

## DEEP RESTORE

Using Image Segmentation and Image In-painting Techniques

Matthias Reiterer, BSc

*Institute of Computer Graphics and Vision  
Graz University of Technology, Austria*

Seminar/Project Computer Vision  
*Supervisor Prof. Dr. Thomas Pock*  
Graz, September 27, 2018

## Abstract

Analog movies, recorded on roll film, have the property that their image quality worsens statically either from usage or by time. To prevent this deterioration from happening, analog movies often get digitalized. The digitalization of analog movies has the advantage that during post-processing steps the image quality of previously analog movies can be improved - this is known as film restoration.

In this project, we will focus on such a post-processing step, namely the removing of impurities in single image frames. Based on previous research in the field of image restoration we came up with a step-wise solution: first, we detect the impurities of each image and generate a segmentation map (image segmentation), and second, based on this segmentation map we fill the impurities by using information from surrounding pixels (image in-painting).

Throughout this project, we mainly focus on modern machine learning approaches, especially convolutional neural networks (CNNs). For the image segmentation we test different architectures and regularization techniques to come up with a suitable architecture which gives the best results for the segmentation map.

After we generated the segmentation map we can perform image in-painting of the impure pixels. Again, we test different CNN-architectures and another method, image in-painting based on Perona-Malik-Diffusion.

We evaluate our approaches on ground truth data given from HS-Art [10]. Due to a lack of training data and poor ground-truth segmentation maps, we do not reach top results for image segmentation. The image in-painting results in visually pleasing images.

**Keywords:** Computer Vision, Film Restoration, Machine Learning, Image Segmentation, Image In-painting

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Related Work</b>	<b>5</b>
2.1	Image Segmentation . . . . .	5
2.2	Image In-painting . . . . .	5
<b>3</b>	<b>Image Segmentation</b>	<b>6</b>
3.1	Single Image . . . . .	6
3.1.1	Simple . . . . .	6
3.1.2	U-Net . . . . .	7
3.1.3	Densenet . . . . .	7
3.2	Flow . . . . .	7
3.2.1	Early . . . . .	8
3.2.2	Late . . . . .	9
3.2.3	Late-Late . . . . .	9
3.2.4	U-Net . . . . .	10
3.2.5	Densenet . . . . .	10
3.2.6	Densenet Combined . . . . .	11
3.2.7	Bottleneck . . . . .	11
3.2.8	Densenet Combined 2 . . . . .	12
3.3	Regularization Techniques . . . . .	12
3.3.1	Early Stopping . . . . .	12
3.3.2	Dropout . . . . .	13
3.3.3	$L_2$ -Regularization (Weight Decay) . . . . .	13
3.3.4	Input Normalization . . . . .	13
3.3.5	Batch Normalization . . . . .	13
<b>4</b>	<b>Image In-painting</b>	<b>14</b>
4.1	CNN . . . . .	14
4.1.1	U-Net . . . . .	14
4.1.2	Densenet . . . . .	14
4.2	Perona-Malik In-painting . . . . .	15
<b>5</b>	<b>Evaluation</b>	<b>16</b>
5.1	Dataset . . . . .	16
5.2	Implementation Details . . . . .	17
5.3	Results for Image Segmentation . . . . .	17
5.3.1	Single Image . . . . .	17
5.3.2	Flow . . . . .	17
5.4	Results for Image In-painting . . . . .	22
<b>6</b>	<b>Conclusion</b>	<b>23</b>

# 1 Introduction

”Deep Restore” - a combination of the two terms ”deep learning” and ”image restore”. As the name of this project suggests, we focus our considerations on image restoration based on deep learning methods. This project is done in cooperation with HS-Art [10], who specialize in fields of archive film restoration and preservation. All the images used in this project are provided by them.

Image restoration describes the process of determining undistorted images from corrupted ones. It is used in a wide range of image and video processing applications and a basic method in computer vision. There exist many different approaches and an overview on some techniques can be found in Section 2.

Deep learning is a part of machine learning methods based on feature learning. This learning can either happen supervised, semi-supervised or unsupervised. There are different deep learning architectures like *e.g.* recurrent neural networks (RNNs), Boltzmann machines, convolutional neural networks (CNNs), *etc.* Deep learning is nowadays present in different fields like: computer vision, speech and audio recognition, bioinformatics and others.

Since we focus on image restoration in this project, we decide to use CNNs to estimate undistorted images. CNNs are feed-forward neural networks and are typically applied to inputs with local structure. Therefore, CNNs are mainly used in combination with images, video or audio. We choose to train our CNNs in a supervised manner. Input images and ground-truth outputs combined, are used to estimate the parameters of the networks.

As already mentioned the goal of this project is to perform image restoration based on deep learning methods. Given a corrupted image (Figure 1a), we want to determine an undistorted image (Figure 1b).

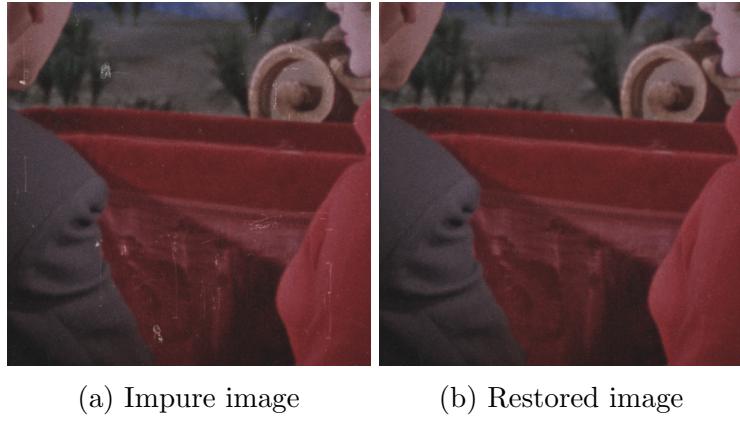


Figure 1: An example of image restoration. Figure (a) shows an image with impurities, whereas Figure (b) shows an undistorted/restored variant of (a).

By investigating our given images and by performing research, we conclude that we can categorize the impurities of corrupted images into: dust (Figure 2), noise (Figure 3), scratches (Figure 4), other damages (like bent frames (Figure 5), *etc.*) and for subsequent image frames flicker (Figure 6). During this theses we just cover and restore: dust, noise, scratches and other damages. (The example images from Figure 2-6 were taken from <https://www.youtube.com/watch?v=neZdImdW6BI>, HS-Art Demo Reel, accessed 24<sup>th</sup> September 2018)



(a) Impure image

(b) Restored image

Figure 2: An example of dust restoration.



(a) Impure image

(b) Restored image

Figure 3: An example of noise restoration.



(a) Impure image

(b) Restored image

Figure 4: An example of scratch restoration.



(a) Impure image

(b) Restored image

Figure 5: An example of restoration for a bent frame.

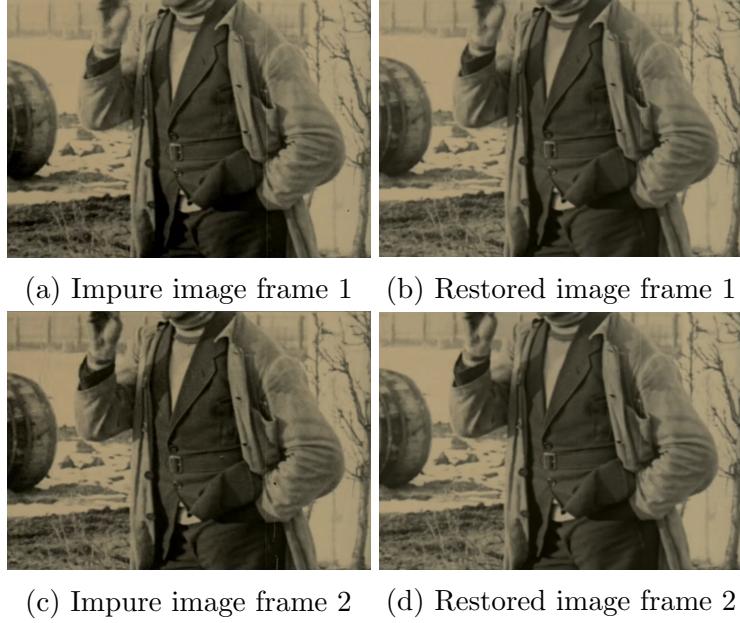


Figure 6: Example of flicker restoration for two following frames.

We split the workflow of image restoration into two consecutive steps: image segmentation and image in-painting.

In the image segmentation step we generate a segmentation map based on the input image. Image segmentation is the process of dividing an image into different classes, where each class represents a segment. In our case we have two classes, one class representing impurities and the other clear pixels. The resulting segmentation map highlights impurities which are then removed.

The next step is image in-painting. Image in-painting describes the process of filling in unknown pixels with information based on neighboring pixels. This is done by using the previously created segmentation map and the input image.

This workflow can be seen in Figure 7. In the first step we determine the segmantation map (Figure 7b), where impurities are highlighted in white, from the input image (Figure 7a). Then, we perform image in-painting (Figure 7c) based on the segmentation map (Figure 7b) and the input image (Figure 7a).

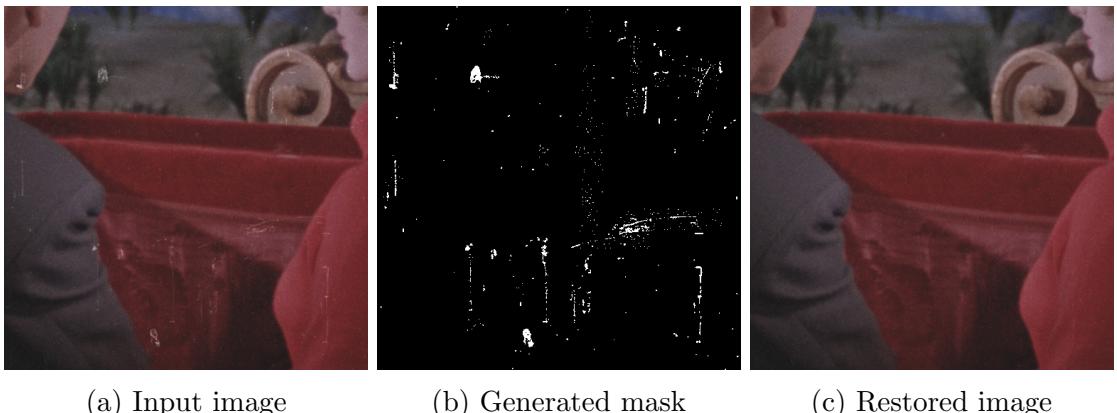


Figure 7: An example workflow of image restoration.

## 2 Related Work

Image restoration is a fundamental part of computer vision. There exist many different approaches and an overview on different techniques is given by Milanfar [15]. Most of the current work on image restoration focuses on developing methods for achieving visual great looking results. This is mostly done by the cost of computational efficiency. Chen et al. ([5] and [6]) propose methods based on diffusion processes, which also take computational efficiency into account.

Since we decided to split up the task of image restoration into image segmentation and image in-painting, the following Sections give an overview of different image segmentation techniques (Section 2.1) and image in-painting techniques (Section 2.2).

### 2.1 Image Segmentation

Driven by a wide range of application, image segmentation received a significant amount of attention from the research community. Image segmentation is used in the medical sector for *e.g.* segmenting cancer cells, in pedestrian detection, face recognition systems and others.

Yuheng et al. [27] categorize segmentation algorithms into: (i) region-based threshold segmentation, (ii) regional growth segmentation, (iii) edge detection segmentation, (iv) segmentation based on clustering, and (v) segmentation based on supervised learning with CNNs. A newer approach from Qazanfari et al. [20] uses evolutionary computation that is inspired from human behavior to perform image segmentation.

We now focus on CNN based segmentation methods since this approach is most interesting for our project. The recent usage of CNNs for various computer vision tasks but also image segmentation led in partly groundbreaking results.

Long et al. [14] use fully convolutional networks to solve image segmentation. They adapt classification networks, such as VGG net [22] or GoogLeNet [23] into fully convolutional networks by transferring learned representations and combining semantic information from deep layers with appearance information from shallow layers.

In their work, Ronneberger et al. [21] propose a CNN architecture for image segmentation that consists of a contracting path to capture context and a symmetric expanding path that enables precise localization. This U-Net architecture achieves good performance on different biomedical segmentation tasks and is also used in other areas.

Huang et al. [11] introduce in their work a dense convolutional network (DenseNet), in which each layer is connected to every other layer in a feed-forward fashion. Feature-maps of preceding layers are used as input for the current layer. An advantage of this architecture is that the number of parameter gets reduced.

The recent work of Chen et al. [4] proposes new methods to further improve image segmentation with CNNs. They introduce convolution with upsampled filters to control the resolution at which feature responses are computed, a spatial pyramid pooling to segment object at multiple scales and to improve object localization, a combination of probabilistic graphical models and CNNs.

During our project we take inspiration from the work of Huang et al. [11] and Ronneberger et al. [21].

### 2.2 Image In-painting

Image in-painting is an ill-posed inverse problem with no well-defined solution. Its application reaches from removing text overlays or scratches, restoring lossy image transmission, object removal and others.

Guillemot et al. [8] give an overview of different image in-painting techniques and categorizes as follows: (i) in-painting using diffusion partial differential equations (PDEs), (ii) variational in-painting, and (iii) exemplar-based methods. Another category is (iv) image in-painting with the help of deep learning and machine learning.

Diffusion based image in-painting, like the work of Boujena et al. [2], rely heavily on the famous work of Perona and Malik [18] or Weickert [25]. Image diffusion is based on PDEs, where diffusion is described as a physical process. The crucial part to establish a visual appealing in-painting, is to incorporate a diffusion tensor into the equations (a detailed description can be found in Section 4.2).

The work of Getreuer [7] and Neri et al. [16] is based on total variation. The use slight variations of the total variation model for improved results.

Related works in the machine learning sectors are *e.g.* the work of Köhler et al. [12] and Gupta et al. [9]. Köhler et al. use a CNN to perform image in-painting. They state that with an additional segmentation mask, highlighting missing pixels, the results of the in-painting can be improved. Gupta et al. use both CNN and RNN approaches.

Another recent trend in machine learning, that is popular in image in-painting, is the usage of an encoder-decoder structure or Generative adversarial network (GANs). Burlin et al. [3], Yeh et al. [26] and van Noord et al. [24] all use this structure to perform state-of-the-art image in-painting.

### 3 Image Segmentation

In this Section follows a description of our methods and approaches for detecting impurities by image segmentation. As already mentioned we focus on a machine learning approach namely CNNs. Section 3.1 gives an overview of architectures that use just the current image to perform the image segmentation. Section 3.2 shows architectures that use additional information, the previous and next image, to perform image segmentation of the current image.

It should be mentioned that if not stated otherwise, we use ReLU activation functions between layers and a sigmoid output activation function to perform image segmentation, which is in general a classification task. We also apply regularization techniques that are described in Section 3.3. The training patches we use are of size  $128 \times 128$  and depending on the architecture an additional padding is also used because of the usage of valid-convolutions. We choose the cross-entropy-error (CEE) for computing the loss/error between predicted and ground-truth segmentation maps.

#### 3.1 Single Image

Our first approach is to train simple architectures with few parameters to perform image segmentation (workflow: from Figure 7a to Figure 7b). We train all following networks on pairs of: single images + ground-truth masks.

The "C" in each Figure 8 - 10, in "image tile C", stands for current and means that we just use information from the current image. The results of the different architectures are listed in Section 5.3.1.

##### 3.1.1 Simple

This network (Figure 8) consists of five blocks with an additional padding of size 5 is added.

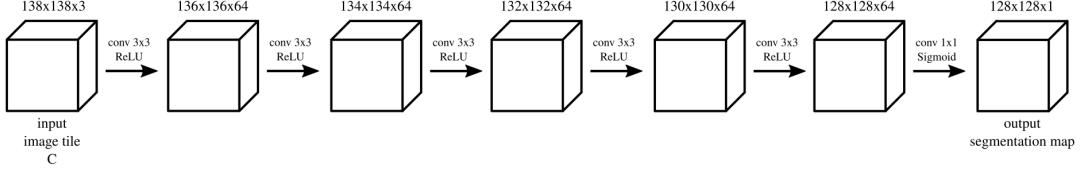


Figure 8: A ”simple”-network with five blocks. An additional padding was added to each  $128 \times 128$  patch because valid-convolution was used.

### 3.1.2 U-Net

This network (Figure 9) is based on the work of Ronneberger et al. [21] and consists of three layers. An additional padding of size 15 is added. This network uses a combination of max-pooling and de-convolution.

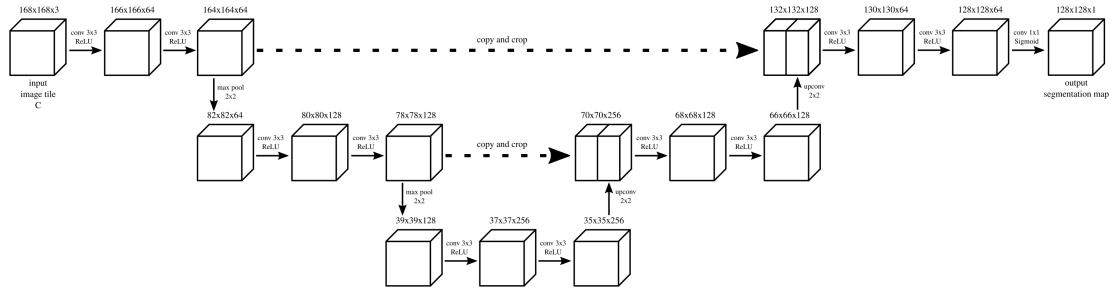


Figure 9: An U-Net-network with three layers. An additional padding was added to each  $128 \times 128$  patch because valid convolution was used.

### 3.1.3 Densenet

This network (Figure 10) is based on the work of Huang et al. [11] and consists of five sub-blocks. Here we do not need an additional padding because of the usage of same-convolutions.

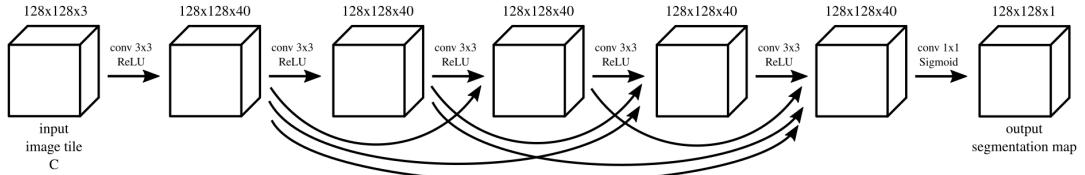


Figure 10: A densenet-network with five sub-blocks, this is later referred to as a single densenet block.

Due to the bad results we achieved with a single input image (see Section 5.3.1), we decided to add additional information to our input, see Section 3.2.

## 3.2 Flow

This section gives an overview of architectures that use additional input information to perform image segmentation (workflow: from Figure 11 and Figure 12a to Figure 12b). We now also use

the previous and next image to perform image segmentation on the current frame, see Figure 11. The "P" and "N" in each Figure 13 - 20, in "image tile P + C + N", stands for previous and next, and means that we use additional information from previous and next image as well.

In Figure 11 we can also see that impurities of the current frame (Figure 11b) are mostly not visible in the previous (Figure 11a) and next frame (Figure 11c). The results of the different architectures are listed in Section 5.3.2.

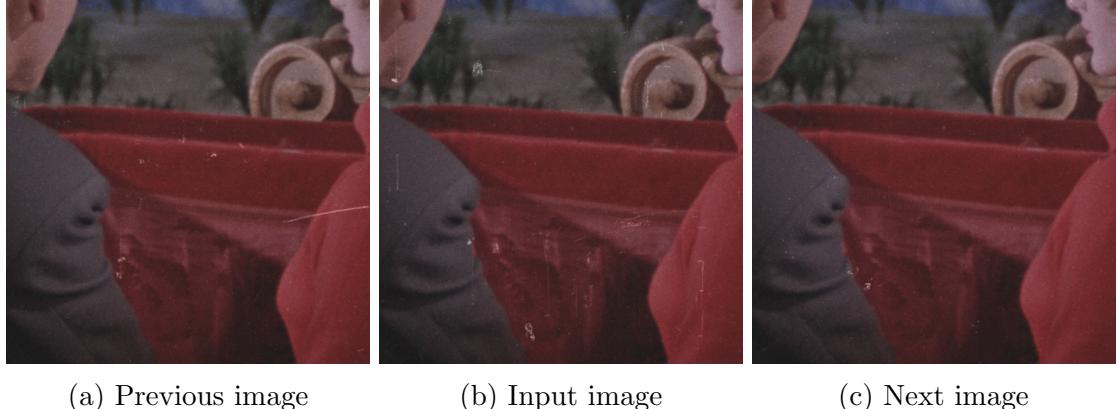


Figure 11: An example of using additional information for image restoration by the usage of previous (a), current input (b) and next image (c).

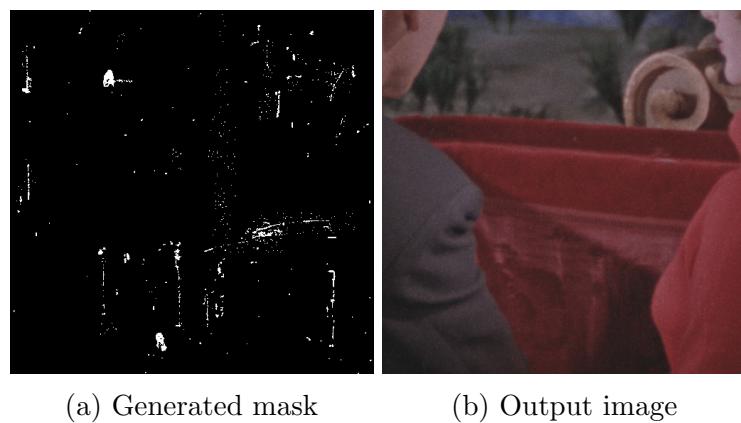


Figure 12: An example workflow of image restoration using additional information (previous, current and next image, *e.g.* Figure 11) for generating the segmentation mask (12a) and inpainting (12b).

### 3.2.1 Early

This network (Figure 13) consists of five blocks with an additional padding of size 5 added. We name this network "early combine" since we combine the information of all images in the beginning.

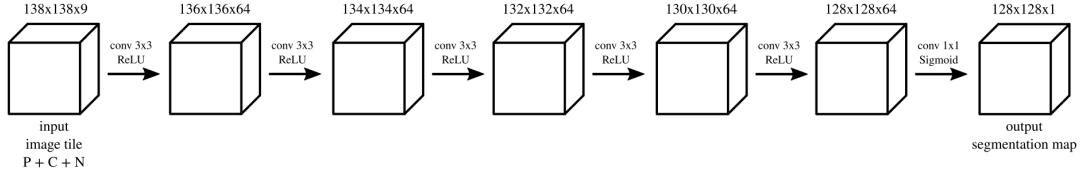


Figure 13: An "early combine"-network with five blocks. An additional padding was added to each  $128 \times 128$  patch because valid-convolution was used.

### 3.2.2 Late

This network (Figure 14) consists of two shared blocks with an additional padding of size 4 added. We name this network "late combine" since we combine the information of all images later. It should be mentioned that the line with the filled circle in Figure 14 means that weights and biases are shared for this block.

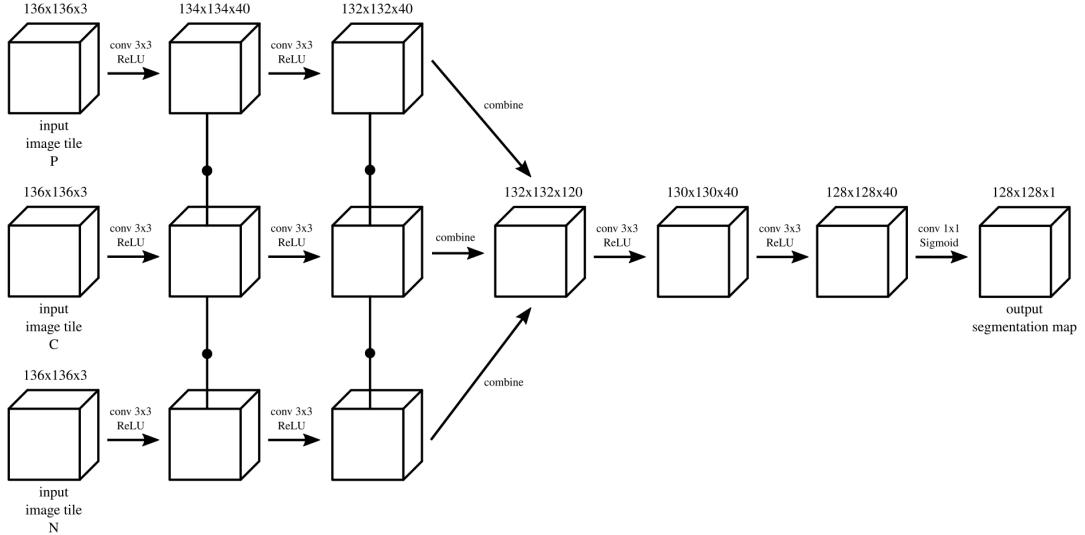


Figure 14: A "late combine"-network with two shared blocks. An additional padding was added to each  $128 \times 128$  patch because valid convolution was used.

### 3.2.3 Late-Late

This network (Figure 15) consists of three shared blocks with an additional padding of size 4 added. We name this network "late late combine" since we combine the information of all images the latest.

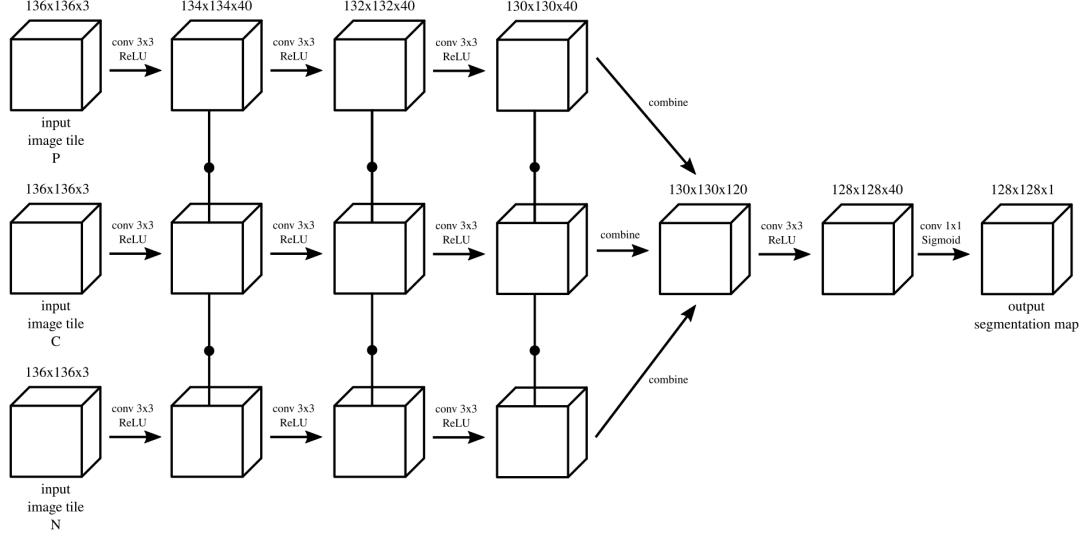


Figure 15: A "late late combine"-network with three shared blocks. An additional padding was added to each  $128 \times 128$  patch because valid convolution was used.

### 3.2.4 U-Net

This network (Figure 16) is based on the work of Ronneberger et al. [21] and consists of three layers. An additional padding of size 15 is added.

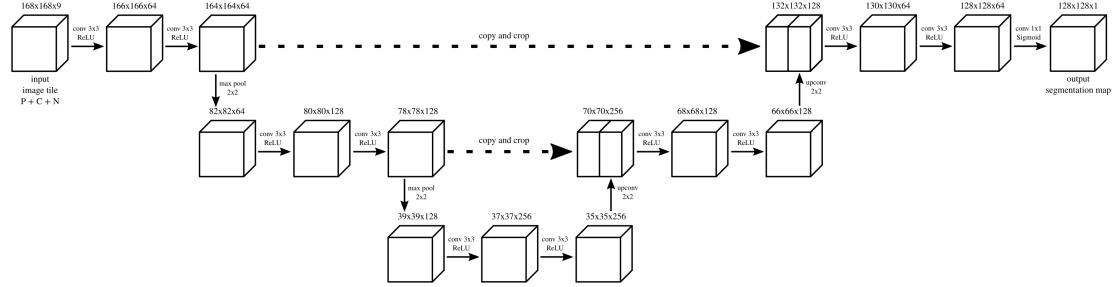


Figure 16: An U-Net-network with three layers. An additional padding was added to each  $128 \times 128$  patch because valid convolution was used.

### 3.2.5 Densenet

This network (Figure 17) is based on the work of Huang et al. [11] and consists of five sub-blocks. Here we do not need an additional padding because of the usage of same-convolutions.

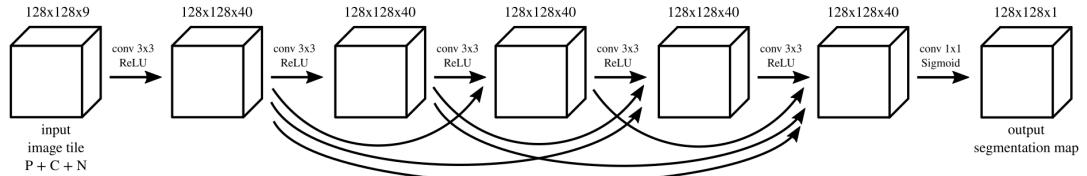


Figure 17: A densenet-network with five sub-blocks.

### 3.2.6 Densenet Combined

This network (Figure 18) is based on the work of Huang et al. [11]. It consists of an upper path and a lower path. The upper path consists of two densenet blocks and the lower path of two shared densenet block. We combine the two paths via concatenating to achieve a final segmentation map. The idea behind this concept is that the upper and lower patch may learn different features that combined achieve better results. Here we do not need an additional padding because of the usage of same-convolutions. Also, we trained this architecture on  $64 \times 64$  sized patches.

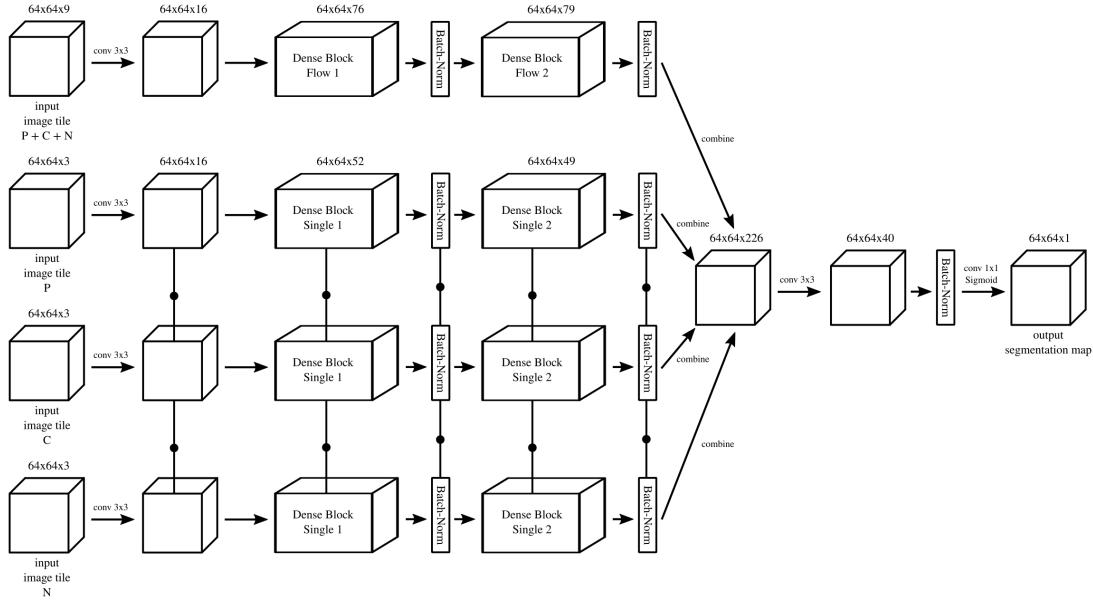


Figure 18: A "densenet combine"-network. The upper path consists of two densenet blocks and the lower path of two shared denesnet block. The two paths get combined via stacking to achieve a final segmentation map.

### 3.2.7 Bottleneck

This network (Figure 19) consists of five blocks with an additional padding of size 5 added. We name this network "bottleneck" because of its structure to reduce the features via a  $1 \times 1$  convolution. This is done to reduce the parameters of the network.

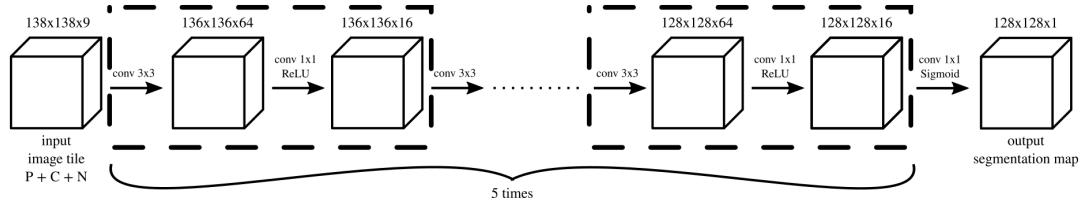


Figure 19: A "bottleneck"-network with five blocks. An additional padding was added to each  $128 \times 128$  patch because valid convolution was used.

### 3.2.8 Densenet Combined 2

This network (Figure 20) is based on the work of Huang et al. [11] and Ronneberger et al. [21]. It consists of an upper path and a lower path. The upper path consists of two densenet blocks and the lower path of an U-Net block with three layers. We subtract the intermediate results of the two paths to achieve a final segmentation map. The idea behind this concept is the same as in Section 3.2.6. Here we do not need an additional padding because of the usage of same-convolutions. Also, we trained this architecture on  $64 \times 64$  sized patches.

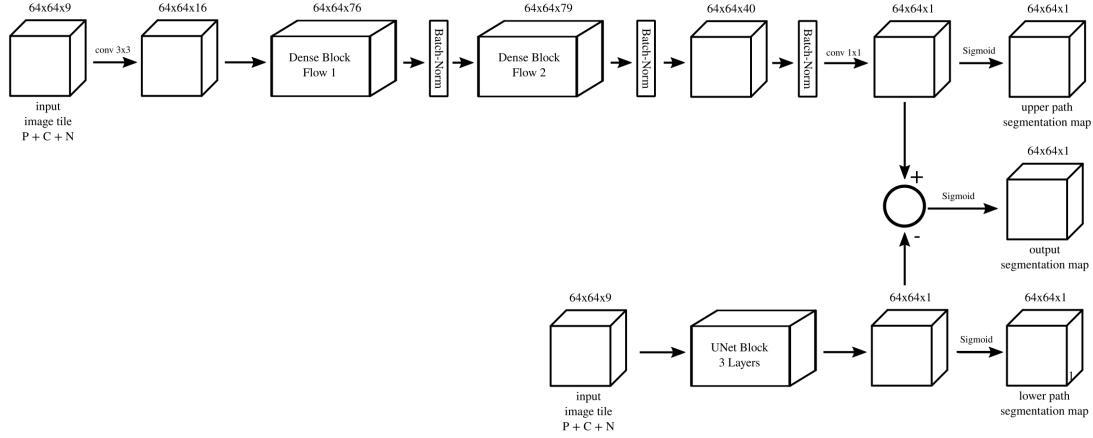


Figure 20: A "densenet combine 2"-network. The upper path consists of two densenet blocks and the lower path of an U-Net block with three layers. The intermediate results of the two paths get subtracted to achieve a final segmentation map.

## 3.3 Regularization Techniques

Techniques to prevent deep learning algorithms from overfitting are called regularizer. The following Sections 3.3.1 - 3.3.5 give a description of the different methods used in our CNN models (Section 3.1, 3.2 and 4.1). It should be mentioned that the following definitions are based on the material from the lectures *Computational Intelligence* [17] and *Neural Networks* [13].

### 3.3.1 Early Stopping

By investigating the trend of test and training error over iterates of deep learning algorithms (Figure 21), one can observe that the training error decreases over iterates, but, the test error first decreases and then increases.

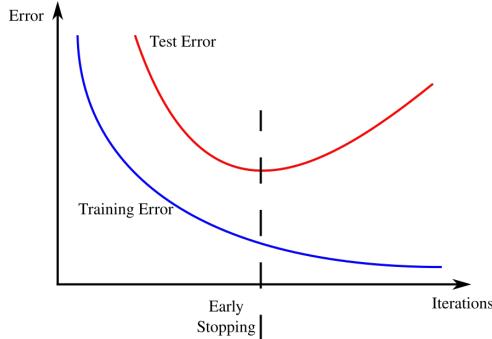


Figure 21: Trend of training and test error over iterates of deep learning algorithms. The training error is increasing. The test error decreases first and then increases.

Early stopping now works as follows: we monitor the error on the test set and store the model parameters every time we get a smaller error, after we determined with the number of iterations, we restore our model parameters with the lowest error on the test set.

Early stopping acts as a regularizer because with increasing iterations the parameters grow and the effective model capacity grows.

### 3.3.2 Dropout

The idea behind is that neurons can not fully rely on the output of other neurons. So, co-specialization is not possible.

With dropout, a neuron/unit is dropped during training with a probability  $1 - p$ , where  $p$  is the so called "keep probability". To drop a unit means, that both in the forward and backward-pass, the unit is ignored (output is set to 0).

For our models we choose a keep probability  $p$  of 0.85.

### 3.3.3 $L_2$ -Regularization (Weight Decay)

$L_2$ -Regularization is a form of adding a parameter norm penalty  $\Omega(\mathbf{W})$  to the error  $E$ , Eq. (1), with  $\mathbf{W}$  being the weights of our model (the biases are not regularized):

$$\tilde{E}(\mathbf{W}, D) = E(\mathbf{W}, D) + \alpha\Omega(\mathbf{W}), \quad (1)$$

with

$$\Omega(\mathbf{W}) = \frac{1}{2}\|\mathbf{w}\|_2^2 \quad (2)$$

and  $\mathbf{w}$  as the vectorized weights. This added  $L_2$ -Regularization , Eq.(2), leads to a weight decay. For the parameter  $\alpha$ , we choose a value of  $10^{-3}$  for our models.

### 3.3.4 Input Normalization

Input normalization is a technique that can speed up learning. We normalize each input feature to have zero-mean and unit-variance, like in Eq. (3):

$$\tilde{x}_i^{(n)} = \frac{x_i^{(n)} - \bar{x}_i}{\sigma_{x_i}}, \quad (3)$$

with  $\bar{x}_i$  as the mean value of the  $i^{th}$ -feature of  $x$ ,  $x_i$ ,

$$\bar{x}_i = \frac{1}{N} \sum_{j=1}^N x_i^{(j)}, \quad (4)$$

and  $\sigma_{x_i}$  as the standard deviation of the  $i^{th}$ -feature of  $x$ ,  $x_i$ ,

$$\sigma_{x_i} = \sqrt{\frac{1}{N} \sum_{j=1}^N (x_i^{(j)} - \bar{x}_i)^2} \quad (5)$$

### 3.3.5 Batch Normalization

The output of a layer can be interpreted as the input of the next layer. In Batch normalization, this normalization (like in Section 3.3.4) is done to the output of each layer for the current batch.

This techniques leads to better learning dynamics. Batch normalization can be seen as an intermediate layer that linearly transforms the output of a layer.

## 4 Image In-painting

In this Section follows a description of our methods and approaches for the removing of impurities by image in-painting. Again, we focus on machine learning approaches namely CNNs but compare these against a diffusion based image in-painting. Section 4.1 gives an overview CNN architectures that are trained for a single color channel. Section 4.2 gives an theoretical overview on diffusion based image in-painting, as well as the exact used approach. The results of the different image in-painting approaches can be seen in Section 5.4.

### 4.1 CNN

Based on the results of Section 5.3.2 we decide to use an U-Net approach and a densenet approach. Additionally, as seen in Section 5.3.2, we use information from previous and next images as well. We train all following networks on quintuples of: previous - current - next images - masks + ground-truth in-painting.

It should be mentioned that we use ReLU activation functions between layers and a sigmoid output activation function to perform image in-painting. We also apply regularization techniques, namely early stopping (Section 3.3.1), dropout (Section 3.3.2) and  $L_2$ -regularization (Section 3.3.3). The training patches we use are of size  $128 \times 128$ . We choose the  $L_2$ -loss for computing the loss/error between predicted and ground-truth in-painting but just over masked pixels. And, we copy valid input pixel values into the predicted in-painting.

Since we train on a single color channel, we also have to perform the testing on each color channel separately and combine the intermediate results.

#### 4.1.1 U-Net

This network (Figure 22) is based on the work of Ronneberger et al. [21] and consists of three layers.

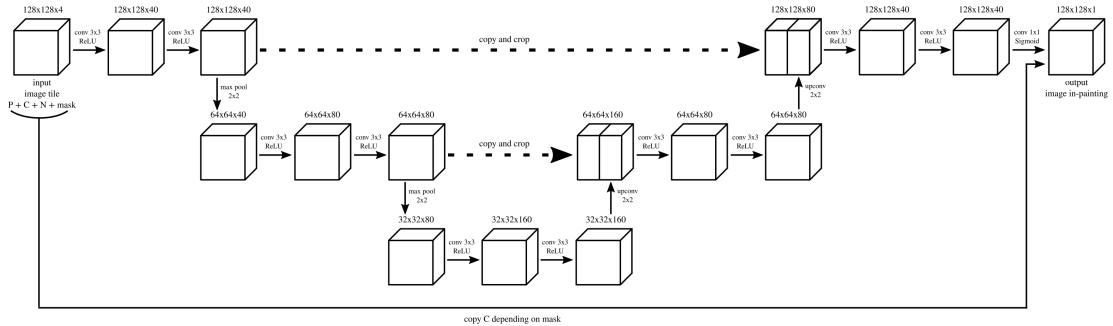


Figure 22: An U-Net-network for image in-painting with three layers.

#### 4.1.2 Densenet

This network (Figure 23) is based on the work of Huang et al. [11] and consists of five sub-blocks.

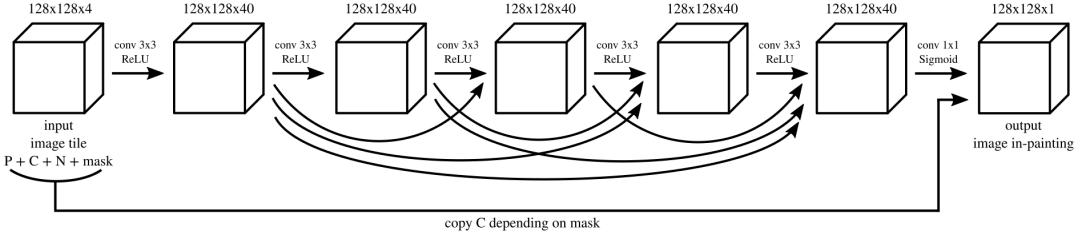


Figure 23: A densenet-network for image in-painting with five sub-blocks.

## 4.2 Perona-Malik In-painting

The second image in-painting approach is based on a modified version of coherence-enhancing image diffusion, which was an assignment in the course *Bildverarbeitung und Mustererkennung* [19]. The theoretical basis are after Perona and Malik [18] or Weickert [25]. Now follows a description of the used method for one color channel.

Let  $u : \Omega \rightarrow \mathbb{R}$  be the image function, with  $\Omega \subset \mathbb{R}^2$  the image domain. The basic diffusion equation is given by:

$$\frac{\partial u}{\partial t} = \operatorname{div}(\nabla u). \quad (6)$$

Eq. (6) is a PDE and is used to describe *isotropic* diffusion. To increase the modeling accuracy, a so-called diffusion tensor  $D$  is introduced into Eq. (6):

$$\frac{\partial u}{\partial t} = \operatorname{div}(D \nabla u), \quad (7)$$

Eq. (7) describes *anisotropic* diffusion and the diffusion tensor  $D$  models rate of diffusion in different directions. If  $D$  is the identity matrix, we again obtain *isotropic* diffusion.

Coherence-enhancing diffusion leads the diffusion along coherent regions, similar structured regions, in our case edge directions. We achieve this, by using a structure tensor  $S$  ( $2 \times 2$  matrix). The structure tensor is able to give for every pixel information about its surrounding structure, and is calculated via:

$$S = G_\sigma * \begin{bmatrix} \tilde{u}_x^2 & \tilde{u}_x \tilde{u}_y \\ \tilde{u}_x \tilde{u}_y & \tilde{u}_y^2 \end{bmatrix}, \quad (8)$$

with subscripts denoting directional derivatives.  $G_\sigma$  is a Gaussian smoothing kernel with  $\sigma_t$ , and  $\tilde{u}$  is a Gauss-filtered version of  $u$  with  $\sigma_g$ . The eigenvalues  $\mu_1, \mu_2$  and eigenvectors  $\nu_1, \nu_2$  of the structure tensor  $S$  give information about the local structure of the image, with the eigenvalues determining the magnitude of the grayvalue variation, the corresponding eigenvector giving the direction of the grayvalue variation:

1.  $\mu_1, \mu_2$  both small: flat region,
2.  $\mu_1 \gg \mu_2$  or vice-versa: strong grayvalue variation in one direction, *i.e.* image edge,
3.  $\mu_1, \mu_2$  both large: strong grayvalue variation in two directions, *i.e.* image corner,

with this information we will guide the diffusion process.

The diffusion tensor  $D$  is calculated as

$$D = [\nu_1 \quad \nu_2] \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} \nu_1^T \\ \nu_2^T \end{bmatrix}, \quad (9)$$

using a small constant  $\alpha > 0$ ,  $\lambda_1$  and  $\lambda_2$  are given as

$$\begin{cases} \lambda_1 = \alpha \\ \lambda_2 = \alpha + (1 - \alpha)(1 - g(|\mu_1 - \mu_2|)) \end{cases} \quad (10)$$

and  $g(s) = e^{-\frac{s^2}{2\gamma^2}}$  with a parameter  $\gamma$ .

We solve Eq. (7) with a semi-implicit approach. Since we are in the discrete world, an image of size  $M \times N$  is represented as a vector  $U \in \mathbb{R}^{MN}$ , the time derivative  $\frac{\partial u}{\partial t}$  and  $K \in \mathbb{R}^{2MN \times MN}$  the equivalent to  $\nabla$  approximated by finite differences,  $D(U^t) \in \mathbb{R}^{2MN \times 2MN}$  the diffusion tensor,  $\tau$  the time discretisation step and since  $\text{div} = -\nabla^T$  we then can write Eq. (7) as

$$\frac{U^{t+1} - U^t}{\tau} = -K^T D(U^t) K U^{t+1}. \quad (11)$$

Re-writting Eq. (11) and the introduction of a diagonal mask matrix  $Q \in \mathbb{R}^{MN \times MN}$  that highlights pixels for in-painting, and  $\hat{Q}$  the inverse of  $Q$  for highlighting preserved regions and solving for  $U^{t+1}$  leads to the image in-painting equation for one channel:

$$\begin{aligned} Q(U^{t+1} - U^t) + \hat{Q}\tau K^T D(U^t) K U^{t+1} &= 0 \\ Q U^{t+1} - Q U^t + \hat{Q}\tau K^T D(U^t) K U^{t+1} &= 0 \\ (Q + \hat{Q}\tau K^T D(U^t) K) U^{t+1} &= Q U^t \\ U^{t+1} &= (Q + \hat{Q}\tau K^T D(U^t) K)^{-1} Q U^t \end{aligned} \quad (12)$$

The extension of this single channel image in-painting to color images is done by applying the diffusion for each channel. To avoid artifacts many different approaches exist. As Åström et al. [1] state one could use a different color space, like HSV instead of RGB, use one diffusion tensor for each channel that is computed as the weighted average sum of the structure tensors for each color channel, use a different structure tensor for each channel, etc. We choose the approach of using the weighted sum over all channels for the structure tensor and use this for each channel.

## 5 Evaluation

We evaluate our different approaches with ground-truth images from HS-Art [10]. A description of this dataset follows in Section 5.1. The needed implementation details for the evaluation are stated in Section 5.2 and in Section 5.3 our results are presented.

### 5.1 Dataset

HS-Art [10] provides us with 96 image sequences. Per sequence we get an input image (current image) that is corrupted, the previous image in the sequence, the next image in the sequence, the ground-truth segmentation map, highlighting impurities in the current image as well as ground truth image in-painting results of the current image. Each image is of size  $512 \times 512$ .

For training and evaluation of our CNN approaches we split the sequences into 80% training sequences and 20% test sequences. The training is done on patches of size  $128 \times 128$ , unless otherwise stated, and we train on patches were atleast 1% and at most 80% pixels are impure.

For the image segmentation task we further categorize our test set into sets that have: (i) ground-truth masks with large areas, (ii) ground-truth masks with small areas, (iii) sequences with large flow, and (iv) sequences with small flow.

We also perform dataset augmentation to get more data for training. Here we perform mirroring of the images and masks as well as a flow reversal by exchanging previous and next image.

## 5.2 Implementation Details

The different parameters used for our CNN approaches are listed either in Section 3.3 or in the corresponding model definitions. Additionally, we use the ADAM-optimizer with a learning rate of  $10^{-3}$  to train other models.

For the Perona-Malik image in-painting approach (Section 4.2) we use the following values for the parameters:  $\alpha = 0.005$ ,  $\gamma = 0.001$ ,  $\tau = 5$ , end time = 15,  $\sigma_g = 0.7$  and  $\sigma_t = 1.5$ .

The project was realized in python and we used the Tensorflow API to implement our CNN approaches. The training and testing was done on a NVidia Geforce Titan X (GM200 series).

## 5.3 Results for Image Segmentation

In the following Sections we list the results of image segmentation using a single image (Section 5.3.1) and using additional information (Section 5.3.2).

### 5.3.1 Single Image

The results in Figure 24 suggest that the information given in just a single current image is not enough for a sufficient learning and estimating of the segmentation mask. Therefore, we add additional information by using the previous and next image frames for all additional approaches.

### 5.3.2 Flow

Following prior results, we decided to not evaluate the "late late combine"-architecture (Section 3.2.3) and the UNet-architecture (Section 3.2.4) but rather work with the other approaches from Section 3.2.

In Figure 25 the accuracy results of the different architectures on the test set are illustrated. One can see that the Densenet Combined model performs best, closely followed by the Densenet, Bottleneck and Densenet Combined 2 approach. The Early and Late models perform the worst.

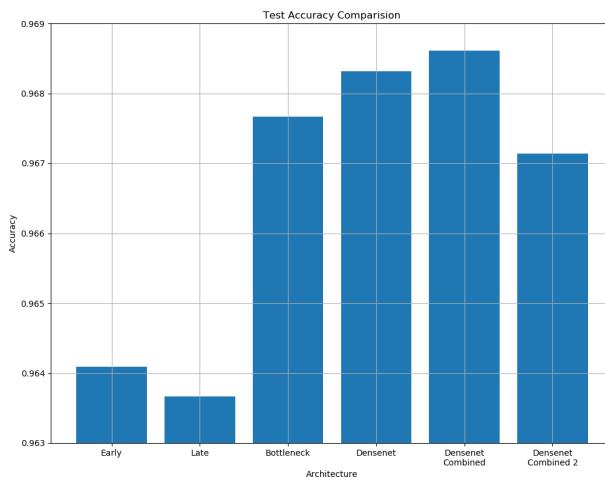


Figure 25: Accuracy for the different image segmentation architectures, evaluated on the test set. For the accuracy higher values mean better results. The best results are achieved with the Densenet Combined approach (Section 3.2.6).

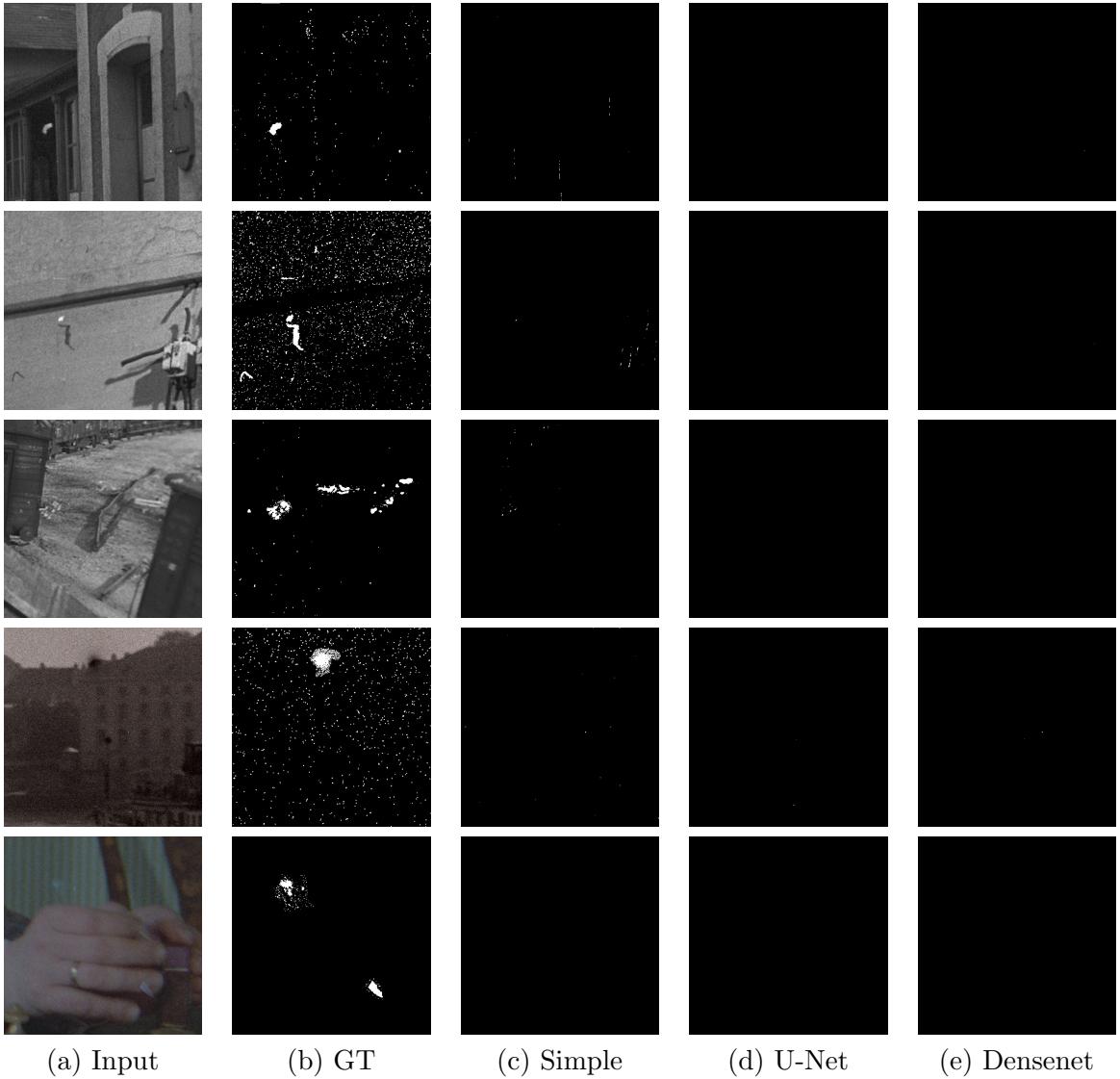


Figure 24: Sample results of single image segmentation for Simple (Section 3.1.1), U-Net (Section 3.1.2) and Densenet (Section 3.1.3).

Figure 26 shows the accuracy of the different architectures on the sub-categorized test sets: large area (Figure 26a), small area (Figure 26b), large flow (Figure 26c) and small flow (Figure 26d).

The overall best architecture, Densenet Combined, also performs well on each of the single sub-test sets. The Bottleneck model also performs well on each of these subsets. The overall poor performing Early model achieves great results on the small area test set. Whereas the Late model performs on each single test set poorly.

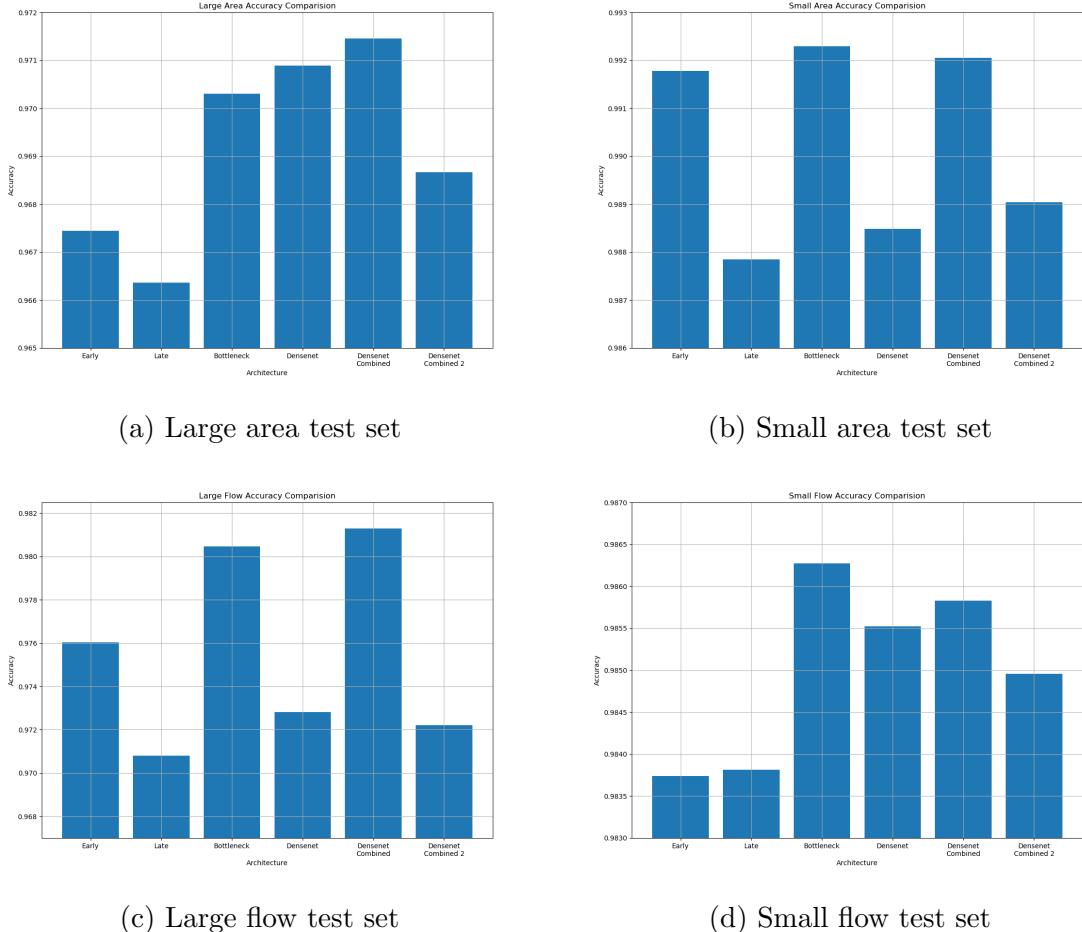
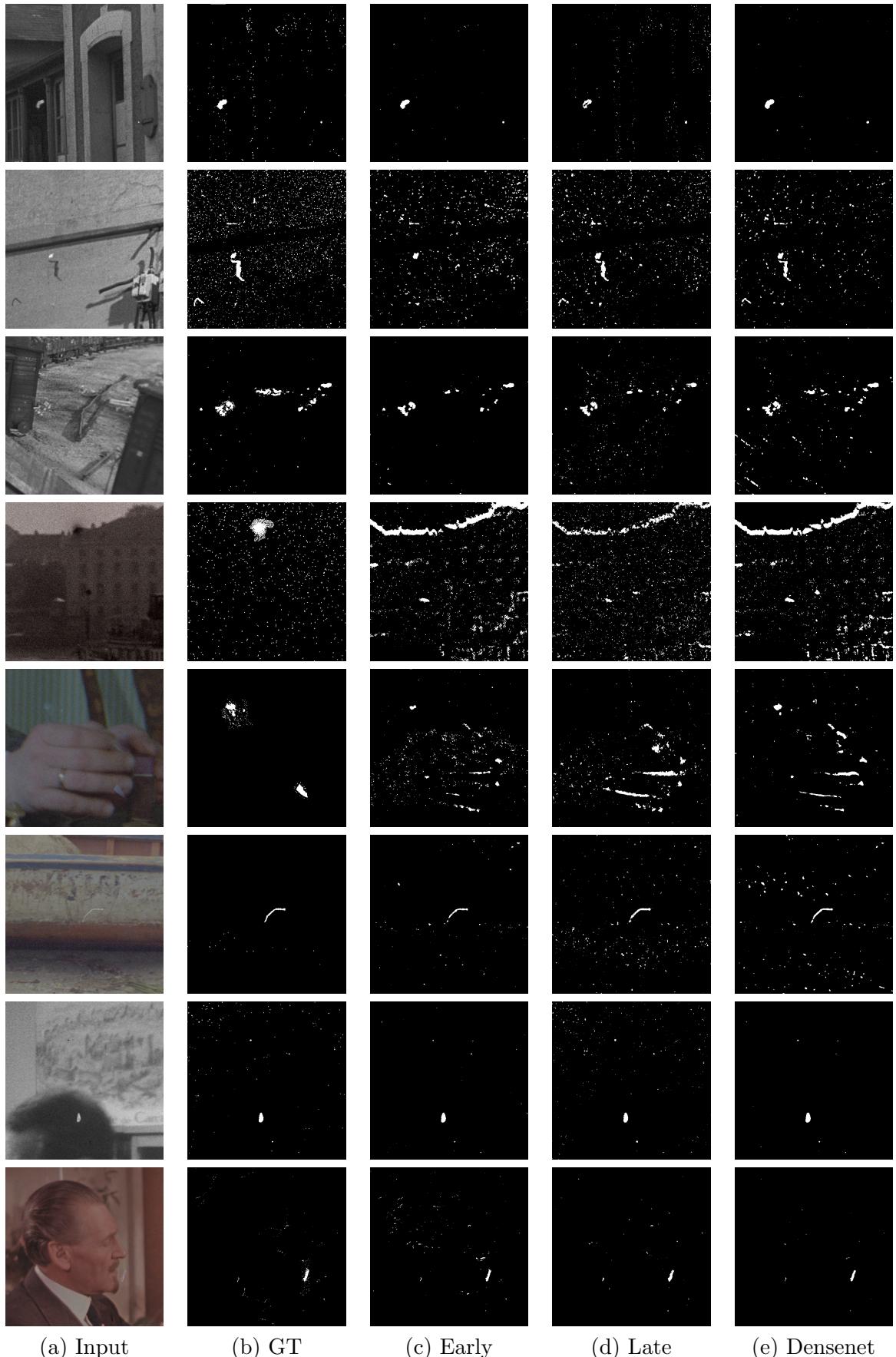


Figure 26: Accuracy for the different image segmentation architectures, evaluated on the sub-categorized test sets.

Figures 27 and 28 give an overview of estimated segmentation maps of the different architectures compared to the ground-truth masks.



(a) Input

(b) GT

(c) Early

(d) Late

(e) Densenet

Figure 27: Sample results 1 of flow image segmentation for Early (Section 3.2.1), Late (Section 3.2.2) and Densenet (Section 3.2.5).

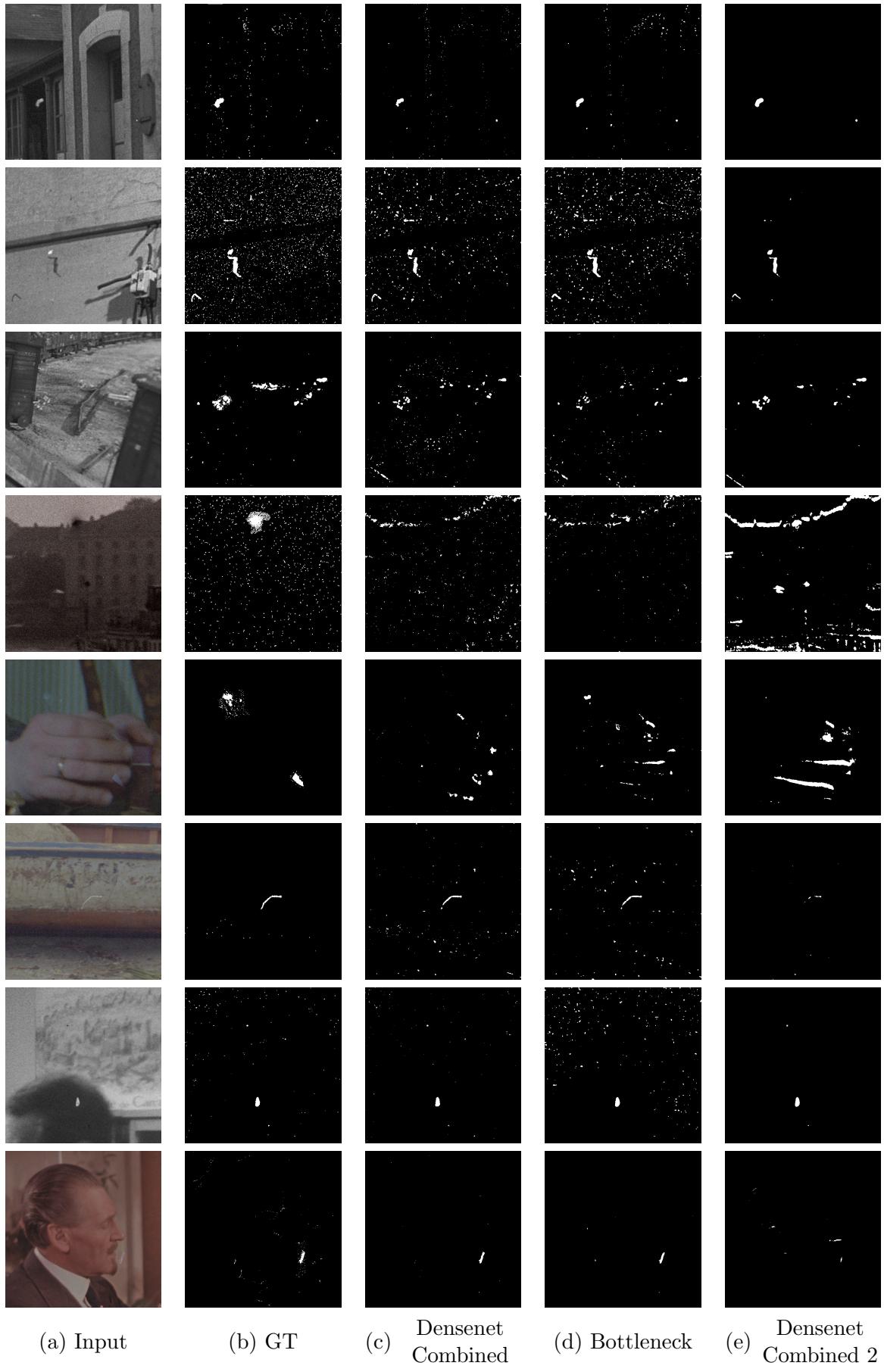


Figure 28: Sample results 2 of flow image segmentation for Densenet Combined (Section 3.2.6), Bottleneck (Section 3.2.7) and Densenet Combined 2 (Section 3.2.8).

## 5.4 Results for Image In-painting

Figure 29 gives an overview of the results of different image in-painting approaches. One can see that the Densenet approach achieves the visual most appealing image in-painting results followed by the Perona-Malik method. The U-Net approach has the visual worst results.

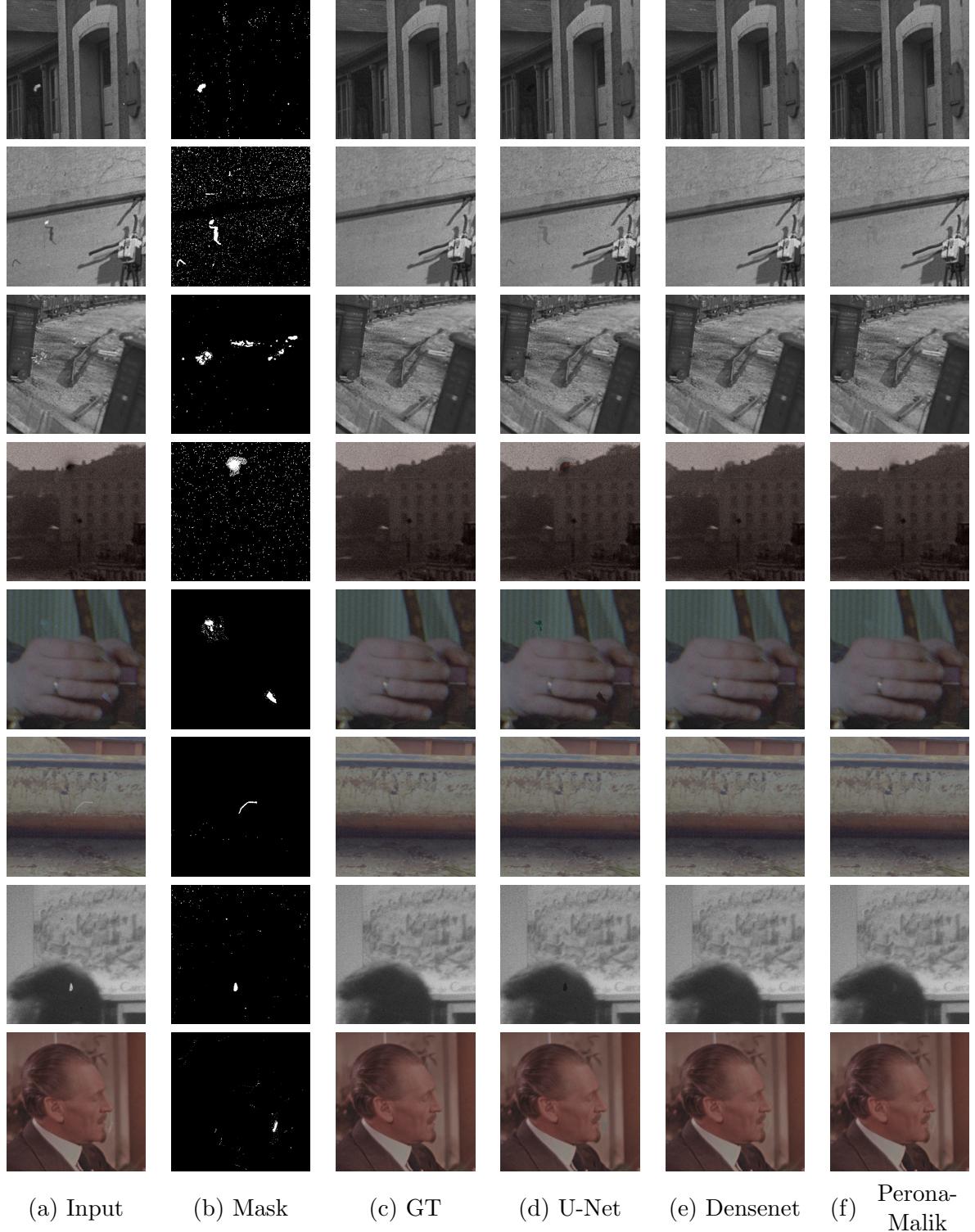


Figure 29: Sample results of image in-painting for U-Net (Section 4.1.1), Densenet (Section 4.1.2) and Perona-Malik (Section 4.2).

## 6 Conclusion

We have introduced an deep learning method to perform image restoration. We showed that a two part approach, split into the detection of the impurities (via image segmentation, Section 3) and the correction of them (via image in-painting, Section 4), is able to perform visually appealing end results.

The impurity detection based on image segmentation works best with the Densenet Combined model (Section 3.2.6) *w.r.t* the achieved accuracy score on the test set. By investigating segmentation map results (Figure 27 and 28) one can see that for sequences with high flow the resulting segmentation map is quite wrong compared to the ground-truth map. This may be improved by taking additional flow information into account during training. Also the segmentation of using just a single input image does not work at all. Adding additional information from previous and next images improves the results, since most of the time the impurities on the current frame are not visible in the previous and next image (Figure 11).

The correction of the impurities based on image in-painting works visually best with the Densenet model (Section 4.1.2). Other approaches, like *e.g.* an encoder-decoder scheme or an approach based on GANs may achieve better results. The Perona-Malik approach may be improved by finding better values for its parameters but as a relative efficient model, it works quite well.

In general a higher number of training and test data would improve the results of our CNN approaches. Supervised deep learning approaches need a huge amount of training data to perform well. Also, some of the provided ground-truth masks are of poor quality, which makes it hard for the networks to learn sufficiently.

## References

- [1] F. Åström, M. Felsberg, and R. Lenz. Color persistent anisotropic diffusion of images. In *Proceeding of the Scandinavian Conference on Image Analysis (SCIA 2011)*, pages 262–272. Springer, 2011. [16](#)
- [2] S. Boujena, K. Bellaj, E. El Guarmah, and O. Gouasnouane. An improved nonlinear model for image inpainting. *Applied Mathematical Sciences*, 9(124):6189–6205, 2015. [6](#)
- [3] C. Burlin, Y. Le Calonnec, and L. Duperier. Deep image inpainting. In *Stanford C231N course project*, 2017. [6](#)
- [4] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–848, April 2018. [5](#)
- [5] Y. Chen and T. Pock. Trainable nonlinear reaction diffusion: A flexible framework for fast and effective image restoration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1256–1272, 2017. [5](#)
- [6] Y Chen, W. Yu, and T. Pock. On learning optimized reaction diffusion processes for effective image restoration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2015)*, 2015. [5](#)
- [7] P. Getreuer. Total Variation Inpainting using Split Bregman. *Image Processing On Line*, 2:147–157, 2012. [6](#)
- [8] C. Guillemot and O. Le Meur. Image inpainting: Overview and recent advances. *IEEE Signal Processing Magazine*, 31(1):127–144, 2014. [6](#)
- [9] K. Gupta, S. Kazi, and T. Kong. Deeppaint: A tool for image inpainting. In *Stanford C231N course project*, 2016. [6](#)
- [10] HS-Art. We are the specialist for digital film restoration solutions. <http://www.hs-art.com/index.php>. [Online; accessed 24<sup>th</sup> September 2018]. [ii](#), [2](#), [16](#)
- [11] G. Huang, Z. Liu, and K. Q. Weinberger. Densely connected convolutional networks. *CoRR*, abs/1608.06993, 2016. [5](#), [7](#), [10](#), [11](#), [12](#), [14](#)
- [12] R. Köhler, C. Schuler, B. Schölkopf, and S. Harmeling. Mask-Specific Inpainting with Deep Neural Networks. In *Proceedings of the German Conference on Pattern Recognition (GCPR 2014)*, 2014. [6](#)
- [13] R. Legenstein. Neural networks. University Lecture, 2017. [12](#)
- [14] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2015)*, 2015. [5](#)
- [15] P. Milanfar. A tour of modern image filtering: New insights and methods, both practical and theoretical. *IEEE Signal Processing Magazine*, 30(1):106–128, 2013. [5](#)
- [16] M. Neri and E. R. R. Zara. Total variation-based image inpainting and denoising using a primal-dual active set method. *Philippine Science Letters*, 7(1):97–103, 2014. [6](#)
- [17] F. Pernkopf, G. Bellec, A. Subramoney, and C. Knoll. Computational intelligence. University Lecture, 2016. [12](#)

- [18] P. Perona and J. Malik. Scale-space and edge detection using anisotropic diffusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(7):629–639, 1990. [6](#), [15](#)
- [19] T. Pock. Bildverarbeitung und mustererkennung. University Lecture, 2016. [15](#)
- [20] K. Qazanfari, R. Aslanzadeh, and M. Rahmati. An efficient evolutionary based method for image segmentation. *CoRR*, abs/1709.04393, 2017. [5](#)
- [21] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI 2015)*, pages 234–241, 2015. [5](#), [7](#), [10](#), [12](#), [14](#)
- [22] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. [5](#)
- [23] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2015)*, June 2015. [5](#)
- [24] N. van Noord and E. O. Postma. Light-weight pixel context encoders for image inpainting. *CoRR*, abs/1801.05585, 2018. [6](#)
- [25] J. Weickert. *Anisotropic diffusion in image processing*, volume 1. Teubner Stuttgart, 1998. [6](#), [15](#)
- [26] R. A. Yeh, C. Chen, T. Y. Lim, A. G. Schwing, M. Hasegawa-Johnson, and M. N. Do. Semantic image inpainting with deep generative models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2017)*, pages 5485–5493, 2017. [6](#)
- [27] S. Yuheng and Y. Hao. Image segmentation algorithms overview. *CoRR*, abs/1707.02051, 2017. [5](#)