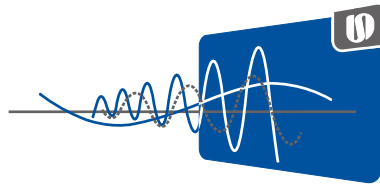


Universität Stuttgart



Universität Stuttgart
Institut für Statik und Dynamik
der Luft- und Raumfahrtkonstruktionen
Prof. Dr.-Ing. Tim Ricken



Near Realtime Damage Detection of Wooden Planks

Lena Scholz, Samuel Lehmköster
& Julian Franquinet

Supervisors:
André Mielke M. Sc.,
Luis Mandl M. Sc.

Pfaffenwaldring 27, 70569 Stuttgart

Report

September 27, 2021

Contents

1	Introduction	1
1.1	Description of the Task	1
1.2	Approach to the Solution	1
2	Pre-trained Encoders with a Simple Decoder	2
2.1	VGG16	3
2.2	MobileNetV2	3
2.3	ResNet50V2	4
3	Advanced Net Architectures and Decoder	4
3.1	SegNet Architecture	5
3.2	U-Net Architecture	5
3.3	Advanced Decoders	6
4	Final Segmentation Architecture with a Pre-trained Encoder	7
5	Evaluation of the Results	8
5.1	Pre-trained Encoders with a Simple Decoder	8
5.2	SegNet	12
5.3	U-Net	14
5.4	Advanced Decoder	19
5.5	The Final Model	23
5.6	Benchmark Studies on a Coral Chip	24
6	Summary & Outlook	26
	Appendices	28
A	Parameter Summary of All Architectures	28
B	Algorithms	29

List of Figures

1	Example of a white coated wooden plank with an error on the top right corner together with its corresponding mask	1
2	Example of a black coated wooden plank with a bar-code label together with its corresponding mask	1
3	Architecture of simple decoder	3
4	Architecture of VGG16	3
5	Architecture of MobileNetV2	4
6	Architecture of ResNet50V2	4
7	Architecture of simplified SegNet	5
8	Architecture of U-Net with complexity of 4	6
9	Architecture of the Advanced Decoder	7
10	Architecture of Final Model	7
11	Losses of different pre-trained encoders for single-classification . .	8
12	Losses of different pre-trained encoders for multi-classification . .	9
13	Losses of different pre-trained encoder for single-classification . .	10
14	Predictions for different encoders for single-classification	10
15	Losses of different pre-trained encoder for multi-classification . .	11
16	Predictions for different encoders for multi-classification	11
17	Predictions of simplified SegNet for ReLU- and SELU-activation	12
18	Loss of simplified SegNet for ReLU- and SELU-activation	13
19	Loss of simplified SegNet for single- and multi-classification . . .	13
20	Losses of U-Net with different complexities using sliced images for single-classification	14
21	Losses of U-Net with different complexities using big images for single-classification	15
22	Predictions of U-Net with different complexities	16
23	Loss of U-Net with complexity of 5 for multi-classification	17
24	Predictions of U-Net with different sizes of images	17
25	Metric of U-Net with complexity of 5 for multi-classification with different losses	18
26	Predictions of U-Net with complexity of 5 for multi-classification with different losses	19
27	Loss for a decoder with Conv2DTranspose layers for multi- & single-classification	20
28	Loss for different decoder with added layers for multi-classification	21
29	Predictions for different decoder with added layers for multi-classification	21
30	Loss for different decoder with added layers for single-classification	22
31	Predictions for different decoder with added layers for single-classification	22
32	Losses of final network for single- & multi-classification	23
33	Fine tuning of final network for single- & multi-classification . . .	23
34	Predictions of final network	24

35	Benchmark study for the Pre-trained models, the U-Net and the SegNet	25
36	Benchmark study for different inputs while using a U-Net	25

List of Algorithms

1	Simplified Decoder	29
2	Jaccard Distance	29
3	Dice Metric	29
4	Simplified SegNet	30
5	U-Net	31
6	Advanced Decoder	32
7	Advanced Decoder with Dropout	33
8	Advanced Decoder with Upsampling	34
9	Advanced Decoder with Transposed Layer	35
10	Final combined version of U-Net with the MobileNetV2	36

Abbreviation

Abbreviation	Description
CNN	Convolutional Neural Network
Param.	Parameter
nS	Trainingsset without shift
wS	Trainingsset with shift
ES	Early stopping method applied
MC	Multi-Classification
SC	Single-Classification
BN	Batch normalization function applied
ReLU	Rectified linear unit
SELU	Scaled exponential linear unit

1 Introduction

1.1 Description of the Task

In this report the segmentation of wooden planks with images of 256x3072 pixels are investigated for a real-time analysis task. The motivation for this report is the creation of a complete hardware and software system for a real-time damage detection on coated wooden planks in a fully automatized factory for built-in kitchen furniture.

The task is to detect damages due to a coating process on wooden planks. In order to be able to adjust the sensitivity of the detection, it is divided into two main tasks. The first task is the segmentation of the image in order to find the wooden plank inside the image. The second task is the actual detecting of the damage. This should be easily accomplished by a simple threshold. Therefore, this work will focus on the more complicated part, creating and training a Convolutional Neural Network (CNN) for the purpose of the segmentation. For the training, real example pictures (1.1) of automatically coated wooden planks were used. The Segmentation Network should find these planks in a picture section and distinguish them from the conveyor belt system in the background.



Figure 1: Example of a white coated wooden plank with an error on the top right corner together with its corresponding mask

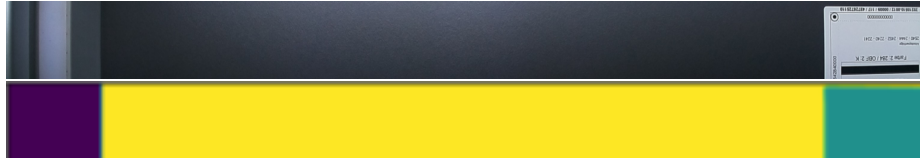


Figure 2: Example of a black coated wooden plank with a bar-code label together with its corresponding mask

1.2 Approach to the Solution

For the approach to the solution of this problem, at first, a pre-trained encoder with a small decoder as segmentaion network was applied. Therefore, several pre-trained encoders of the **Tensorflow** library were investigated and later on different decoder structures tested. Additional to this simple Fast-forward CNNs [4] smaller, self-edited codes for more complex Segmentation Network Architectures were researched, like the U-Net[5] and the SegNet[1].

The original example set contains 126 pictures of the size 256x3072, which is further referred to as 'big images'. In a data generation process, the images were flipped and shifted, in order to generate more data for the training. Because of the usage of pre-trained encoders, which have a build-in input (no further layers were included on the input-side of the encoders), the original image (1.1) with a size of 256x3072 pixels is split into 12 images of 256x256 pixels, further referred to as 'sliced images'.

Since some of the planks have bar code labels, two setups are considered. At first we only aim to distinguish planks from the background (single-class setup). In order to get rid of the confusion around labels we separate planks, labels, and background (multi-class setup).

In all test setups, unless stated otherwise, `BinaryCrossentropy` is used as loss function for single-classification while `SparseCategoricalCrossentropy` is applied for the multi-classification. The used optimizer for the compiling of the models is the `adam` optimizer, based on stochastic gradient descent scheme. The learning rate is not set for in any of these investigations, so it is always set to 0.001 as a standard setting in Keras (<https://keras.io/api/optimizers/adam/>). As activation function `SELU` was applied in nearly all of the layers, except when different activations should be tested (as `ReLU` in section 5.2).

The amount of total parameters and trainable parameters for each model trained in this project can be found in the appendix (appendix A).

2 Pre-trained Encoders with a Simple Decoder

In this first part of the project the pre-trained encoders of the `VGG16`, the `MobileNetV2` and the `ResNet50V2` were chosen out of the available models in the keras software (<https://keras.io/api/applications>). These CNNs were selected due to their broad use in other open projects and their further development. Furthermore, the encoder of the three CNN differ a lot in the amount of weights (see appendix A and it could thereby be evaluated if deeper encoder and therefore deeper CNN perform better in this task. In order to compare these different encoders, a simple decoder (code found in alg. 1) as shown in fig. 3 was applied to each of the encoders in order to yield a CNN for the first tests in this project.

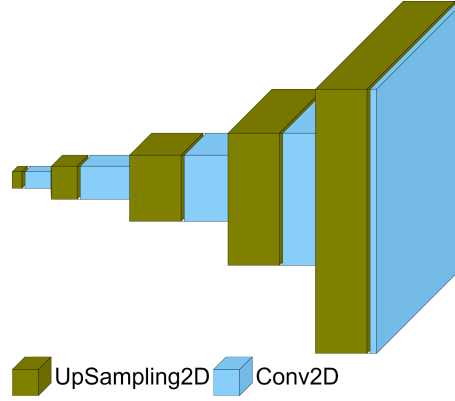


Figure 3: Architecture of simple decoder

2.1 VGG16

This CNN was build for the ImageNet Challenge 2014 for image classification and localisation [8]. This is why it is the the oldest architecture. The objective was to develop deeper CNN architectures, here with 23 weigh layer and over 14.700.000 parameters. Its size lays between the `ResNet50V2` and the `MobileNetV2` regarding the amount of parameters.

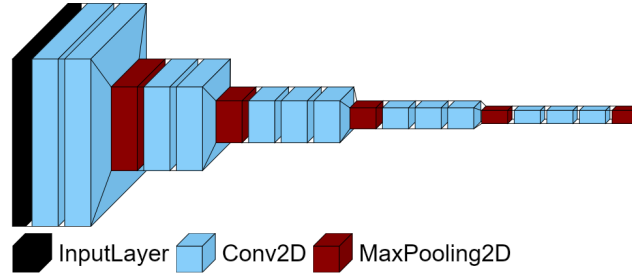


Figure 4: Architecture of VGG16

2.2 MobileNetV2

The `MobileNetV2`, introduced in [6], was build and improved in order to obtain a small mobile model for semantic segmentation. It is the smallest pre-trained model and also the latest model used in this project. It contains around 2.400.000 parameters in 88 layers.

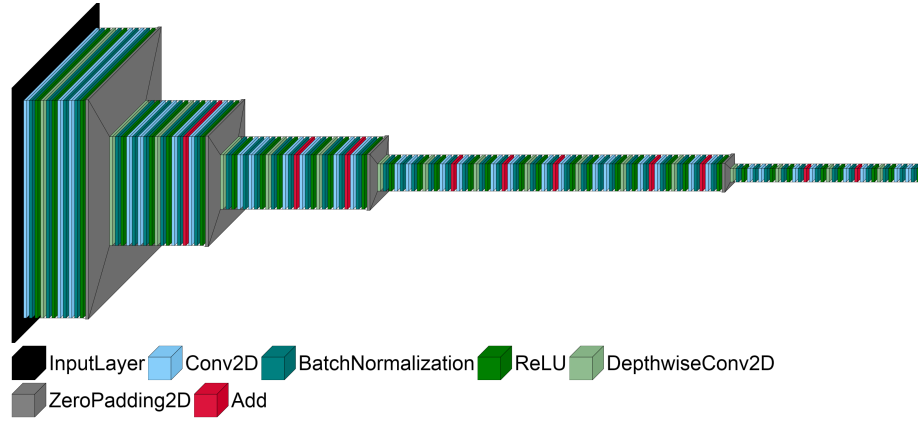


Figure 5: Architecture of MobileNetV2

2.3 ResNet50V2

The ResNet50V2 architecture was built in order to ease the training of deep CNN (see [2]) by using layers as residual functions regarding the input. The ResNet50V2 is the biggest segmentation net encoder used in this project regarding the amount of parameters. It contains around 23.700.000 parameters.

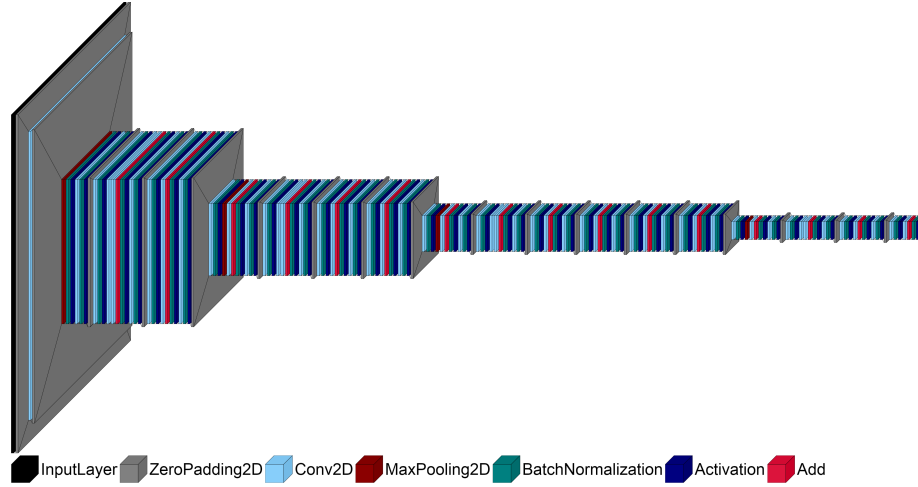


Figure 6: Architecture of ResNet50V2

3 Advanced Net Architectures and Decoder

In this second part of the project, different architectures of segmentation nets are tested. Because of the different difficulties in every architecture, not all CNN

performed the same test cases. So, several different settings in each investigation are applied.

3.1 SegNet Architecture

The SegNet Architecture in [1] was created for a good segmentation performance while using fewer parameters. This is achieved by using max-pooling indices from the encoder for the up-sampling in the decoder. It was motivated for scene understanding applications.

In order to accelerate the training of the SegNet, a reduced architecture with fewer parameters than in the setup of [1] is considered (see fig. 7). The SegNet architecture can be extended by adding more blocks of `Conv2D`, `BatchNormalization` and `Activation` before the `MaxPoolingWithIndices` is performed. The presented implementation below is based on existing setups using keras (see <https://github.com/danielenricocahall/Keras-SegNet> and <https://github.com/Runist/SegNet-keras>).

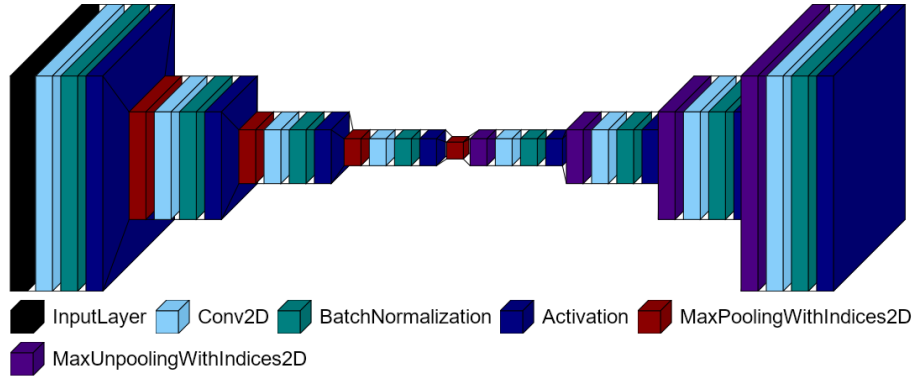


Figure 7: Architecture of simplified SegNet

In contrast to the other models that are presented within the scope of this project, the SegNet requires the use of layers that are not predefined in Keras for the use of the mentioned max-pooling indices. This concerns `MaxPoolingWithIndices2D` and `MaxUnpoolingWithIndices2D`.

3.2 U-Net Architecture

The U-Net architecture was motivated by biological image segmentation and the winning segmentation net in the ISBI challenge 2015. Similar to the SegNet, it applies encoder feature maps to the decoder in order to use trained feature maps more efficiently. With this approach, the U-Net presented in [5] could achieve very good results by using only little data for the training.

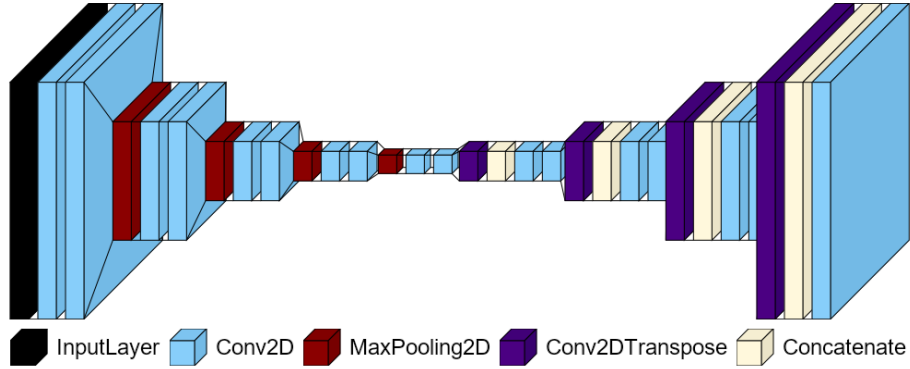


Figure 8: Architecture of U-Net with complexity of 4

3.3 Advanced Decoders

Starting with the simple decoder in the first part of the project (alg. 1), more complex decoders were created and tested. In order to have more parameters for the learning, additional layers like **Upsampling**, **Conv2D**, **Conv2DTranspose**, **BatchNormalization**, **Activation**, **Add** and **Dropout** are applied in the decoder architecture. Different combinations are tested as well as an architecture inspired by the examples for Computer Vision (https://keras.io/examples/vision/oxford_pets_image_segmentation/). The encoder used for these investigations is the **MobileNetV2**, because of its fewer parameters and therefore faster training and prediction. The most important test are shown in this report:

- Multi- & Single-Class Segm. with **MobileNetV2** using **Conv2DTranspose** with and without **BatchNormalization**
- Multi- & Single-Class with **MobileNetV2** using two **Conv2D** while **Add** an **Upsampling**, seen in alg. 8
- Multi- & Single-Class with **MobileNetV2** using **Conv2DTranspose** while **Add** an **Upsampling**, seen in alg. 9
- Multi- & Single-Class with **MobileNetV2** using **Dropout** while **Add** an **Upsampling**, seen in alg. 7
- Multi- & Single-Class with **MobileNetV2** using two paths with different layer structures, seen in alg. 6.

In alg. 6 the decoder structure for the most advanced case is shown. In all of these tests of advanced decoders, five layers were applied in the decoder. These layers were built always exactly the same and similarly but simpler to

the here shown structure. They all consisted of an addition of an **Upsampling** and another path including more parameters. Before the paths were added a **BatchNormalization** and an **Activation** were computed.

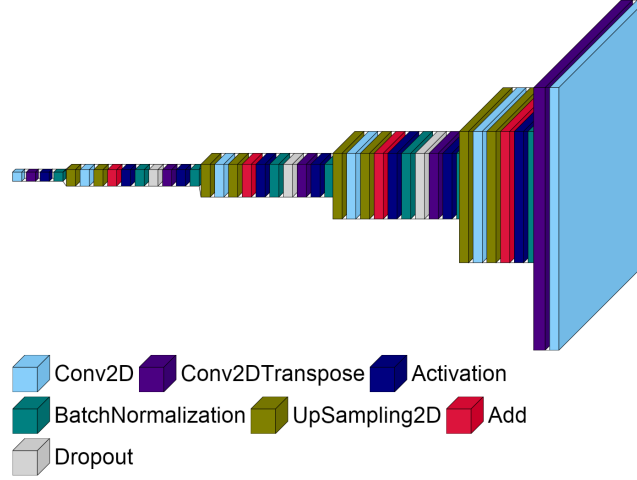


Figure 9: Architecture of the Advanced Decoder

4 Final Segmentation Architecture with a Pre-trained Encoder

For the final segmentation net, the U-Net architecture was combined with a **MobileNetV2** pre-trained encoder. Also not all layers of the **MobileNetV2** were used in order to keep the network small and slim. The skipped connections from the encoder were exported and added in the different decoder layers as it can be seen in alg. 10. This combination was trained in 50 epochs for a single-classification and a multi-classification.

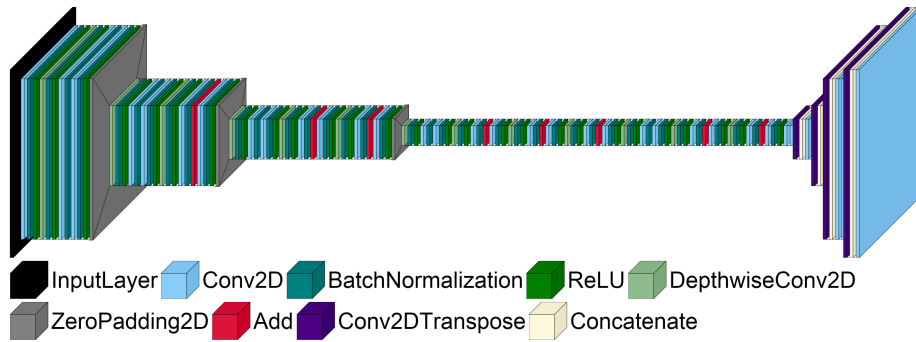


Figure 10: Architecture of Final Model

5 Evaluation of the Results

5.1 Pre-trained Encoders with a Simple Decoder

In the first investigations, three pre-trained encoders were applied together with a simple decoder architecture. Only the decoder was trained. The **VGG16**, the **MobileNetV2** and the **ResNetV2** were to be tested.

The results for single- & and multi-classification can be seen in the results plots in fig. 11 and fig. 12. The results were achieved by using an early stopping as **callback** in training and with a shifting of the images in the data generation of the training data.

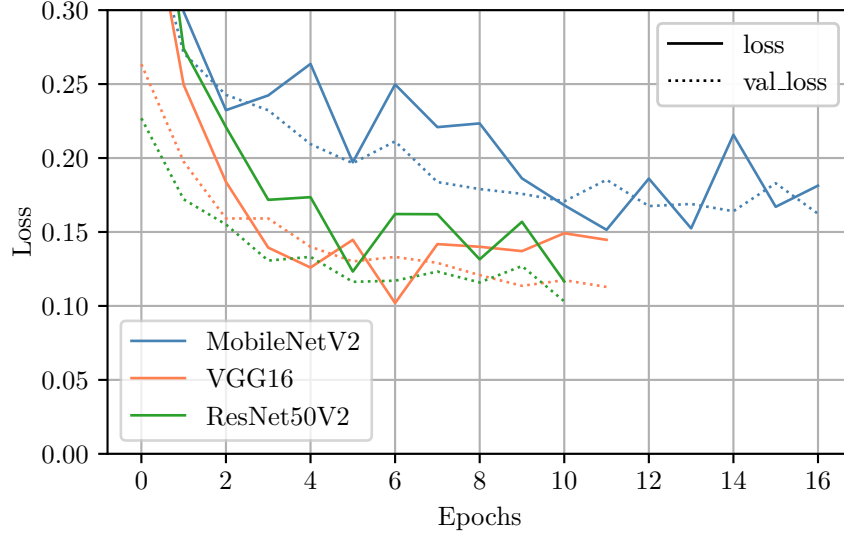


Figure 11: Losses of different pre-trained encoder for single-classification with training settings: Shift = True; Multi-Class = False; EarlyStopping = True; Epochs = 50; Image size = 256x256

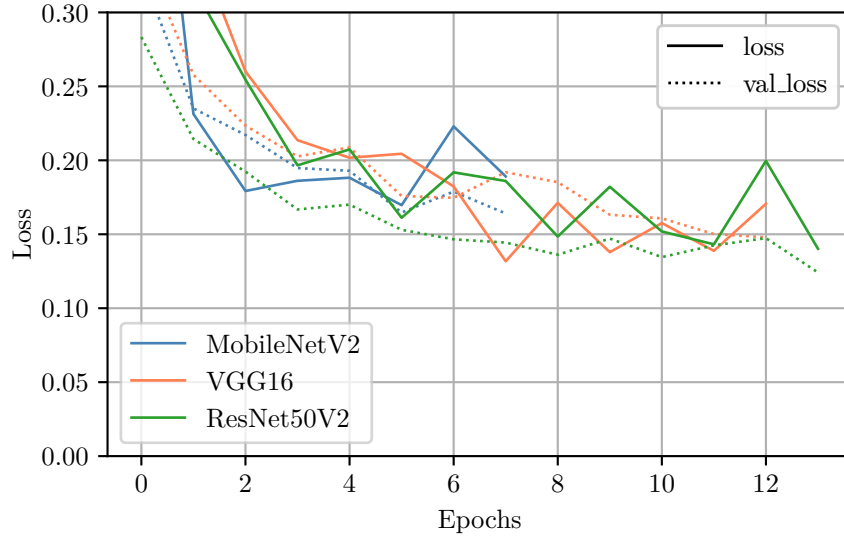


Figure 12: Losses of different pre-trained encoder for multi-classification with training settings: Shift = True; Multi-Class = True; EarlyStopping = True; Epochs = 50; Image size = 256x256

It was found that the shifting of the images for data generation could yield more training data, but hindered the CNN to achieve finer results. Also the early stopping ended the training too early, without being sure that convergence of the loss-curve was already reached.

Consequently, another training of the pre-trained encoders was performed with these improvements. The loss-curves for single- (fig. 13) and multi-classification (fig. 15) show this improvement of the learning regarding the previous results (fig. 11 & fig. 12).

As one can see in the predictions of these different encoders (in fig. 14 and fig. 16), all three of them performed as weak in the segmentation task here. Although the multi-classification models in fig. 16 found white planks and the bar-code label relatively accurate, their overall performance is not sufficient for the segmentation task.

For the single-classification models in fig. 14, the black plank is poorly detected and the edges are largely blurred.

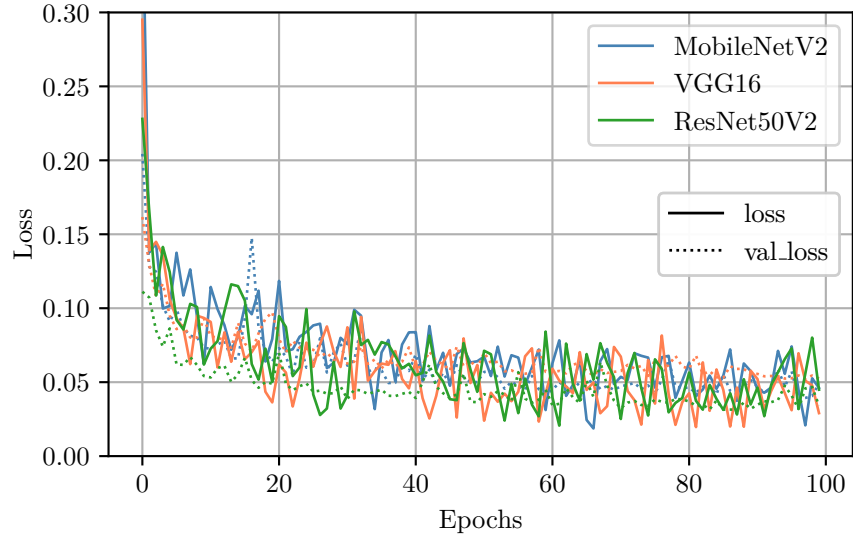


Figure 13: Losses of different pre-trained encoder for single-classification with training settings: Shift = False; Multi-Class = False; EarlyStopping = False; Epochs = 100; Image size = 256x256

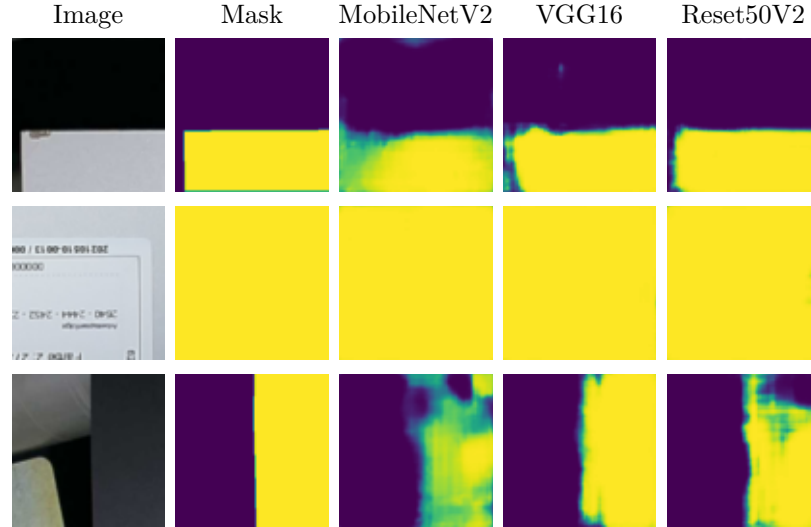


Figure 14: Predictions for different encoders for single-classification with training settings: Shift = True; Multi-Class = False; EarlyStopping = True; Epochs = 50; Image size = 256x256

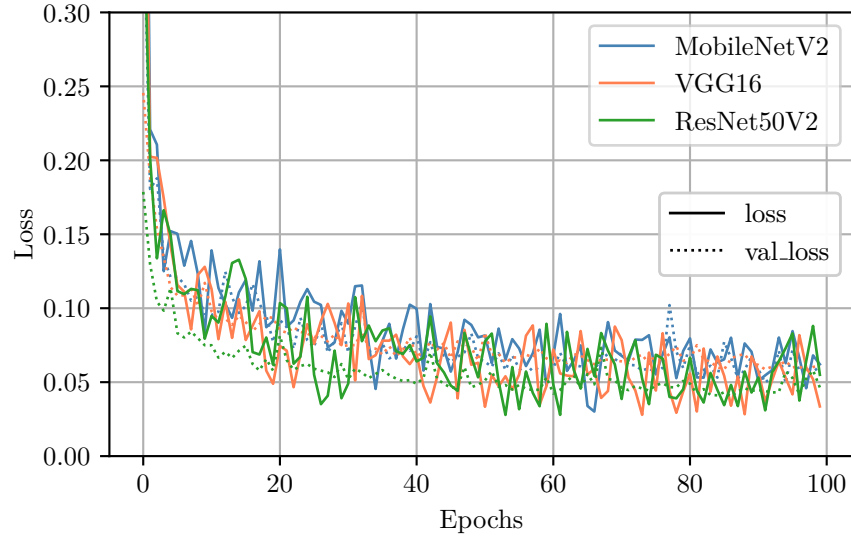


Figure 15: Losses of different pre-trained encoder for multi-classification with training settings: Shift = False; Multi-Class = True; EarlyStopping = False; Epochs = 100; Image size = 256x256

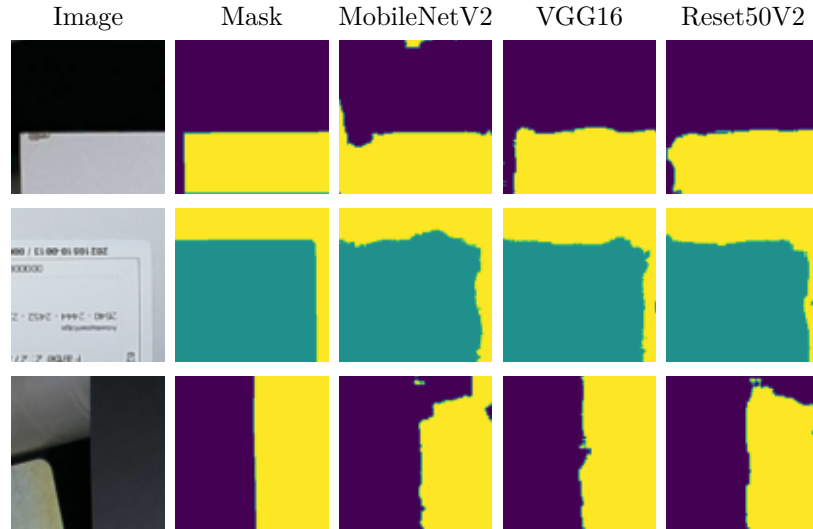


Figure 16: Predictions for different encoders for multi-classification with training settings: Shift = True; Multi-Class = True; EarlyStopping = True; Epochs = 50; Image size = 256x256

5.2 SegNet

As one can see from fig. 17, even the simplified version of the SegNet convinces with its ability to accurately detect boundaries. However, relevant boundaries in the background are detected as well, especially in the case of dark-colored planks. This causes the prediction within the homogeneous surfaces to be significantly blurred compared to other architectures considered within this project.

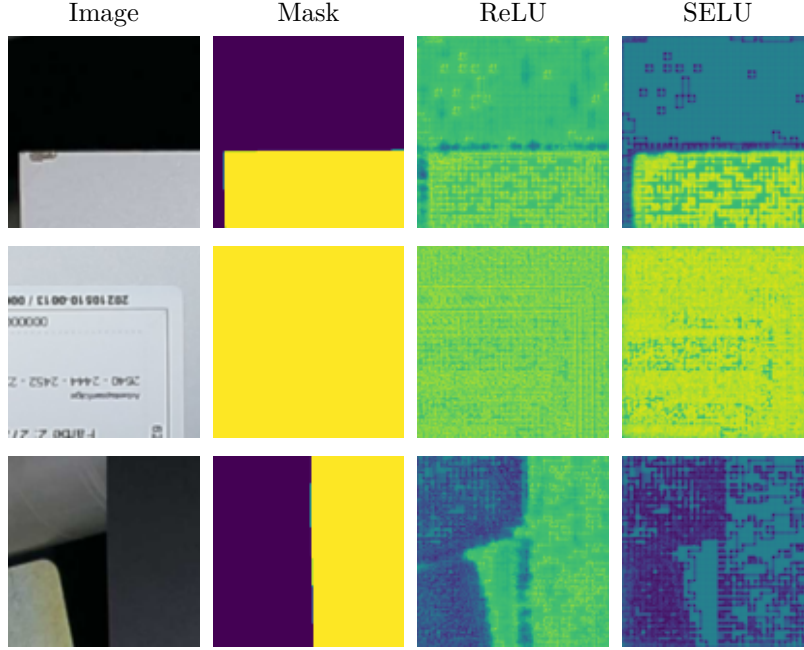


Figure 17: Predictions of simplified SegNet for ReLU- and SELU-activation with training settings: Shift = False, Multi-Class = False; EarlyStopping = False; Epochs = 50; Image size = 256x256

Despite the use of **ReLU-Activation** in [1] we considered a replacement by **SELU-Activation**. The respective losses of the training in fig. 18 with both activation functions do not indicate significant differences. In the presented images the detection of light-colored planks still seems to be clearer when **SELU** is used.

For the multi-class segmentation task the simplified SegNet does not perform as well as before. However, the loss plot in fig. 19 suggests that further training might lead to better results in this case.

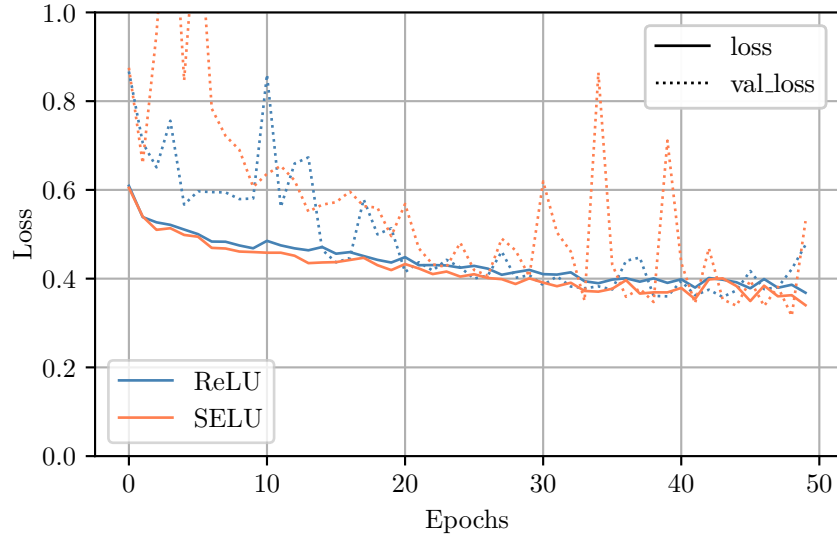


Figure 18: Loss of simplified SegNet for ReLU- and SELU-activation with training settings: Shift = False; Multi-Class = False, EarlyStopping = False; Epochs = 50; Image size = 256x256

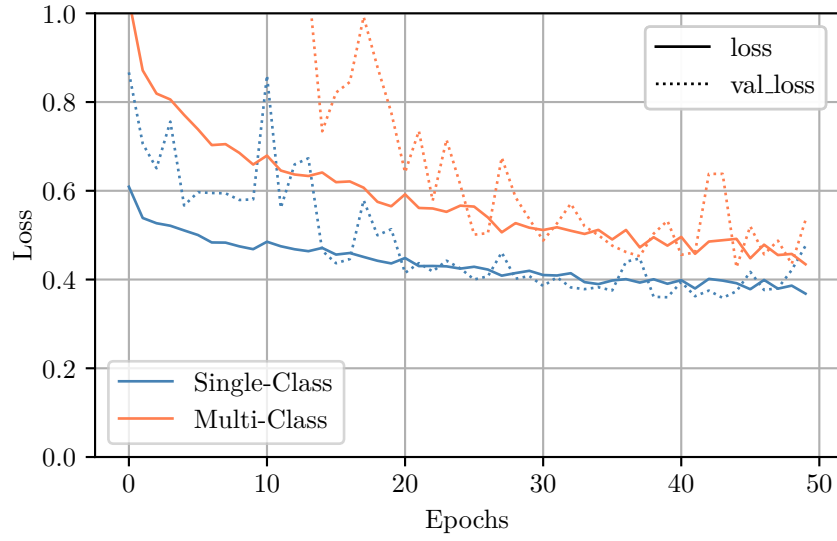


Figure 19: Loss of simplified SegNet for single- and multi-classification with ReLU-activation and training settings: Shift = False; Multi-Class = \times , EarlyStopping = False; Epochs = 50; Image size = 256x256

5.3 U-Net

The results of the U-Net architecture differ with a change in complexity regarding fig. 20 and fig. 21 for single-classification. In alg. 5 it can be seen that the complexity is used as a measure for the applied layer packages in the U-Net. Small complexities like a complexity of 2 do not seem to manage the training as well as higher complexities like 4 or 5. However, higher complexities (complexity of 6 in fig. 20) also do not improve the training. So a suitable complexity had to be found. The best results were achieved while using a complexity of 5 as it can be seen in fig. 21 and in fig. 23. Therefore, all the predictions (fig. 22 and fig. 24) are computed using U-Nets of a complexity of 5.

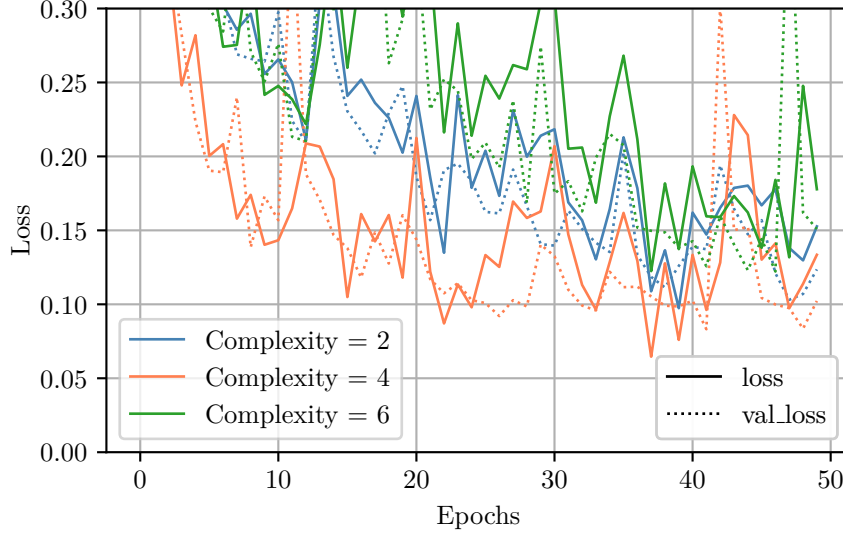


Figure 20: Losses of U-Net with different complexities using sliced images for single-classification with training settings: Shift = True; Multi-Class = False; EarlyStopping = False; Epochs = 50; Complexity = x ; Image size = 256x256

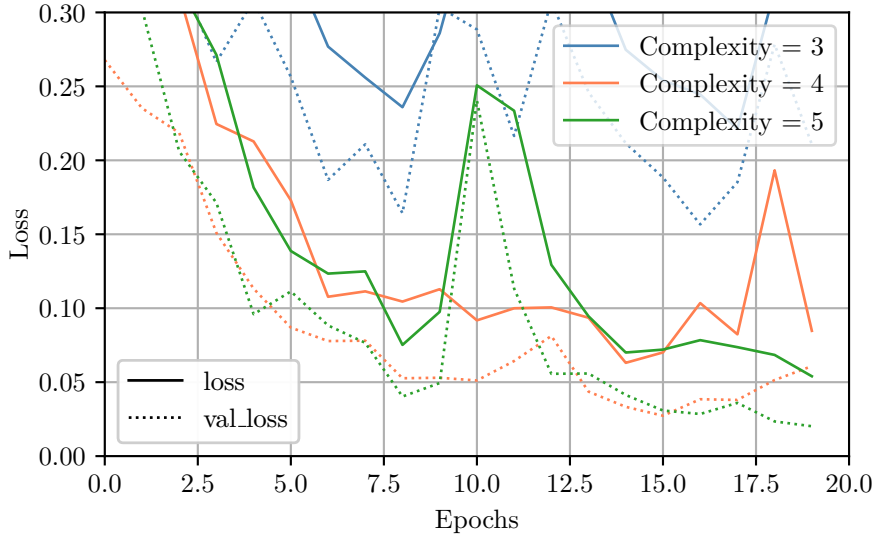


Figure 21: Losses of U-Net with different complexities using big images for single-classification with training settings: Shift = False; Multi-Class = False; EarlyStopping = False; Epochs = 20; Complexity = x ; Image size = 256x3072

Regarding the predictions in fig. 24 for multi-classification and in fig. 22 for single-classification, one can see a clear reproduction of the plank, even more precise than the actual mask for the single-classification.

Differences can be found when increasing the complexity. It leads to a visibly stronger consideration of other features in the image, so that is already at a complexity of 6, the U-Net found parts of the conveyor belt system in the background as plank. This may be caused by a strong dependence of the learning on white plates due to a ratio of 3:1 of white to black planks in the data. In addition it might be influenced by the darker background, reinforcing problems with segmentation of the black planks.

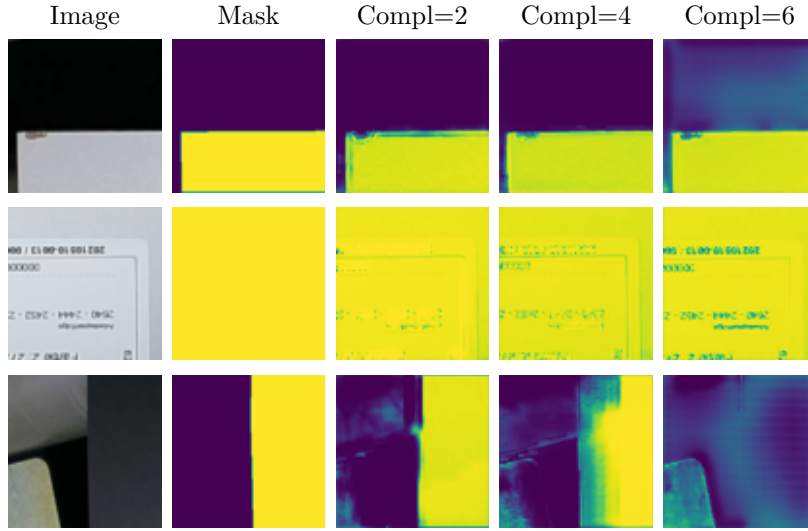


Figure 22: Predictions of U-Net with different complexities with training settings: Shift = True, Multi-Class = False; EarlyStopping = False; Epochs = 50; Complexity = x ; Image size = 256x256

The results in the loss-curves for the single-classification U-Net in fig. 20 cover the conclusions drawn before from the predictions in fig. 14. In the multi-classification (as well as the single-classification) U-Net the predictions are more accurate when the U-Net was trained on big-sized images compared to a U-Net trained on sliced images. This was possible here, because of the self edited architecture of the U-Net that even non-squared images could be processed as input. The additional information that was given by the context of the big-sized image might have led to the better performance of this model, as it can be found in the loss plot in fig. 23.

Regarding the predictions in fig. 24 one can see that the segmentation is very precise if the whole big-sized images are used within the U-Net. If sliced images are considered, then the result is significantly worse.

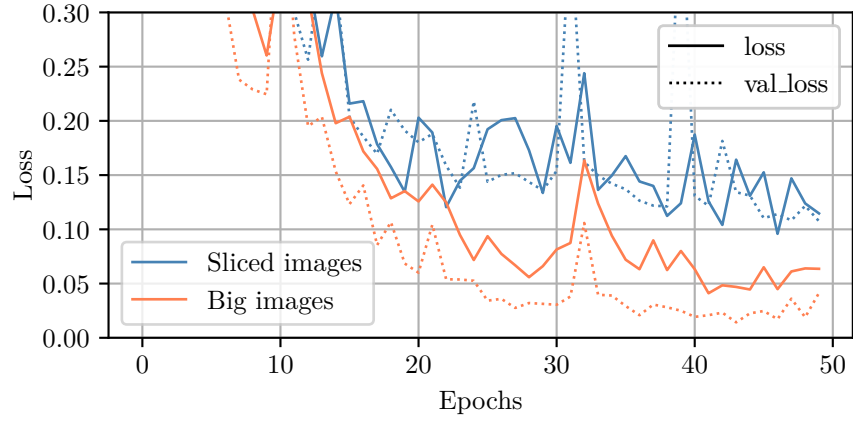


Figure 23: Loss of U-Net with complexity of 5 for multi-classification with training settings: Shift = False; Multi-Class = True; EarlyStopping = False; Epochs = 50; Complexity = 5; Image size = \mathbf{x}

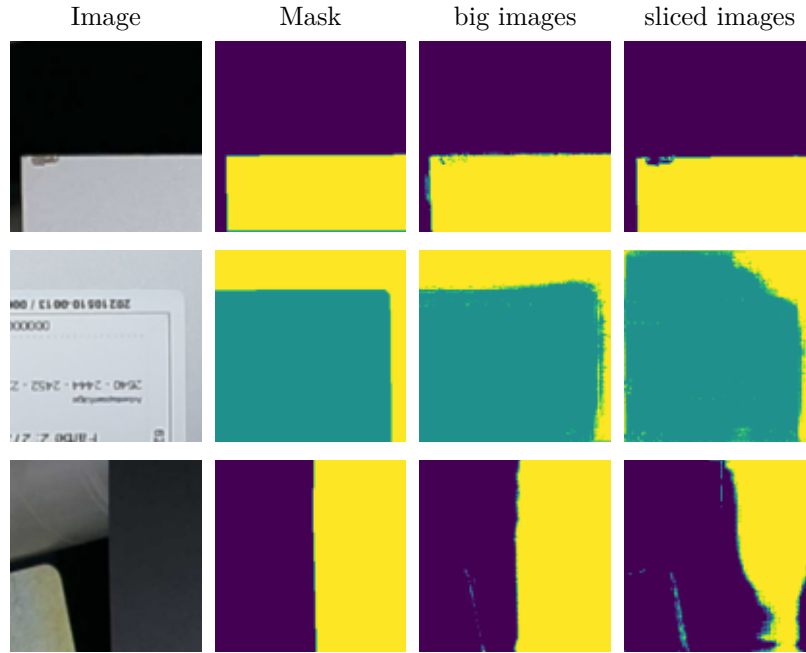


Figure 24: Predictions of U-Net with different sizes of images with training settings: Shift = False, Multi-Class = True; EarlyStopping = False; Epochs = 50; Complexity = 5; Image size = \mathbf{x}

Further the **jaccard**-loss was implemented in the training of the U-Net (see alg. 2) and compared to the standard loss **BinaryCrossentropy** for single-classification, seen in fig. 25. For the comparison of these losses the **dice_metric** was used (seen in alg. 3). It is visible that the **BinaryCrossentropy** yields better results the further the training is performed.

In the predictions of these two models, shown in fig. 26, one can see that the model using the **jaccard**-loss is not able to detect black planks at all while the other model at least finds a light shadow of this plank. Additionally, the model using the **BinaryCrossentropy** could even differ the background from the plank in the small gap in the lower left corner.

Regarding these results, further models will always contain the **BinaryCrossentropy** when it comes to single-classification.

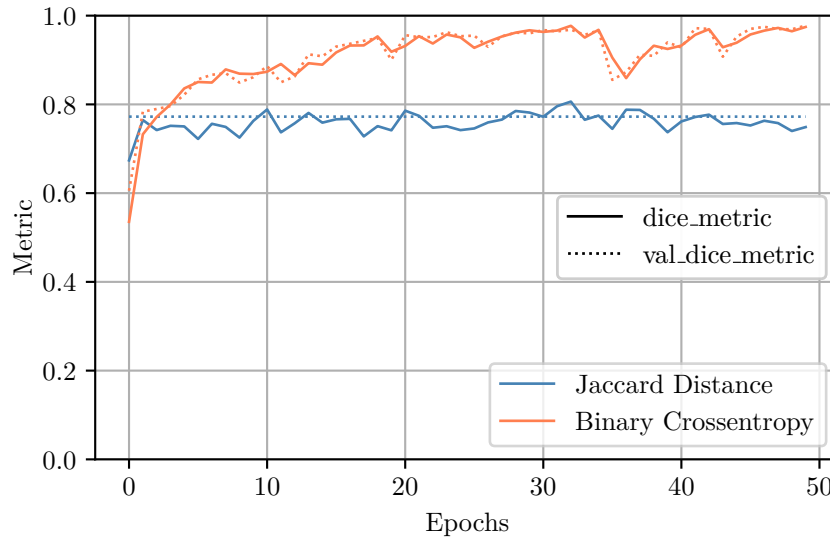


Figure 25: Metric of U-Net with complexity of 5 for multi-classification with different losses with training settings: Shift = False; Multi-Class = False; EarlyStopping = False; Epochs = 50; Complexity = 5; Image size = 256x256

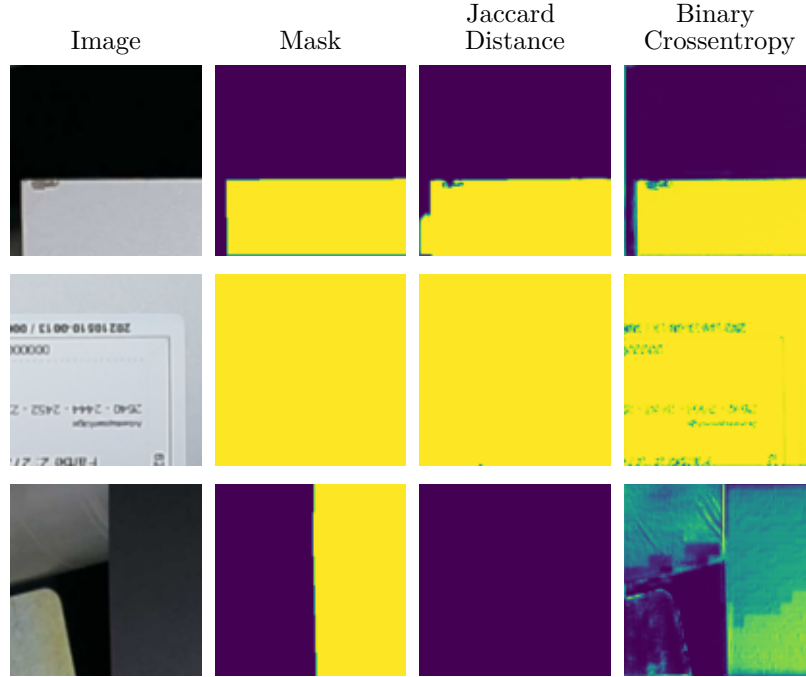


Figure 26: Predictions of U-Net with complexity of 5 for multi-classification with different losses with training settings: Shift = False; Multi-Class = False; EarlyStopping = False; Epochs = 50; Complexity = 5; Image size = 256x256

5.4 Advanced Decoder

From the results for the advanced decoder architecture one can see in fig. 27 that the **BatchNormalization** doesn't lead to a big improvement for the learning in case of single- and multi-classification. In the case of multi-classification the loss-curves with the **BatchNormalization** lie even above the loss-curve without normalization.

However, it must be considered that only one **BatchNormalization** was used (alg. 9) and not in a reasonable position. This could disturb the result and so the conclusions. It was repaired in the advanced decoder architecture in alg. 6 later on.

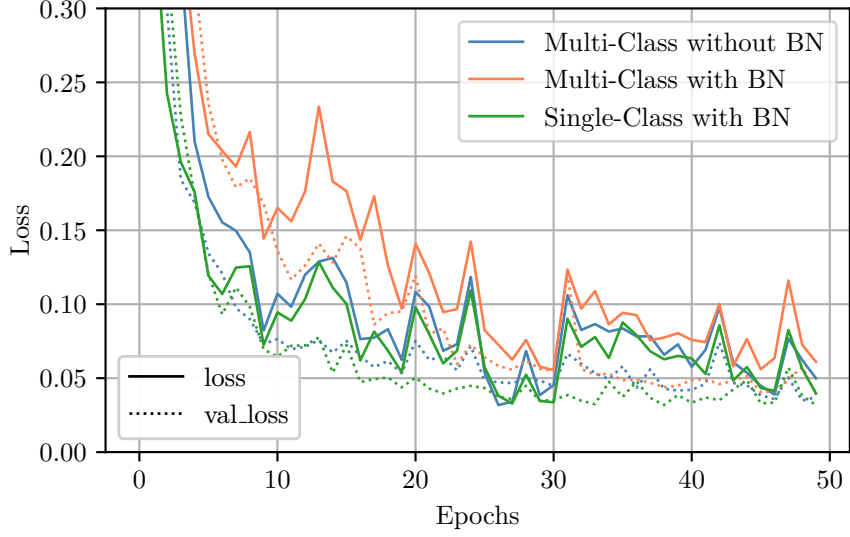


Figure 27: Loss for a decoder with Conv2DTranspose layers for multi- & single-classification with training settings: Shift = False; Multi-Class = \mathbf{x} ; EarlyStopping = False; Epochs = 50; Image size = 256x256

For the four advanced decoder structures (alg. 6) tested here, the losses are plotted in fig. 28 together with the predictions (fig. 29) for the multi-classification. All different decoder models perform in a similar way.

In the loss-plot fig. 28 the decoder models reach a minimal loss of around 0.05 whereby the decoder with the transposed convolutional layer can be found as slightly better compared to the three other models. The predictions are nevertheless only roughly similar to the mask, the results are not usable for the segmentation purpose.

A similar conclusion can be found for the single-classification task while using these four decoder architectures. There, the losses in fig. 30 end up at around 0.1. As well as with the multi-classification the predictions for the single-classification in fig. 31 are not usable for the segmentation task. At the edges the predictions are still blurry and coarse and not sharp lines as it would needed to be. Regarding this point, the multi-classification yields clearer edges although these edges are not as straight as required.

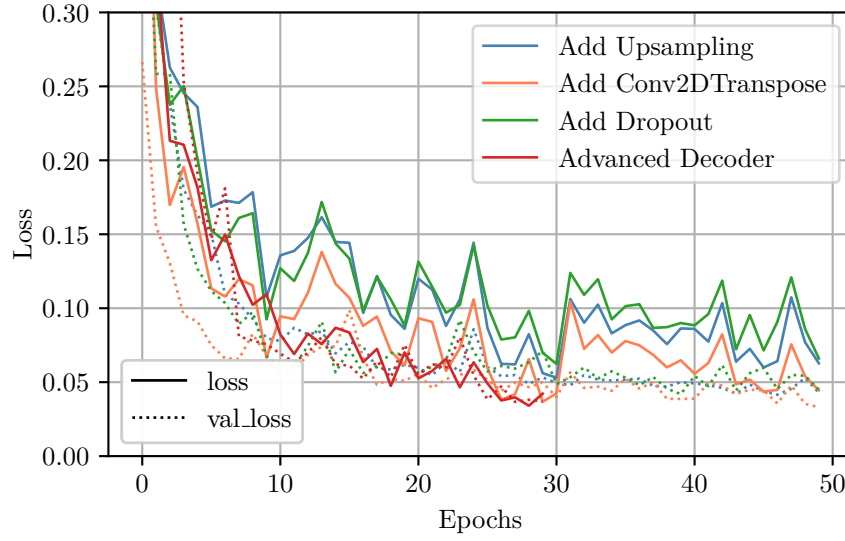


Figure 28: Loss for different decoder with added layers for multi-classification with training settings: Shift = False; Multi-Class = True; EarlyStopping = False; Epochs = 50; Image size =256x256

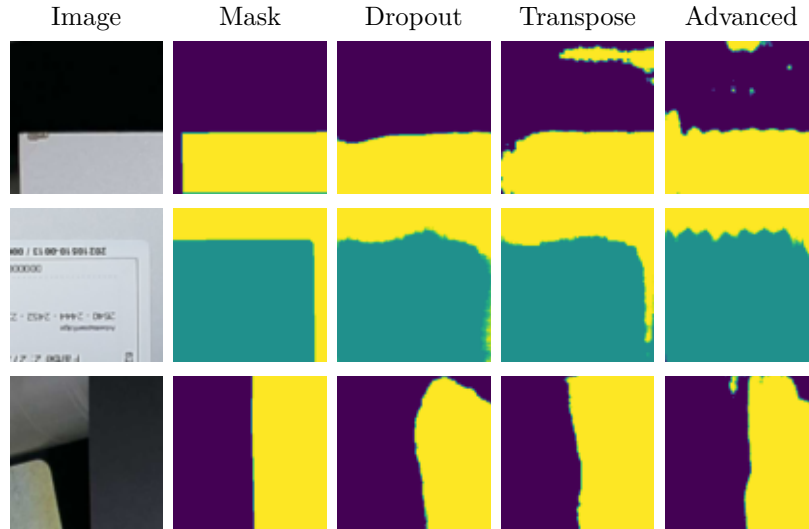


Figure 29: Predictions for different decoder with added layers for multi-classification with training settings: Shift = False; Multi-Class = True; EarlyStopping = False; Epochs = 50; Image size =256x256

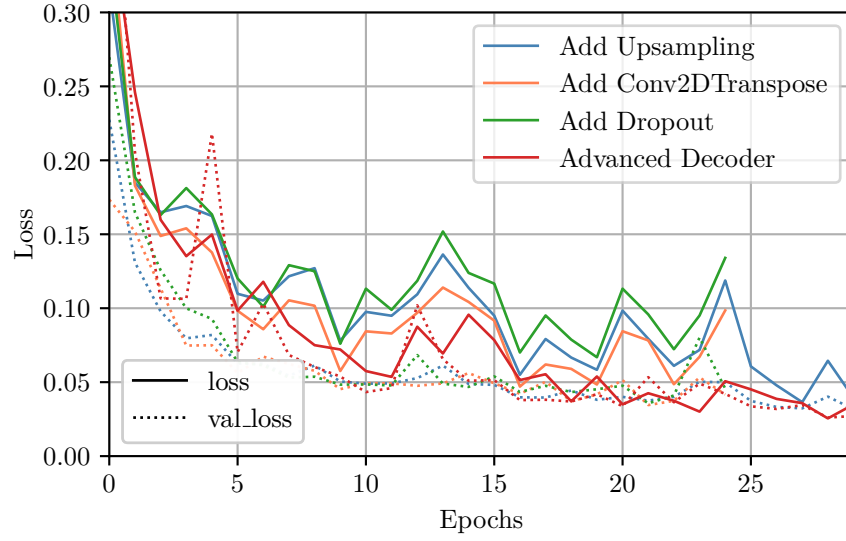


Figure 30: Loss for different decoder with added layers for single-classification with training settings: Shift = False; Multi-Class = False; EarlyStopping = False; Epochs = 25; Image size =256x256

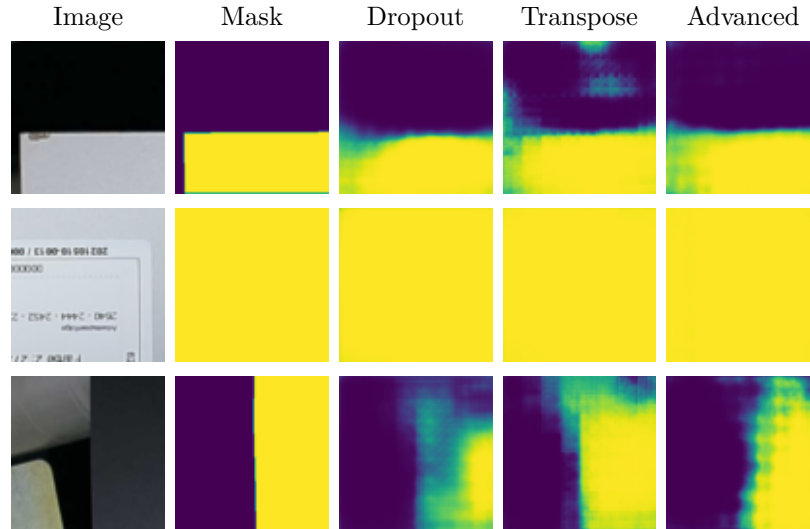


Figure 31: Predictions for different decoder with added layers for single-classification with training settings: Shift = False; Multi-Class = False; EarlyStopping = False; Epochs = 25; Image size =256x256

5.5 The Final Model

With the final model configuration a minimum loss of 0.01 for multi-classification and of even under 0.01 for single-classification could be reached as it can be seen in fig. 32. In order to yield finer results a fine-tuning was performed, meaning that the pre-trained model was now included in the training for another 20 epochs. This didn't lead to greater improvements, only for the single-classification a slightly decrease can be observed in fig. 33.

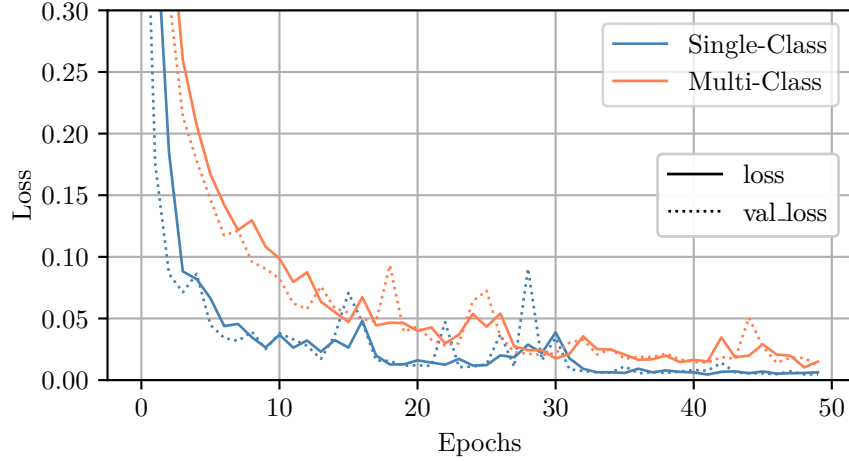


Figure 32: Losses of final network for single- & multi-classification with training settings: Shift = False; Multi-Class = x; EarlyStopping = False; Epochs = 50; Image size = 256x256

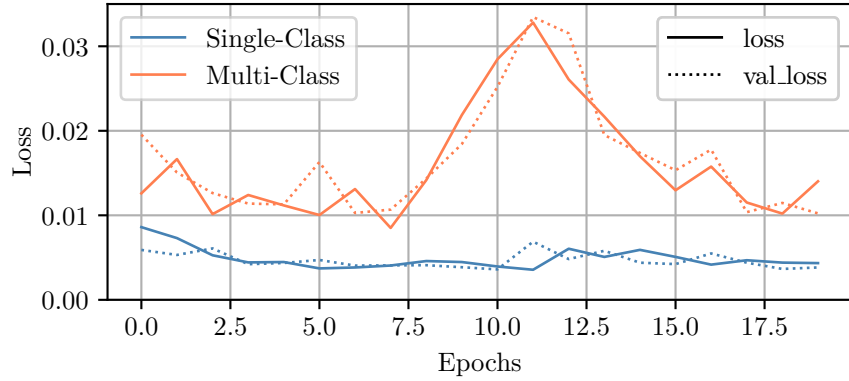


Figure 33: Fine tuning of final network for single- & multi-classification with training settings: Shift = False; Multi-Class = x; EarlyStopping = False; Epochs = 20; Image size = 256x256

Regarding the prediction plots for this final model in fig. 34, a very accurate segmentation of the edges of the planks can be seen for the white as well as for the black planks. The fine-tuning could improve the single-classification model in order to remove small mistakes in detection of the black planks. However, the segmentation of the bar-code-label still causes difficulties for the segmentation net. So, only the rough shape of the bar-code can be found. This may be sufficient to scan the bar-codes as the label also contains several centimeters around the code.

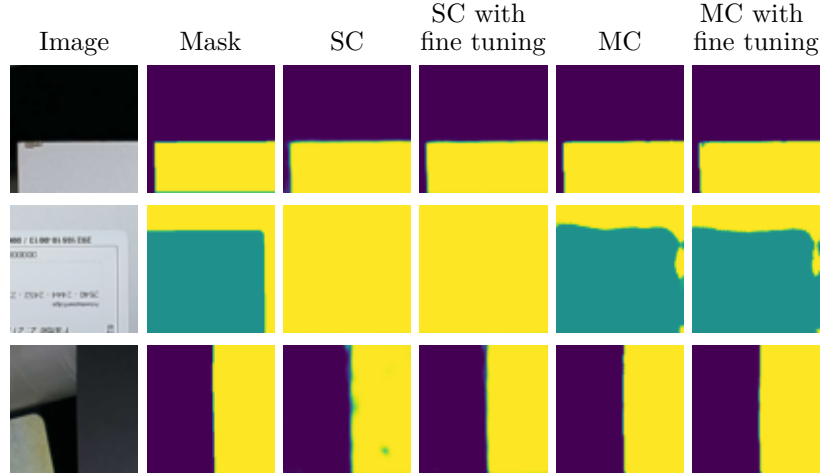


Figure 34: Predictions of final network single- & multi-classification with training settings: Shift = False, Multi-Class = **x**; EarlyStopping = False; Epochs = 50; Image size = 256x256

5.6 Benchmark Studies on a Coral Chip

As final step, the segmentation nets are implemented on a coral chip which can later on be used with a raspberry pi 4 for real-time segmentation of wooden planks. In this report only a benchmark study is discussed.

Therefore the models containing **Keras** software based on the **tensorflow** package are compressed into **tensorflowlite**-models. For these 'lite' models only few layer-types are implemented and a nightly version of **tensorflowlite** had to be used. Since the SegNet contains layer structures that are not predefined in **Keras**, significantly more effort would be needed to prepare its architecture for the coral chip. Therefore, the SegNet is not considered within this benchmark study.

The results are shown in fig. 35 for all pre-trained models, the U-Net and the SegNet. The results for different types of the U-Net are shown in fig. 36 .

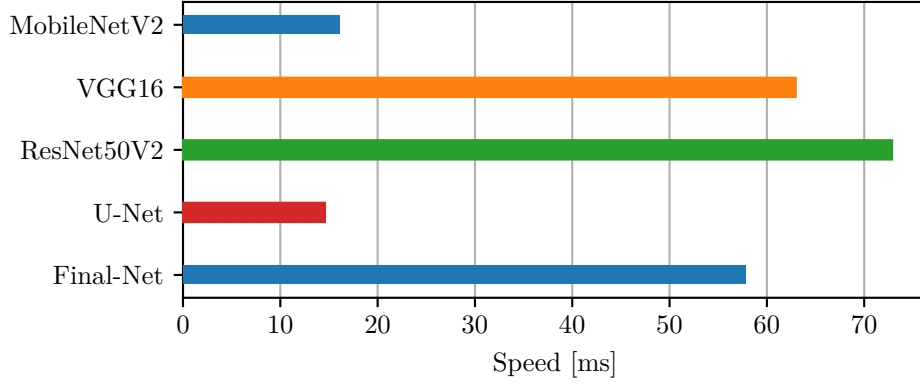


Figure 35: Benchmark study for the Pre-trained models, the U-Net and the SegNet

Although the U-Net trained on big-sized images performed better (fig. 24), it cannot be used for fast predictions on a coral chip due to its non-squared input. But with using bigger images with the U-Net, the more information led to better segmentation. So as further improvement, a pre-processing step could be applied which produces a usable input while containing as much information as the big-sized, but non-squared images.

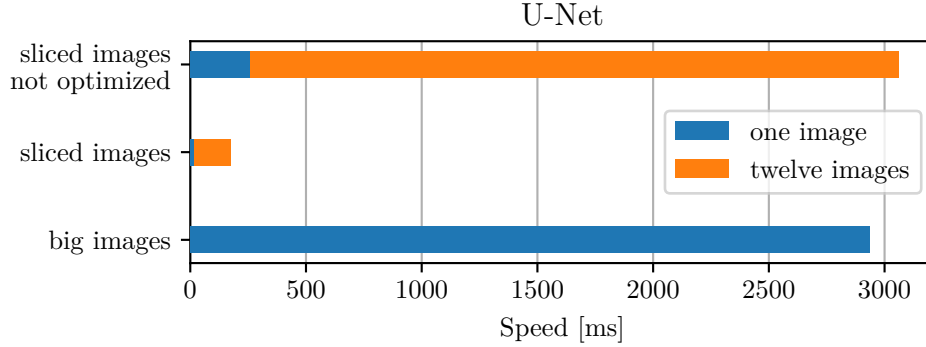


Figure 36: Benchmark study for different inputs while using a U-Net:
Comparing a U-Net trained with not optimized, sliced images,
a U-Net trained on optimized, sliced images and a U-Net trained
on full-scale images (256x3072).

6 Summary & Outlook

In this project different pre-trained encoders from CNNs like the `MobileNetV2`, the `VGG16` or the `ResNet50V2` have not proven to be suitable solutions for the segmentation task of coated wooden planks. Also with different decoder architectures the results were not sufficient. Advanced segmentation architectures as the U-Net and the SegNet had been proven as better approaches for the segmentation here. The best results could be achieved with a U-Net applied on a pre-trained encoder of a `MobileNetV2` with a self-designed decoder.

For the further development, a lot of improvements can be realized. In the post-processing, now a threshold function could be applied to the results in order to detect the damages on the plank while controlling the damages size over a threshold value. Alternatively, another R-CNN could be used for the damage detection on the plank-segmentation.

A problem that appeared in the SegNet and the U-Net predictions was the detection of black planks. This could be improved with more data of black planks and maybe with two CNNs, each trained on only white or only black planks. Also a deeper SegNet architecture as used in the first proposal of the SegNet [1] could lead to better results. The noise in predictions of the SegNet could be suppressed using a post-processing function. Also the color and the background remain the same, so developing a pre-processing function to insert this kind of information into the images could also improve the performance.

For the advanced decoder structures also deeper structures and a longer training could yield better results. Nevertheless these decoders will not perform as good as advanced architectures like the SegNet and the U-Net given the same training and data. The new findings of implementing these decoders can be used for improving the SegNet and the U-Net architecture further as already done in the final model of this project.

References

- [1] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(12):2481–2495, 2017.
- [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks, 2016.
- [4] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. *Artificial Intelligence Review*, 34:1–54, 2015.
- [5] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation.
- [6] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks, 2019.
- [7] Adnan Saood and Iyad Hatem. Covid-19 lung ct image segmentation using deep learning methods: U-net versus segnet. *BMC Medical Imaging*, 21(1), 2021.
- [8] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015.
- [9] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net.
- [10] Zhi Tian, Tong He, Chunhua Shen, and Youliang Yan. Decoders matter for semantic segmentation: Data-dependent decoding enables flexible feature aggregation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

Appendices

A Parameter Summary of All Architectures

Model	Total parameter	Trainable parameter
Simpler Decoder Architecture		
MobileNetV2	2.356.011	98.027
VGG16	14.757.419	42.731
ResNet50V2	23.718.123	153.323
SegNet Architecture		
SegNet	261.263	260.361
U-Net Architecture		
U-Net Complexity = 3	141.475	141.475
U-Net Complexity = 4	484.115	484.115
U-Net Complexity = 5	1.757.835	1.757.835
U-Net Complexity = 5; big images	1.757.835	1.757.835
Advanced Decoder Architecture		
Added Upsampling	2.451.907	193.795
Added Transpose	2.549.891	291.779
Added Dropout	2.549.891	291.779
Advanced Decoder	3.049.139	790.099
Final Model Architecture		
U-Net out of a MobileNetV2	1.095.604	479.348

Table 1: Listing of all models in this project with their total amount of parameters and the amount of trainable parameters. In this table only multi-classification models are taken into account.

B Algorithms

```

def generate_model(img_size, multiclass = True):

    encoder = MobileNetV2(include_top=False, weights='imagenet',
                          input_shape=(img_size),
                          classifier_activation =None)

    encoder.trainable = False

    d1 = UpSampling2D(size=(2, 2))(encoder.layers[-1].output)
    c1 = Conv2D(8, kernel_size=(3, 3), activation='selu', padding='same')(d1)
    d2 = UpSampling2D(size=(2, 2))(c1)
    c2 = Conv2D(16, kernel_size=(3, 3), activation='selu', padding='same')(d2)
    d3 = UpSampling2D(size=(2, 2))(c2)
    c3 = Conv2D(16, kernel_size=(3, 3), activation='selu', padding='same')(d3)
    d4 = UpSampling2D(size=(2, 2))(c3)
    c4 = Conv2D(16, kernel_size=(3, 3), activation='selu', padding='same')(d4)
    d5 = UpSampling2D(size=(2, 2))(c4)

    if multiclass:
        c5 = Conv2D(3, kernel_size=(1, 1), activation='softmax', padding='same')(d5)
    else:
        c5 = Conv2D(1, kernel_size=(1, 1), activation='sigmoid', padding='same')(d5)

    output = c5

    model = Model(inputs=encoder.inputs, outputs=output)
    return model

```

Algorithm 1: Simplified Decoder

```

def jaccard_distance_loss(y_true, y_pred, smooth=100):

    intersection = K.sum(K.abs(y_true * y_pred), axis=-1)
    sum_ = K.sum(K.abs(y_true) + K.abs(y_pred), axis=-1)
    jac = (intersection + smooth) / (sum_ - intersection + smooth)

    return (1 - jac) * smooth

```

Algorithm 2: Jaccard Distance

```

def dice_metric(y_pred, y_true):
    intersection = K.sum(K.sum(K.abs(y_true * y_pred), axis=-1))
    union = K.sum(K.sum(K.abs(y_true) + K.abs(y_pred), axis=-1))

    return 2*intersection / union

```

Algorithm 3: Dice Metric

```

def create_segnet(img_size, num_filters, multiclass = True):

    inputs = Input(shape=img_size)

    # Encoder
    conv_1 = Convolution2D(num_filters, (3, 3), padding="same",
                           kernel_initializer='he_normal')(inputs)
    conv_1 = BatchNormalization()(conv_1)
    conv_1 = Activation("relu")(conv_1)
    pool_1, mask_1 = MaxPoolingWithIndices2D(pool_size=(2, 2))(conv_1)
    ### Add two equivalent blocks ###
    conv_4 = Convolution2D(4 * num_filters, (3, 3), padding="same",
                           kernel_initializer='he_normal')(pool_1)
    conv_4 = BatchNormalization()(conv_4)
    conv_4 = Activation("relu")(conv_4)
    pool_4, mask_4 = MaxPoolingWithIndices2D(pool_size=(2, 2))(conv_4)

    # Decoder
    unpool_1 = MaxUnpoolingWithIndices2D(size=(2, 2))([pool_4, mask_4])
    conv_5 = Convolution2D(2 * num_filters, (3, 3), padding="same",
                           kernel_initializer='he_normal')(unpool_1)
    conv_5 = BatchNormalization()(conv_5)
    conv_5 = Activation("relu")(conv_5)
    unpool_2 = MaxUnpoolingWithIndices2D(size=(2, 2))([conv_5, mask_3])
    ### Add two equivalent blocks ###
    conv_8 = Convolution2D(n_labels, (1, 1), padding="same",
                           kernel_initializer='he_normal')(unpool_2)
    conv_8 = BatchNormalization()(conv_8)
    outputs = Activation(output_mode)(conv_8)

    # Classification
    if multiclass:
        outputs = Activation("softmax")(conv_8)
    else:
        outputs = Activation("sigmoid")(conv_8)

    # Model creation
    segnet = Model(inputs=inputs, outputs=outputs)

    return segnet

```

Algorithm 4: Simplified SegNet

```

def generate_model(img_size, complexity, multiclass = True):

    x = Input((img_size))
    inputs = x

    # Encoder
    f = 8
    layers = []

    for i in range(0, complexity):
        x = Conv2D(f, 3, activation='relu', padding='same')(x)
        x = Conv2D(f, 3, activation='relu', padding='same')(x)
        layers.append(x)
        x = MaxPooling2D()(x)
        f = f*2
    ff2 = 64

    # Bottleneck
    j = len(layers) - 1
    x = Conv2D(f, 3, activation='relu', padding='same')(x)
    x = Conv2D(f, 3, activation='relu', padding='same')(x)
    x = Conv2DTranspose(ff2, 2, strides=(2, 2),
                        padding='same')(x)
    x = Concatenate(axis=3)([x, layers[j]])
    j = j - 1

    # Decoder
    for i in range(0, complexity-1):
        ff2 = ff2//2
        f = f // 2
        x = Conv2D(f, 3, activation='relu', padding='same')(x)
        x = Conv2D(f, 3, activation='relu', padding='same')(x)
        x = Conv2DTranspose(ff2, 2, strides=(2, 2),
                            padding='same')(x)
        x = Concatenate(axis=3)([x, layers[j]])
        j = j - 1

    # Classification
    if multiclass:
        outputs = Conv2D(3, 1, activation='softmax')(x)
    else:
        outputs = Conv2D(1, 1, activation='sigmoid')(x)

    # Model creation
    model = Model(inputs=[inputs], outputs=[outputs])

    return model

```

Algorithm 5: U-Net

```

def generate_model(img_size, multiclass = True):
    # Encoder
    encoder = MobileNetV2(include_top=False, weights='imagenet',
                          input_shape=(img_size),
                          classifier_activation =None)

    encoder.trainable = False
    x = encoder.layers[-1].output

    # Decoder
    f = [16,32,64,128,256]
    d = [16,16,32,48,64]
    for i in range(len(f)):
        dUp = UpSampling2D(size=(2, 2))(x)
        dUp = Conv2D(d[-i], kernel_size=(3, 3),padding='same')(dUp)
        dUp = BatchNormalization()(dUp)
        dUp = Activation("selu")(dUp)
        x = Conv2DTranspose(f[i], (3, 3), strides=2, activation="selu",
                           padding="same")(x)
        x = Conv2D(d[-i], kernel_size=(3, 3),padding='same')(x)
        x = BatchNormalization()(x)
        x = Activation("selu")(x)
        x = Conv2D(d[-i], kernel_size=(3, 3),padding='same')(x)
        x = BatchNormalization()(x)
        x = Activation("selu")(x)
        x = add([x,dUp])
        x = Activation("selu")(x)

    # Classification
    if multiclass == True:
        c5 = Conv2D(3, kernel_size=(1, 1), activation='softmax',padding='same')(x)
    else:
        c5 = Conv2D(1, kernel_size=(1, 1), activation='sigmoid',padding='same')(x)
    output =c5

    # Model creation
    model = Model(inputs=encoder.inputs, outputs=output)

    return model

```

Algorithm 6: Advanced Decoder

```

def generate_model(img_size, multiclass = True):
    # Encoder
    encoder = tf.keras.applications.MobileNetV2(include_top=False,
                                                weights='imagenet', input_shape=(img_size),
                                                classifier_activation =None)
    encoder.trainable = False
    x = Conv2D(16, kernel_size=(3, 3), activation='selu',
              padding='same')(encoder.layers[-1].output)

    # Decoder
    f = [16,32,64,128,256]
    for i in range(len(f)):
        dUp = UpSampling2D(size=(2, 2))(x)
        x = Conv2DTranspose(16, (3, 3), strides=2, activation="selu",
                          padding="same")(x)
        x = Conv2D(16, kernel_size=(3, 3), activation='selu', padding='same')(x)
        x = add([x,dUp])
        x = Activation("selu")(x)
        x = BatchNormalization()(x)
        x = Dropout(0.3)(x)
    d5 = Conv2DTranspose(256, (3, 3), strides=2, activation="selu", padding="same")(x)

    # Classification
    if multiclass == True:
        c5 = Conv2D(3, kernel_size=(1, 1), activation='softmax', padding='same')(d5)
    else:
        c5 = Conv2D(1, kernel_size=(1, 1), activation='sigmoid', padding='same')(d5)
    output =c5

    # Model creation
    model = Model(inputs=encoder.inputs, outputs=output)

    return model

```

Algorithm 7: Advanced Decoder with Dropout

```

def generate_model(img_size, multiclass = True):
    # Encoder
    encoder = tf.keras.applications.MobileNetV2(include_top=False,
                                                weights='imagenet', input_shape=(img_size),
                                                classifier_activation =None)
    encoder.trainable = False
    x = Conv2D(16, kernel_size=(3, 3), activation='selu',
              padding='same')(encoder.layers[-1].output)

    # Decoder
    f = [16,32,64,128,256]
    for i in range(len(f)):
        dUp = UpSampling2D(size=(2, 2))(x)
        x = UpSampling2D(size=(2, 2))(x)
        x = Conv2D(16, kernel_size=(3, 3), activation='selu', padding='same')(x)
        x = add([x,dUp])
        x = Activation("selu")(x)
        x = BatchNormalization()(x)
        x = Dropout(0.3)(x)
    d5 = Conv2DTranspose(256, (3, 3), strides=2, activation="selu", padding="same")(x)

    # Classification
    if multiclass == True:
        c5 = Conv2D(3, kernel_size=(1, 1), activation='softmax', padding='same')(d5)
    else:
        c5 = Conv2D(1, kernel_size=(1, 1), activation='sigmoid', padding='same')(d5)
    output =c5

    # Model creation
    model = Model(inputs=encoder.inputs, outputs=output)

    return model

```

Algorithm 8: Advanced Decoder with Upsampling

```

def generate_model(img_size, multiclass = True):
    # Encoder
    encoder = tf.keras.applications.MobileNetV2(include_top=False,
                                                weights='imagenet', input_shape=(img_size),
                                                classifier_activation =None)

    encoder.trainable = False
    x = Conv2D(16, kernel_size=(3, 3), activation='selu',
              padding='same')(encoder.layers[-1].output)

    # Decoder
    f = [16,32,64,128,256]
    for i in range(len(f)):
        dUp = UpSampling2D(size=(2, 2))(x)
        x = Conv2DTranspose(16, (3, 3), strides=2, activation="selu",
                          padding="same")(x)
        x = Conv2D(16, kernel_size=(3, 3), activation='selu', padding='same')(x)
        x = add([x,dUp])
        x = Activation("selu")(x)
        x = BatchNormalization()(x)
    d5 = Conv2DTranspose(256, (3, 3), strides=2, activation="selu", padding="same")(x)

    # Classification
    if multiclass == True:
        c5 = Conv2D(3, kernel_size=(1, 1), activation='softmax', padding='same')(d5)
    else:
        c5 = Conv2D(1, kernel_size=(1, 1), activation='sigmoid', padding='same')(d5)
    output =c5

    # Model creation
    model = Model(inputs=encoder.inputs, outputs=output)

    return model

```

Algorithm 9: Advanced Decoder with Transposed Layer

```

def generate_model(img_size, multiclass = True):
    inputs = Input(shape=(img_size), name="input_image")

    # Encoder
    encoder = MobileNetV2(include_top=False,
                          weights='imagenet',
                          input_tensor=inputs,
                          classifier_activation =None)
    encoder.trainable = False
    skip_connection_names = ["input_image",
                             "block_1_expand_relu",
                             "block_3_expand_relu",
                             "block_6_expand_relu"]

    encoder_output = encoder.get_layer("block_13_expand_relu").output
    x = encoder_output

    # Bottleneck
    f = 64
    ff2 = 8
    x = Conv2D(f, 3, activation='relu', padding='same')(x)
    x = Conv2D(f, 3, activation='relu', padding='same')(x)
    x = Conv2DTranspose(ff2, 2, strides=(2, 2), padding='same')(x)
    x_skip = encoder.get_layer(skip_connection_names[-1]).output
    x = Concatenate(axis=3)([x, x_skip])

    # Decoder
    for i in range(2, 5):
        ff2 = ff2 * 2
        f = f // 2
        x = Conv2D(f, 3, activation='relu', padding='same')(x)
        x = Conv2D(f, 3, activation='relu', padding='same')(x)
        x = Conv2DTranspose(ff2, 2, strides=(2, 2), padding='same')(x)
        x_skip = encoder.get_layer(skip_connection_names[-i]).output
        x = Concatenate(axis=3)([x, x_skip])

    # Classification
    if multiclass:
        outputs = Conv2D(3, 1, activation='softmax')(x)
    else:
        outputs = Conv2D(1, 1, activation='sigmoid')(x)

    # Model creation
    model = Model(inputs=[inputs], outputs=[outputs])

    return model

```

Algorithm 10: Final combined version of U-Net with the MobileNetV2