??) / todo add example image. Other approaches are only able to extract a few characteristics, like the orientation of a brushstroke, which proves insufficient as well.

Other approaches will usually not get parameters but again pixels, which they analyze.

show some of the extracted brushstrokes

## Stroke Based Rendering

Stroke Based Rendering or Painterly Rendering also uses algorithms to approximate an image through brushstrokes with the original goal to achieve some stylization along the way. While early works relied on interactive approaches, later publications were then able to automate this process entirely. Judging from tests by the authors, such an approach would take between and hours per image. As figures ?? and ?? show, these approaches obtain similar results to genetic algorithms and tend to draw an image from a coarse scale to a more delicate scale over time instead of a locally coherent manner.

this number

this number

These approaches were not intended to be used to obtain brushstrokes from a painting but focus on stylizing images. This renders such an approach unfit for the goal of this thesis as well.

## Drawing Networks

Lastly, drawing networks are the newest iteration of approaches in this field. Beginning with , there have been approaches that make use of feed-forward or recurrent neural networks combined with either supervised training or deep reinforcement learning. The best results of these approaches are shown in figure ??.

put much of this into  
the related work section

which one was first?

Noticeably, all of these images have a maximum resolution of 256 pixels along their longest edge. There have not been any approaches yet, which can go significantly beyond this limitation. Notably, the high computational costs of training recurrent neural networks seem to be an obstacle when shooting for higher resolutions. Such resolutions can not be deemed sufficient when looking at individual brushstrokes, as outlined in 6.3.

Consequently, drawing networks must also be ruled out.

## Combined Approach

Even though there have been plenty of previous approaches to the task of extracting brushstrokes from images or similar tasks like rendering images through brushstrokes, none of these quite meet all the requirements of this task. Namely, high-resolution input images, brushstroke focussed rendering or retrieval, limited resources, and realistic depiction of brushstrokes.

Although brushstroke extraction is the topic of this chapter, stroke based rendering and drawing networks both seem capable of making brushstroke extraction more feasible with their inverse approach. Drawing networks introduced differentiable renderers as a tool to facilitate training. Stroke based renderers achieve a parametrization implicitly by focussing on replicating the entire image instead of extracting individual brushstrokes.

They can be combined into a unified approach that uses the differentiable renderer and the objective of recreating an entire image. The resulting approach does not use re-sampling like genetic algorithms but relies on gradient descent to converge to a solution significantly faster.

The limited capabilities of such a renderer would also guarantee that any approximation is composed of only valid brushstrokes instead of single pixels as it would be the case with normal generative models.

All in all, this approach is probably the best suited for approximating a target image only through brushstrokes as closely as possible.

The decision of whether to use an optimization-based approach is also linked to another question: Whether to sequentially or parallelly place brushstrokes.

As an example, most drawing networks and stroke-based renderers rely on a sequential approach. Intuitively, it makes sense to use a sequential approach as artists also place their

brushstrokes sequentially on the canvas. However, there is a significant difference in how existing computer vision approaches place strokes compared to artists. As neural networks search for ways to reduce the loss function as much as possible and as fast as possible, images tend to be made up of large canvas-filling brushstrokes at early stages, which in turn reduce the L2 loss.

This is contrary to how an artist usually works. Artists tend to fill the background based on content (*e.g.* sea or sky is often painted first) and will often leave some blank canvas to paint foreground objects later on. They also do not use giant brushes for this but normal-sized brushes with which they place many strokes. Thus, sequential approaches differ significantly from real-world paint processes.

In contrast, parallel approaches predict the whole painting as one set of parameters. They are not as widespread due to the computational pitfalls that come with predicting brushstrokes for a whole painting at once compared to predicting only a few at a time. Also, parallel approaches still need to predict in which order brushstrokes should be placed on the canvas (this will be explained in section 6.3) such that brush strokes can overlap. Nonetheless, the advantage to sequential approaches is the fact, that foreground and background strokes can influence one another.

An example would be a stroke in the foreground, changing its path and revealing the canvas beneath. Then another stroke in the background should cover this up if the color matches. For sequential approaches, the background stroke can not be changed afterward, and it is very complicated to formulate a loss that propagates this information. As artists already plan their future brushstrokes when painting the background, it would require a drawing network to plan far ahead. At the same time, an optimization-based approach will not need to do that.

Another argument for parallel optimization is the focus on actually visible brushstrokes. As artists cover up previous

strokes, again and again, it becomes a shot in the dark to guess what these first brush strokes might have looked like. Subsequently, any such guess is ill-posed and only introduces noise to the problem. A parallel approach does not care for how the background could have been drawn but only focus on what is visible in the given image.

### Pitfalls of Feed-Forward Approaches

In the development of this thesis, there were experiments targeting a feed-forward approach before ultimately deciding on an optimization-based approach that is now presented. As it would exceed the scope of this thesis, the arising problems and cause for discarding this approach shall be discussed.

First, the computational burden of a feed-forward approach is very high. Existing feed-forward drawing networks compromise image resolution to realize their implementation. As it was possible to implement this for small scale data like the cMNIST data set, such an approach seems feasible at first. However, since compromising image resolution is not a viable option (see Section ??, one must find ways around this problem. Fully convolutional architectures are a popular solution, which allows training on small scale data and inference on large scale images. Unfortunately, the problem of predicting many spatially brushstrokes with complex parameters has proven too difficult for fully convolutional approaches.

This goes hand in hand with another problem that occurred: the placement of brushstrokes on the canvas. As artists are not bound to the same pixel grid as computers typically are, they can place brushstrokes freely on the canvas. More so, they can pack brushstrokes densely in one area while distributing them broadly in another. Classic CNNs are not able to allow for similar behavior as they always require a grid layout. Experiments with either displaceable grid cells or stacked grids have proven did work on a small scale.

margin explain cmnist

add image that were  
drawn by ff network

However, together with a fully convolutional architecture, the approach did not seem to scale to larger images.

Thus, an optimization-based approach became favorable as it offers good approximations at high resolutions with manageable computational overhead.

## Optimization Procedure

The previous section ?? already specified why an optimization-based approach is preferable over a feed-forward or recurrent approach. This section aims to give more details about the optimization procedure.

### Rendering Layout

Fundamentally, the optimization procedure is inspired by stroke-based rendering procedures. It could also be compared to the style transfer approach by ~gatys [5], with parameters optimized instead of pixels. The difference to normal stroke-based renderers, though, is the limited size of a rendered brushstroke in this work's renderer.

This poses a significant challenge that might not be obvious at first.

Ideally, the optimization procedure should be able to place strokes freely on the canvas, as this allows for an unbiased approximation. Furthermore, it would allow allocating many small strokes in areas where the artist placed many strokes and use fewer and wider strokes in other areas. However, due to various limitations, which were explained in ??, the renderer is not able to render single brushstrokes in a 1MP frame. Similar approaches only perform on relatively small canvas sizes, likely due to this issue .

[add ref](#)

Dissecting the target image into many small patches and then running the optimization procedure on these individual patches is a workaround for this. After the patches are

add figure to this

rendered, they are then fused along their edges to give the full image. A good comparison is a grid of renders where each grid cell is the center of many renderings at the same time. The major issue of such an approach is the adjacent edges where the grid cells are joined. Also, a grid structure will almost always differ from the inherent distribution of brushstrokes in an image; grid edges will, more often than not, separate brushstrokes between two grid cells. A simple solution to this problem is transitioning from a stacked grid structure to an overlapping grid.

graphics!!!!

By overlapping the grid cell, obvious edges between render windows are hidden as every edge lies within the frame of a different render window. Choosing a lattice vector size smaller than the render window's dimensions such a grid can easily be realized.

Still, this kind of initialization requires a very even distribution of brushstrokes, ideally with a brushstroke at the center of each cell. As this is not the case, and stroke densities will vary locally, the grid layout is prone to erroneously enforce a grid-like layout of strokes where there is none. This is due to the inability of the grid to account for local changes in density and the following propagation of error. Starting with a single region of high-density strokes in the vicinity of one grid cell, this cell would ideally render a narrow stroke to achieve high accuracy. Neighboring cells then have to shift their strokes towards the center of that grid cell to account for the free space that is not covered by the narrow brushstrokes. This shift must then be accounted for by the next neighboring cells and so on, which will cause all strokes in a row or column to shift towards this one spot with a high stroke density.

Now, a painting usually has many such high-density regions, which possibly cancel the shift that is caused by other areas. As a result, the renderings are likely to not shift at all, as shifting mostly cancels out over the whole image. Subsequently, an area of high stroke density will not have enough strokes available in its local region and thus will be covered up by a single broad stroke as this minimizes the L2 loss.

The core of the problem is the previously imposed lattice structure that propagates local density shifts along its principal axes.

One possible solution is getting rid of the lattice structure and replacing it with a more random structure that also covers the image sufficiently. This can be accomplished by using **superpixels [superpixels]**. Superpixels were popular in the pre-neural network era of computer vision and were often used in image segmentation tasks ([img segmentation with SP]). Nevertheless, superpixels are also a popular starting point for brushstroke extraction algorithms [**brushs stroke extraction**].

Basically, superpixels are pairwise disjoint groups of pixels in an image that would usually join pixels with similar colors in a local region. Straight away, it is obvious why this is interesting for brushstroke extraction. The distribution of superpixels will not follow a grid-layout as the previous approach, and the location of superpixels should relate to the given color distribution in the image. It is easy to imagine that the location of the superpixel centers would be a good prior for locations of render windows as well. Also, as the colors of pixels inside a superpixel should be similar, one can use the mean color of a superpixel as an initialization for the color of the brushstroke.

Ultimately, a superpixel segmentation will be used to infer positions for render windows as well as the color initialization of each stroke.

## Rendering Order

Another problem that will come up during the optimization procedure is the order in which strokes are rendered (already mentioned in Section ??). Real-world brushstrokes are also subject to the same issue as the current brushstroke will always be placed on top of brushstrokes painted earlier, (see section 6.3).

As the optimization-based approach relies on parallel optimization of brushstrokes, it must decide which strokes are in the foreground and which are in the background. Otherwise, as this would randomly change, edges in the image might be obstructed, and optimization could oscillate between solutions where different strokes lie in the foreground. It could also prevent brushstrokes from overlapping such that they cover disjoint areas. All of these outcomes would be unfavorable as it tends to produce worse results in the end.

The solution which is presented in this thesis is an additional parameter that describes a brushstroke's accuracy. The accuracy is defined as the L2 distance of each stroke's pixel to the corresponding pixel in the target image multiplied by this each pixel's alpha value. This removes any pixels which are of no interest from the loss and focusses only on the rendered pixels.

$$\text{accuracy} = 1 - \frac{1}{N} \sum_{p \in \text{pixels}} \|p - p'\|_2^2 \text{ with } \mathbf{p} = \mathbf{p}' \quad (6.6)$$

improve this, especially  
the  $\mathbf{p}=\mathbf{p}'$  part

The resulting value describes how well the pixels of the rendered stroke match their respective pixels in the target image. Consequently, any brushstroke with higher accuracy will be more faithful to the target image than strokes with lower accuracy. Placing these brushstrokes in the foreground should thus result in a smaller L2 loss than the other way around. Vice versa, brushstrokes that connect two same colored areas will aggregate a lower accuracy as the brushstroke is compared at the intersection as well. A rendered brushstroke that fits the foreground brushstroke will not be affected by this, thus keeping its high accuracy and laying on top of the other brush stroke.

Notably, the accuracy should not be included in the brushstroke's loss, as this would prohibit background strokes from covering larger areas and result in behavior that is similar to

the non-overlapping issue previously described. Thus, the accuracy of each stroke will be calculated as it is rendered.

## Initialization

The following section will focus on initialization details for all parameters of a brushstroke, their position, and the confidence value.

Besides the original ten parameters of each brushstroke, which were explained in section 6.2, the previous two sections introduced an accuracy parameter for ordering and two translation parameters that define the position of the render window along each axis. All of these parameters must be initialized before the optimization procedure starts. Ideally, the initialization should not introduce any bias to the optimization process. At the same time, an initialization should facilitate training and accelerate convergence in the early stages of optimization.

Unfortunately, the placement of the render window will surely enforce a bias on the optimization (see section 6.3). Thus, a superpixel initialization was motivated for the translation parameters as well as the color of the brushstrokes. Subsequently, the translation parameter for each render window will be equal to the position of the weighted mass center of its respective superpixel.

The initial color will be taken from the mean color value of the superpixel.

The brush size will be initialized with the minimum possible value. This will let brushstrokes not overlap at the beginning of optimization. Only when the brushstrokes already roughly fit their local region, they shall intersect and be ordered by their accuracy. Therefore, the initial accuracy will be 0 everywhere, as the accuracy is recalculated after every optimization step.

Other patch parameters, notably  $\mathbf{s}$ ,  $\mathbf{e}$  and  $\mathbf{c}$ , will be initialised using a narrow Gaussian distribution with  $\sigma = .1$  and values clipped to  $[-1, 1]$  because there is no prior information available on how the brushstrokes are oriented. Instead, this approach relies on the optimization procedure to be minimally biased by this initialization of the path variables.

### Partial Updates

When building an optimizer based on the information provided until now, GPU memory limitations become a guaranteed issue. Even as it might not be evident at first, the optimization procedure imposes a considerable requirement for memory on the graphics card, due to two parts of the training:

First, the number of brushstrokes can easily become very high with large images as input. With a render window size of 64x64, the brushstrokes are relatively small, and paintings can consist of a few thousand brushstrokes. This would equate to a batch size of a few thousand for the brushstroke renderer. Tests have shown that on the latest hardware with 12GB memory, the maximum number of brushstrokes rendered in parallel is  $\approx 256$ . Obviously, this is one to two orders of magnitude smaller than what would be needed to optimize all strokes in the painting in parallel. Still, there is a way around this bottleneck, by optimizing the image not as a whole but as smaller patches consisting of 256 brushstrokes at a time, giving a partial update routine.

Each patch comprises the 256 nearest render windows to a randomly sampled location on the canvas. These 256 brushstrokes are then rendered from their parameters in order to obtain a gradient later on. Then strokes are ordered according to their accuracy, placed on canvas ('padded') and blended ('stitched'). At last, the loss is calculated and back-propagated to update the parameters of the patch.

When restricting the number of brush strokes, bordercases do become an issue. Brushstrokes that lie at the perimeter of the patch are not fully surrounded by other patches. The result is a discrepancy of what the gradient to the brushstroke would look like if the brushstroke had laid in the middle of the patch. A solution to this is laying a ring of already rendered brushstrokes around the patch. The ring guarantees that all brush strokes that brushstrokes at the border of the patch are surrounded by neighboring brushstrokes the same way, as if they were in the middle of the patch. Notably, the brushstrokes that belong to the ring around the patch, serve as 'dummy' brushstrokes and are optimized. Without any need for optimization, these brushstrokes do not need to be rendered in this optimization step. Instead, they could have been rendered earlier and the rendered image is just reused. For this reason rendered brushstrokes for all parameters are saved in a **render image catalog**. Equally, the collection of brushstrokes that shall later compose the image shall be called a **parameter catalog**. The parameters as well as the images in these catalogs are updated whenever the respective brushstroke parameter has been updated. With the render image catalog at hand, it is possible to use the previously rendered brushstrokes to stitch the image as a whole with the newly rendered brushstrokes of the image patch embedded.

As this will cause the border strokes to be surrounded by other brushstrokes to all sides (even if not all of them are freshly rendered), the effect of the partial update routine vanishes.

The other problem for huge input images is that the stitching of brushstrokes takes up a considerable amount of memory. In detail, each brushstroke must be placed on the virtual canvas individually, where the canvas' size is that of the input image. This would equate to a couple of thousand 1MP images being stored in memory before they are stitched to a single 1MP image. As a single 1MP image carries roughly  $1000 \times 1000 \times 4 \text{channels} \times 8 \text{Bit} = 32 \times 10^6 \text{Bit} = 4 \times 10^6 \text{Byte} =$

Images are typically saved as unsigned 8-Bit integers, thus 8-Bit per channel.

4MB of information, a few thousand of these will easily exceed the memory of most graphics cards.

Luckily, the previous workaround for optimizing only 256 brushstrokes will work as well without rendering the full image at every step. Since most brushstrokes are not re-rendered and thus will not be supplied with a gradient, their main task is to regulate losses for edge strokes of the image patch. As this does not need far away strokes, but only those close to the optimized strokes, a ring of pre-rendered strokes around the re-rendered patch will suffice.

approximate the number  
of brushstrokes needed  
to surround the patch

This reduces the number of involved brushstrokes per optimization step from a few thousand down to a couple hundred. Besides allowing for the partial update routine to be performed at all, it should also increase performance significantly compared to an approach that involves all brushstrokes.

## Image Placing & Blending

As the update and optimization procedure has now been explained thoroughly, it is now time to explain the process of placing and blending a rendered brushstroke a bit further.

After each stroke has been rendered, it needs to be placed according to the translation parameters (see 6.3). This requires dynamically placing each brushstroke inside a zero-filled tensor.

By calculating each pixel's global positioning in the rendered image individually, it is possible to scatter the pixels of the original rendered image into the larger zero-filled tensor and obtain a globally placed brushstroke.

The task of blending the resulting canvas-sized rendered brushstrokes together, or stitching them, is more complicated, though. Due to the canvas' alpha channel, it is possible to blend only relevant information while the rest of the image will be ignored.

As far as conventional alpha blending goes, two images are blended by multiplying each pixel value with the alpha value of the top-layer image while the background image is multiplied with the complement to the alpha value:

$$p_{x,y} = p_{x,y}^{\text{top}} \times \alpha_{x,y}^{\text{top}} + p_{x,y}^{\text{bottom}} \times (1 - \alpha_{x,y}^{\text{top}}) \forall (x, y) \in \mathcal{D}(\text{image}) \quad (6.7)$$

correct this equation and make it nicer

For multiple layers, this process can be repeated in various fashions, after the strokes are ordered according to their accuracy. Either one could start from the bottom and blend the two back-most strokes, followed by the next third last strokes and so one, or one could start this process from the front with the two strokes in the very front being blended at first, then the stroke with the third-highest accuracy *etc..*

Both of the approaches would yield the same result but differ only in the order in which they were blended. Subsequently, both methods will have  $(n - 1)$  blending operations to compute per pixel.

Blending brushstrokes in a position-aware manner can reduce this number. The majority of pixels for each padded brushstroke is non-informative, as the alpha value is zero (see Figure ??). This opens the possibility of going from blending *all* pixels of *all* brush strokes to blending just those pixels with non-zero alpha values. As before there were many layers that represented one single brush stroke each, few layers that do not represent brush strokes can achieve the same result. Instead of representing brush strokes, these few layers represent the order in which pixels should be blended. This means going from a layer-focussed approach to a rather pixel-focussed approach to alpha-blending.

The easiest way to accomplish this, is to first find the maximum blending-depth over all pixels, where the blending-depth  $k$  is the number of layers where the alpha value is not

zero.

$$k = \arg \max_{p \in \text{pixels}} \sum_{i \in \#\text{layers}} \mathbb{1}(\alpha_i > 0) \quad (6.8)$$

Then the top  $k$  layer indices for each pixel are picked, which reduces the number of blending operations from  $(n - 1)$  to  $(k - 1)$ . Importantly, the top  $k$  indices should not be ordered by their alpha values but remain in the order that was imposed by sorting according to the accuracy value. Otherwise, the order will most certainly be mixed up, and the pixel with the highest alpha value will always lie on top instead of the pixel that belongs to the most accurate brushstroke. Especially, as brushstroke renderings fade out towards their edges, this makes a significant difference.

Another way of accelerating the process of alpha-blending is **vectorizing**. Instead of iteratively applying the computations, a vectorized operation can perform these computations in parallel. Vectorizing makes it necessary to construct a tensor with the following properties:

For  $\mathbf{I} \in [0, 1]^{H \times W \times 4}$  the image target, the shape will be defined as  $\mathcal{S}(\mathbf{I}) = (H, W, 4)$ . Each alpha channel will have the values  $\alpha^{hw} \in [0, 1]$  for  $h = 0, \dots, H$  and  $w = 0, \dots, W$ .

The set of rendered and padded brushstrokes  $\mathbf{J}$  will have the shape  $(N, H, W, 4)$  with  $N$  depicting the number of brushstrokes that ought to be stitched simultaneously.

Now, looking at each individual pixel in  $\mathbf{J}$ , which is described by  $(z_n^{hw}, \alpha_n^{hw})$  for  $n = 1, \dots, N$  and  $z_n^{hw} \in [0, 1]^3$ ,  $z^{hw}$  describes the RGB values and  $\alpha^{hw}$  the alpha-channel for a pixel at  $(h, w)$ .

A blending operation can then be defined by

$$z'^{hw} = \tilde{\alpha}^{hw} \cdot z^{hw} \quad (6.9)$$

$$\text{or} \quad (6.10)$$

$$z'^{hw} = \sum_{n=1}^N \tilde{\alpha}_n^{hw} z_n^{hw} \quad (6.11)$$

$$(6.12)$$

with  $z'^{hw}$  the resulting RGB values of the blended pixel and  $\tilde{\alpha}^{hw}$  a vector that holds the merged alpha values for each pixel:

$$\tilde{\alpha}^{hw} = \begin{pmatrix} \alpha_1^{hw} \\ \alpha_2^{hw} & (1 - \alpha_1^{hw}) \\ \alpha_3^{hw} & (1 - \alpha_2^{hw}) & (1 - \alpha_1^{hw}) \\ \vdots \end{pmatrix} \quad (6.13)$$

$$= \alpha^{hw} \odot \begin{pmatrix} 1 \\ (1 - \alpha_1^{hw}) \\ (1 - \alpha_2^{hw}) & (1 - \alpha_1^{hw}) \\ \vdots \end{pmatrix} \quad (6.14)$$

$$\rightarrow \tilde{\alpha}_n^{hw} = \alpha_n^{hw} \prod_{i=1}^{n-1} (1 - \alpha_i^{hw}) \quad (6.15)$$

with  $\odot$  the element-wise product

What is left, is to find a way to construct  $\tilde{\alpha}^{hw}$  from  $\alpha^{hw}$ .

For this an auxiliary matrix  $\beta^{hw}$  is constructed:

$$\beta^{hw} = \alpha^{hw} \times \mathbb{1}_{1 \times N} = \begin{pmatrix} \alpha_1^{hw} & \alpha_2^{hw} & \dots & \alpha_n^{hw} \\ \alpha_1^{hw} & \alpha_2^{hw} & \dots & \alpha_n^{hw} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_1^{hw} & \alpha_2^{hw} & \dots & \alpha_n^{hw} \end{pmatrix} \quad (6.16)$$

with

$$\mathbb{1}_{1 \times N} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}^T$$

Then  $\beta^{hw}$  is strictly triangulated such that:

$$\gamma^{hw} = \beta^{hw} \odot \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ 1 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 1 & 1 & \dots & 1 & 0 \end{pmatrix} \quad (6.17)$$

$$= \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ \alpha_1^{hw} & 0 & 0 & \dots & 0 \\ \alpha_1^{hw} & \alpha_2^{hw} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \alpha_1^{hw} & \alpha_2^{hw} & \dots & \alpha_{n-1}^{hw} & 0 \end{pmatrix} \quad (6.18)$$

$$\rightarrow \delta^{hw} = 1 - \gamma^{hw} = \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ (1 - \alpha_1^{hw}) & 1 & 1 & \dots & 1 \\ (1 - \alpha_1^{hw}) & (1 - \alpha_2^{hw}) & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ (1 - \alpha_1^{hw}) & (1 - \alpha_2^{hw}) & \dots & (1 - \alpha_{n-1}^{hw}) & 1 \end{pmatrix} \quad (6.19)$$

By multiplying the elements along each row in  $\delta^{hw}$  one gets:

$$\epsilon_i^{hw} = \prod_{j=1}^N \delta_{ij}^{hw} = \begin{pmatrix} 1 \\ (1 - \alpha_1^{hw}) \\ \vdots \\ \prod_{j=1}^{N-1} (1 - \alpha_j^{hw}) \end{pmatrix} \quad (6.20)$$

$$\rightarrow \tilde{\alpha}^{hw} = \epsilon^{hw} \odot \alpha^{hw} = \begin{pmatrix} \alpha_1^{hw} & & \\ \alpha_2^{hw} & (1 - \alpha_1^{hw}) & \\ \alpha_3^{hw} & (1 - \alpha_2^{hw}) & (1 - \alpha_1^{hw}) \\ & \vdots & \\ \alpha_N^{hw} & \prod_{j=1}^{N-1} (1 - \alpha_j^{hw}) & \end{pmatrix} \quad (6.21)$$

This vectorized version of alpha blending will introduce a new possible bottleneck as it is, since  $\beta^{hw}$  will be a tensor of shape  $(N, N, H, W)$ , which will equate to

$$256 \times 256 \times 256 \times 256 \times 4B = 2^{36}B = 64\text{GiB}$$

alone.

This is where the previous position-aware alpha blending tricks becomes useful. By computing  $\beta^{hw}$  only through the top  $k$  values of  $\alpha^{hw}$  instead of the full tensor  $\alpha^{hw}$ , the size will be reduced to

$$k \times k \times 256 \times 256 \times 4B = k^2 2^{22}B = k^2 \times 4\text{MiB}$$

as the shape is reduced to  $(k, k, H, W)$ .

Ultimately, this accelerates optimization by a factor of  $2 - 3$ , as it will be shown in ??.

It must be mentioned that the upper boundary for computational complexity in using this kind of alpha-blending is  $\mathcal{O}(N \log k)$ , since the top  $k$  search is bound by this complexity.

calcualte this

## Losses

In this section, the different kinds of losses for the optimization procedure will be discussed.

First off, the L2 loss or **mean squared error** is an obvious choice for this task.

$$L_{\text{MSE}} = \frac{1}{HW} \sum_{p \in \text{pixels}} \|z(p) - z'(p)\|_2^2 \quad (6.22)$$

Since the MSE loss focuses on minimizing the pixel-wise error between the target image and the fully rendered approximation by brushstrokes, it will cause the rendered image to match the target image mainly in color. At the same time, MSE loss is prone to blurring, which results in washed out edges in the rendered image. Thus, this loss is must be accompanied by additional losses to make up for the shortcoming of MSE loss.

A popular choice for preserving the content in an image, which is associated with preserving edges, is **perceptual loss**. Perceptual loss is based on a VGG Network [VGG] that is pre-trained on ImageNet [ImageNet]. To compute the loss, the activations of deep layers (usually the fourth convolutional block) of the pre-trained VGG network are inferred and then compared using MSE loss.

$$L_{\text{perceptual}} = \frac{1}{H_f, W_f} \sum_{p \in \text{pixels}} \|f(p) - f'(p)\|_2^2 \quad (6.23)$$

The resulting distance is meant to capture how well edges between the two input images are preserved which should be equal to whether the content in both images is the same. Together with MSE loss, a perceptual loss is often used to get better reconstructions than with MSE loss alone, as edges of objects in the image are better preserved, prohibiting blurriness that would occur otherwise.

As an extension to perceptual loss, **lips** introduced perceptual similarity or **LPIPS loss**, which weighs the different

layers of the pre-trained VGG-network differently in order to increase the effectiveness of perceptual distance between two images [lips]. LPIPS loss is meant to preserve edges even better than perceptual loss does with a similar computational overhead.

Besides losses that operate in pixel space, it is also necessary to restrict the action space for each brushstroke. As explained in section ??, the renderer has been trained on a limited data set, which puts constraints on how curved brushstrokes may be and how the ratio between length and width ought to look like. These constraints must be enforced in the optimizing process as well. Because the renderer has not been trained on data outside of the generated data set, the renderer will likely break if the input parameters lie too far outside the training space. The results would then be renderings with no output, distorted brushstrokes, or just noisy output, as seen in ??.

Thus, one must think of an additional loss to confine the parameter space to the same space as the generated brushstroke data during optimization. There are two ways of achieving this:

- ▶ Discriminators
- ▶ Explicitly coded losses

**Discriminators** are a popular choice in this context because even if the data distribution is not known beforehand, a discriminator is still able to learn the distribution from data and thus point out wrong parameter combinations in this case. Still, a discriminator comes with a few compromises, as the target distribution will never be entirely learned rather than well approximated by it. This leaves room for weaknesses as well as local minima in the discriminator's prediction, which would result in worse quality for this task. Usually, these weaknesses are made up for during adversarial training as if the generator overfits to such weaknesses, and the discriminator will quickly penalize such a solution. In the

optimization routine, which is employed for this problem, it is not possible to train the discriminator online as the limited amount of brushstrokes will allow the discriminator to overfit the problem easily. Thus only a pre-trained discriminator with its said weaknesses can be used in this case.

**Handpicked losses** As the data distribution for the generated data set is actually known in this case (see 6.2), it is also possible to manually define losses that confine the brushstrokes. The width constraint – as a first example – can easily be enforced by penalizing whenever the brushstroke's width  $w(x)$  is more than half the length  $l(x)$  between the start point  $\mathbf{s}(x)$  and the end point  $\mathbf{e}(x)$  of the brushstroke  $x$ :

$$L_{\text{bs}} = \frac{1}{|X|} \sum_{x \in X} \max(0, 2w(x) - l(x)) \quad (6.24)$$

$$= \frac{1}{|X|} \sum_{x \in X} \max(0, 2w(x) - \|\mathbf{s}(x) - \mathbf{e}(x)\|_2) \quad (6.25)$$

This is a bit more complicated regarding the limitation that is introduced to the control point  $\mathbf{c}(x)$ . As  $\mathbf{c}(x)$  was sampled from a multivariate Gaussian with fixed parameters in the data set, it should now follow a similar distribution in relation to the direction  $\mathbf{s}(x) - \mathbf{e}(x)$  of each stroke.

This can be achieved by first defining two orthonormal basis vectors which are either parallel  $\mathbf{n}_{se}^{\parallel}(x)$  or orthogonal  $\mathbf{n}_{se}^{\perp}(x)$  to the directional vector  $\mathbf{s}(x) - \mathbf{e}(x)$ :

$$\mathbf{n}_{se}^{\parallel}(x) = \frac{\mathbf{s}(x) - \mathbf{e}(x)}{\|\mathbf{s}(x) - \mathbf{e}(x)\|_2} \quad (6.26)$$

$$\mathbf{n}_{se}^{\perp}(x) = R_{\pi/2} \frac{\mathbf{s}(x) - \mathbf{e}(x)}{\|\mathbf{s}(x) - \mathbf{e}(x)\|_2} \quad (6.27)$$

$$\text{with } R_{\pi/2} = \begin{pmatrix} \cos \pi/2 & -\sin \pi/2 \\ \sin \pi/2 & \cos \pi/2 \end{pmatrix} \quad (6.28)$$

Then  $\mathbf{c}(x)$  can be projected into the coordinate system spanned

by  $\mathbf{n}_{se}^{\parallel}(x)$  and  $\mathbf{n}_{se}^{\perp}(x)$ :

$$c^{\parallel}(x) = (\mathbf{c}(x) - \mathbf{a}(x)) \cdot \mathbf{n}_{se}^{\parallel}(x) \quad (6.29)$$

$$c^{\perp}(x) = (\mathbf{c}(x) - \mathbf{a}(x)) \cdot \mathbf{n}_{se}^{\perp}(x) \quad (6.30)$$

$$\mathbf{a}(x) = \frac{\mathbf{s}(x) + \mathbf{e}(x)}{2} \quad (6.31)$$

Now, the axes of the original multivariate distribution co-align with  $c^{\parallel}(x)$  and  $c^{\perp}(x)$ . By calculating the mean and standard deviation along these projections, they can be compared to the parameters of the original data distribution.

$$\mu = \frac{1}{|X|} \sum_{x \in X} \begin{pmatrix} c^{\parallel}(x) \\ c^{\perp}(x) \end{pmatrix} \quad (6.32)$$

$$\Sigma = \left( \frac{1}{|X|} \sum_{x \in X} \begin{pmatrix} c^{\parallel}(x) \\ c^{\perp}(x) \end{pmatrix} - \mu \right) \left( \begin{pmatrix} c^{\parallel}(x) \\ c^{\perp}(x) \end{pmatrix} - \mu \right)^T \frac{1}{2} \quad (6.33)$$

$$(6.34)$$

Using the Kullback-Leibler divergence for multivariate Gaussian distributions, the compliance with the data sets distribution can be checked:

$$\mathcal{L}_{\text{KL}} = \frac{1}{2} \left[ \log \frac{|\Sigma|}{|\tilde{\Sigma}|} - d + \text{tr}(\Sigma^{-1} \tilde{\Sigma}) + (\mu - \tilde{\mu})^T \Sigma^{-1} (\mu - \tilde{\mu}) \right] \quad (6.35)$$

with

$$\tilde{\mu} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (6.36)$$

$$\tilde{\Sigma} = \begin{pmatrix} \frac{1}{200} & 0 \\ 0 & \frac{1}{25} \end{pmatrix} \text{ see (6.4)} \quad (6.37)$$

since the origin of the projection is at  $\mathbf{a}(x) = \frac{\mathbf{s}(x)+\mathbf{e}(x)}{2}$  and coincides with the center of the data distribution.

In theory,  $L_{\text{KL}}$  and  $L_{\text{bs}}$  should be able to capture any deviation from the source data and ensure that parameters stay within the training space of the renderer. One problem that obviously could arise with this formulation is for patches with very similar brushstrokes that show a mean other than  $\tilde{\mu}$  as well as a very low values inside the covariance matrix and non-zero off-diagonal values. Thus it will be favorable to include as many strokes as possible when calculating this loss, ideally all strokes in the global parameter catalog.

## Style Transfer

One question that arises when discussing the optimization procedure is whether style transfer could be performed with this approach. Especially, since the optimization procedure has explicitly been compared to the approach by Gatys *et al.* [5], it is natural to assume that such an approach could also be applicable in this case.

This mainly requires to introduce a style loss as it has been done by Gatys *et al.* [5] since a content loss is already in place (see (6.23)). A style loss can then be implemented similarly by aggregating the activations of more layers and then calculating the gram matrices:

$$\mathcal{L}_{\text{style}} = \sum_{l=0}^L \frac{w_l}{4N_l^2 M_l^2} \sum_{ij} (G_{ij}^l - A_{ij}^l)^2 G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l \quad (6.38)$$

The problem, which arises when trying to apply style transfer with this particular approach is the partial update routine. Since the gram-matrices are meant to catch global second-order statistics of the image, a small patch would be object to a wrongful assumption that the patch represents the whole image. The only chance of dealing with a local patch in a global context in image space is a cached version of the whole current image in which the patch is embedded.

## Optimization Details

As it has been explained in Section ??, it is not possible to optimize brushstrokes for the whole image in parallel. This is partly due to the neural renderer's memory requirements, which scales with the number of rendered patches. The memory requirements of placing and blending, on the other hand, scales with the number of rendered brushstrokes as well as the number of surrounding fixed brushstrokes and the patch window size. Experiments have shown that a combination of 256 rendered brushstrokes with 128 surrounding brushstrokes paired with a patch window size of 320x320 pixels occupies around 9.5GiB of memory which leaves enough space for more advanced losses and tweaks to the network architecture.

The learning rate for the optimization procedure can be significantly larger than for training the neural renderer. With a learning rate of 0.01, the optimization procedure will converge significantly faster to a solution without any instability issues.

One choice, which has to be made individually per target image is how many brushstrokes will cover the image. As larger images obviously require more brushstrokes than smaller images, what should remain the same is the **brush-strokes density**. The brushstroke density will decide how many pixels on average should be covered by each brushstroke and thus be used during initialization. A number of 100pixels/brushstroke has produced the best results during the experiments.

Another important choice that goes along the choice of how dense the brushstrokes should be distributed is that of how many optimization steps each brushstroke will be object to. As for too few steps, the training will not have converged, and for too many, optimization will take an unnecessary amount of time. As this can vary between images, since some images require more time to converge, about 1500steps/brushstroke have been empirically chosen. There is no direct enforcement

that each brushstroke is updated exactly this often, but it can be expected that due to the uniform sampling of render patches, there will be no major deviation for some brushstrokes. Since 256 brushstrokes will be optimized in every step, the total number of optimization steps can be calculated together with the brushstroke density and the image's size.

Another minor detail is the fact that for each render patch, five consecutive optimization steps are performed as this safes memory bandwidth since data must be written and read from memory each time a different render patch is optimized.

## Results

Results can be seen in Figure ?? for Starry Night as the standard reference image of this thesis. It took approximately 15,000 optimization steps, which corresponds to about 2h. The rendering consists of roughly 10,000 brushstrokes.

Figure ?? shows the result for a natural photo as target image which took also 2h to compute with 8.000 brushstrokes and 12.000 optimization steps.

# Evaluations & Ablation Experiments

# 7

In this chapter, the ablation experiments for the approaches presented in Chapter ?? shall be presented. This is meant to show some weaknesses of each method as well as give a visual understanding of the effects of some losses and parameters. Together with an interpretation it should help to get a better understanding of previous motivations.

7.1 Neural Renderer . . . . .	107
7.2 Optimization Procedure . . . . .	107

## 7.1 Neural Renderer

The neural renderer will be evaluated with respect to a few aspects. First, the performance for larger resolution images will be evaluated. Then, the effect of the discriminator will be shown by training the renderer without the discriminator loss. At last, the weakness of the renderer, which is data points outside the training space, will be analyzed.

## 7.2 Optimization Procedure

The optimization procedure gives many starting points for ablation experiments. At first, the effects of each loss will be visually compared. Then, different initializations will be shown as well as various choices for training parameters.



# **CONCLUSION**



# Discussion 8



# Outlook | 9



## **APPENDIX**



# Appendix A



# Bibliography

Here are the references in citation order.

- [1] Jia Li *et al.* ‘Rhythmic Brushstrokes Distinguish van Gogh from His Contemporaries: Findings via Automated Brushstroke Extraction’. In: IEEE Transactions on Pattern Analysis and Machine Intelligence 34.6 (2012), pp. 1159–1176. doi: [10.1109/tpami.2011.203](https://doi.org/10.1109/tpami.2011.203) (cited on pages 44, 45).
- [2] Fabrizio Lamberti, Andrea Sanna, and Gianluca Paravati. ‘Computer-assisted analysis of painting brushstrokes: digital image processing for unsupervised extraction of visible features from van Gogh’s works’. In: EURASIP Journal on Image and Video Processing 2014.1 (2014), p. 53. doi: [10.1186/1687-5281-2014-53](https://doi.org/10.1186/1687-5281-2014-53) (cited on pages 44, 45).
- [3] Aaron Hertzmann *et al.* ‘Image analogies’. In: (2001), pp. 327–340. doi: [10.1145/383259.383295](https://doi.org/10.1145/383259.383295) (cited on page 46).
- [4] Hochang Lee *et al.* ‘Directional texture transfer’. In: (2010), pp. 43–48. doi: [10.1145/1809939.1809945](https://doi.org/10.1145/1809939.1809945) (cited on pages 46, 47).
- [5] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. ‘A Neural Algorithm of Artistic Style’. In: (2015) (cited on pages 47–50, 52, 53, 62, 87, 104).
- [6] Yanghao Li *et al.* ‘Demystifying Neural Style Transfer’. In: arXiv (2017) (cited on page 51).
- [7] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. ‘Perceptual Losses for Real-Time Style Transfer and Super-Resolution’. In: (2016) (cited on pages 52, 53).
- [8] Xun Huang and Serge Belongie. ‘Arbitrary Style Transfer in Real-time with Adaptive Instance Normalization’. In: (2017) (cited on pages 52, 53).
- [9] Yijun Li *et al.* ‘Universal Style Transfer via Feature Transforms’. In: (2017) (cited on pages 53, 54).
- [10] Jun-Yan Zhu *et al.* ‘Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks’. In: (2017) (cited on pages 54, 55).
- [11] Arsiom Sanakoyeu *et al.* ‘A Style-Aware Content Loss for Real-time HD Style Transfer’. In: (2018) (cited on pages 55, 56).
- [12] Paul Haeberli. ‘Paint by numbers: abstract image representations’. In: ACM SIGGRAPH Computer Graphics 24.4 (1990), pp. 207–214. doi: [10.1145/97879.97902](https://doi.org/10.1145/97879.97902) (cited on pages 57–59).
- [13] A. Hertzmann. ‘A survey of stroke-based rendering’. In: IEEE Computer Graphics and Applications 23.4 (2003), pp. 70–81. doi: [10.1109/mcg.2003.1210867](https://doi.org/10.1109/mcg.2003.1210867) (cited on page 58).

- [14] Siddharth Hegde, Christos Gatzidis, and Feng Tian. ‘Painterly rendering techniques: a state-of-the-art review of current approaches’. In: Computer Animation and Virtual Worlds 24.1 (2013), pp. 43–64. doi: [10.1002/cav.1435](https://doi.org/10.1002/cav.1435) (cited on page 58).
- [15] Peter Litwinowicz. ‘Processing images and video for an impressionist effect’. In: (1997), pp. 407–414. doi: [10.1145/258734.258893](https://doi.org/10.1145/258734.258893) (cited on pages 58, 59).
- [16] Bill Baxter *et al.* ‘DAB: interactive haptic painting with 3D virtual brushes’. In: (2001), pp. 461–468. doi: [10.1145/383259.383313](https://doi.org/10.1145/383259.383313) (cited on page 60).
- [17] David Ha and Douglas Eck. ‘A Neural Representation of Sketch Drawings’. In: arXiv (2017) (cited on page 60).
- [18] Alex Graves. ‘Generating Sequences With Recurrent Neural Networks’. In: arXiv (2013) (cited on page 60).
- [19] Yaroslav Ganin *et al.* ‘Synthesizing Programs for Images using Reinforced Adversarial Learning’. In: (2018) (cited on pages 61, 69, 71).
- [20] Ningyuan Zheng, Yifan Jiang, and Dingjiang Huang. ‘Strokenet: A neural painting environment’. In: (2019) (cited on pages 61, 62).
- [21] Reiichiro Nakano. ‘Neural Painters: A learned differentiable constraint for generating brushstroke paintings’. In: (2019) (cited on page 62).
- [22] Zhewei Huang, Wen Heng, and Shuchang Zhou. ‘Learning to Paint With Model-based Deep Reinforcement Learning’. In: (2019) (cited on page 62).
- [23] Biao Jia *et al.* ‘PaintBot: A Reinforcement Learning Approach for Natural Media Painting’. In: arXiv (2019) (cited on pages 62, 63).
- [24] Biao Jia *et al.* ‘LPaintB: Learning to Paint from Self-Supervision’. In: (2019) (cited on page 62).

## **List of Figures**

6.7 Visualization of the generator and discriminator architectures . . . . .	79
--	----

## **List of Tables**