

# **Deep Learning for Instance Segmentation of Agricultural Fields**

**Master Thesis in Geoinformatics**

by

**Christoph Rieke**

15.09.2017

Friedrich-Schiller-University Jena

Supervisors:

Prof. Dr. Christiane Schmullius (FSU Jena)

Dr. Adam Erickson (MPI BGC Jena)



**Friedrich-Schiller-Universität Jena**

## **Abstract**

This thesis aims to delineate agricultural field parcels from satellite images via deep learning instance segmentation. Manual delineation is accurate but time consuming, and many automated approaches with traditional image segmentation techniques struggle to capture the variety of possible field appearances. Deep learning has proven to be successful in various computer vision tasks, and might be a good candidate to enable accurate, performant and generalizable delineation of agricultural fields. Here, a fully convolutional instance segmentation architecture (adapted from Li et al., 2016), was trained on Sentinel-2 image data and corresponding agricultural field polygons from Denmark. In contrast to many other approaches, the model operates on raw RGB images without significant pre- and post-processing. After training, the model proved successful in predicting field boundaries on held-out image chips. The results generalize across different field sizes, shapes and other properties, but show characteristic problems in some cases. In a second experiment, the model was trained to simultaneously predict the crop type of the field instance. Performance in this setting was significantly worse. Many fields were correctly delineated, but the wrong crop class was predicted. Overall, the results are promising and prove the validity of the deep learning approach. Also, the methodology offers many directions for future improvement.

# Contents

<b>1. Introduction</b>	<b>1</b>
1.1. Agricultural field parcels	1
1.2. Challenges for automated delineation methods	2
1.3. Traditional automated approaches	3
1.4. Deep learning for computer vision and remote sensing	4
1.5. Aims and thesis structure	8
<b>2. Deep learning for instance segmentation</b>	<b>9</b>
2.1. Neural network basics	9
2.2. Convolutional neural networks for image recognition	11
2.3. Fully convolutional neural networks for semantic segmentation	16
2.4. Object detection via Regional CNN	19
2.5. Instance segmentation	23
2.5.1. Overview of task and models	23
2.5.2. Deepmask and Multipathnet	25
2.5.3. MNC	26
2.5.4. FCIS	27
2.5.5. Mask R-CNN	29
<b>3. Methodology</b>	<b>31</b>
3.1. Study area and datasets	31
3.2. Preprocessing	32
3.3. Model implementation details	35
3.4. Evaluation metrics	38
<b>4. Results and evaluation</b>	<b>40</b>
4.1. Experiment 1: Single-class	40
4.1.1. Metrics	40
4.1.2. Examples	43
4.1.3. Training process	45

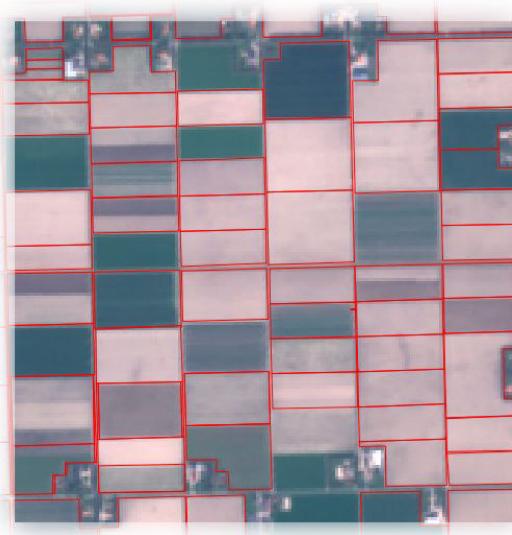
4.1.4. Field size	45
4.1.5. Overlap	46
4.1.6. Problematic training data	47
4.1.7. Empty image chips	47
4.2. Experiment 2: Multi-class	48
4.2.1. Metrics	48
4.2.2. Examples	49
<b>5. Future work</b>	<b>51</b>
<b>6. Conclusion</b>	<b>53</b>
<b>Literature</b>	<b>54</b>
<b>Figures</b>	<b>59</b>
<b>Tables</b>	<b>61</b>
<b>Appendix</b>	<b>62</b>
<b>Acknowledgments</b>	<b>64</b>
<b>Statement of Authorship</b>	<b>65</b>

# 1. Introduction

## 1.1. Agricultural field parcels

Geodata of field boundaries and parcels (fig. 1) is essential for many agriculture-related applications. Examples include crop type and yield monitoring, subsidy management as well as research and decision support for food security efforts. Initially, such data was primarily utilized by governmental agencies. Today, commercial businesses in the booming agricultural technology sector also require field parcels in various application areas like farm management, yield forecasting and precision farming.

Initiatives for the delineation and mapping of agricultural field parcels include the United States Common Land Unit (CLU) system and the European Union Land Parcel Identification System (LPIS) (Leo and Lemoine 2001). High quality data on agricultural parcels is primarily manually traced by experts in aerial or satellite imagery. In high resolution imagery, single field objects are distinguishable from the surrounding area by meaningful transitions in color, brightness or texture caused by land use changes. If available, the task can be supported by property or reference parcels, additional geospatial data provided by the farmer as well as information on crop type, field usage status etc. Although manual segmentation of the crop blocks can be very accurate, the task is time consuming for larger agricultural areas and prone to intra- and inter-observer variability due to varying experience, performance and diligence of the human operators.



**Figure 1:** Polygons of agricultural field parcels overlaid on satellite imagery.

The increasing demand for consistent, regularly updated field parcel data for agricultural regions throughout the world emphasizes the need for reproducible automatization of this

task. However, most traditional image segmentation models for the automated delineation of agricultural parcels are designed for small-scale applications and tuned to the specific environmental settings of the study area. An ideal automatic method would be able to deliver satisfying speed and accuracy for large areas with variable environmental settings. The thesis investigates whether this goal can be achieved via deep learning instance segmentation, a novel approach for the segmentation and classification of image objects.

## 1.2. Challenges for automated delineation methods

Agricultural fields in satellite images are relatively simple image objects when compared to most objects in natural photos. E.g., a picture of a human body appears much more complex because it consists of a variety of subobjects with different textures and shapes (face, hands, clothes, ...). Nevertheless, designing an algorithm for the delineation of agricultural parcel objects is non-trivial, both via traditional image segmentation as well as with deep learning techniques. The model needs to correctly dismiss undesired image objects, differentiate between adjacent fields with nearly similar spectral properties or barely visible borders, and correctly delineate field objects with multiple homogenous areas (e.g. caused by varying soil properties within a single field). Overall, the three main challenges for such an algorithm are: 1) the heterogeneous and volatile geography of agricultural fields, 2) the specific properties of satellite imagery, 3) inaccuracies in the available ground truth datasets. These challenges will be discussed in more detail below.



**Figure 2:** Heterogeneous geography of agricultural areas (Planet 2016).

**Field geography.** The physical appearance (shape, size, texture, color, etc.) of agricultural fields can vary depending on many physical and human geographical factors, both from the ground perspective as well as in satellite imagery. Possible sources of variation include the cultivated species, plant status, topography, soil properties, weather conditions, cultivation methods and other human influences. Due to such local variations, the task of accurately differentiating between several field objects can become ambiguous. This makes it difficult to create robust algorithms that recognize all kinds of field instances without overfitting to the specific scene. In addition to the variation within a single scene, satellite images can show even stronger intra-class variations. This applies to images from different agricultural regions, from different dates or growing seasons, or under different atmospheric and illumination conditions (fig. 2).

**Satellite imagery.** In high resolution satellite images, problems can also be caused by local background clutter objects and ground shadows, e. g. from farm machines or utility poles. Locally different atmospheric and lightning conditions (e.g. cloud shadows) can have similar effects. However, unlike for natural photos, image distortions due to viewpoint and distance variations of the imaging sensor are usually insignificant: Satellite image acquisitions are mostly tasked at nadir angle and the satellite height variation is mostly negligible due to the large distance to the imaged earth surface.

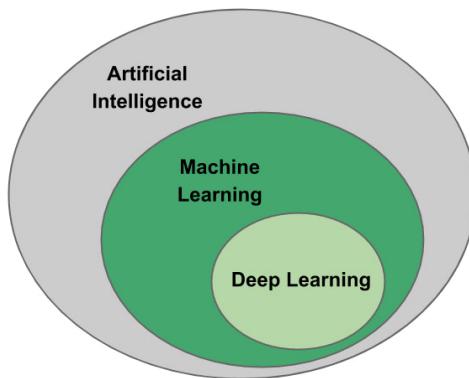
**Ground truth data.** Another challenge originates from the accuracy of the ground truth data sets that are used as validation and/or training data for automated algorithms. These datasets usually stem from manual tracing of field boundaries. However, manual tracing of image objects is heavily dependent on the used imagery and supplementary data. It also is a highly subjective task, inevitably leading to inaccuracies and ambiguities depending on the priorities of the operator. The ground truth dataset can be complemented by existing property parcel information that could be indistinguishable based on the imagery alone. Furthermore, a single field parcel polygon can potentially show several subfields due to unreported land use changes within the growing season. Given these constraints, even a very powerful automated algorithm can only become as good as the ground truth data, but will never be in perfect alignment with reality.

### 1.3. Traditional automated approaches

Several automated and semi-automated methods for the delineation of agricultural parcels from remote sensing imagery have been proposed in the literature. Most of these traditional

methods rely on image segmentation techniques applied to high-resolution, multispectral imagery of different band specifications and geographic regions. This includes image segmentation based on edge detection (Rydberg & Borgefors 2001, Mueller et al. 2004, Turker & Kok 2013), image value gradients (Butenuth et al. 2014), deformable “snake” algorithms (adapting to vector contours) (Torre & Radeva 2000) and textural properties (Da Costa et al. 2007, Tiwari et al. 2009, Yalcin et al. 2016). Another approach leverages multitemporal data by combining vegetation indices and edge detection (Yan and Roy 2014). These studies generally use manually selected features or parameters, which requires a priori knowledge of the scale, physical appearance or distribution of the fields in the scene. No training data is required. Many studies additionally employ region growing as well as post-processing techniques to refine the field detection process. García-Pedrero et al. (2017) use a machine learning approach to iteratively merge selected adjacent superpixels, with the merge criteria determined by a supervised Random Forest classifier via various spectral indices and texture features. The methodology requires training data, manual feature selection and potentially further post-processing for the completed delineation of independent field parcels.

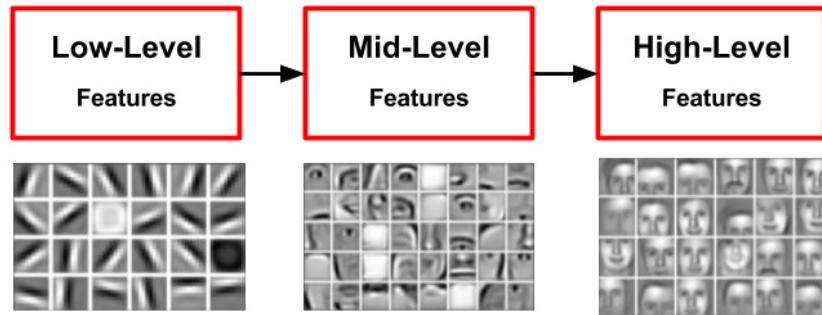
#### 1.4. Deep learning for computer vision and remote sensing



**Figure 3:** The field of deep learning.

Starting in 2012, deep learning, a new area of machine learning (fig. 3), lead to breakthrough advances, especially in computer vision research. Computer vision aims to enable computers to extract and analyse image information in order to gain higher-level understanding of an image. In recent years, deep learning has also enabled groundbreaking results in other fields of machine learning, such as machine translation (Wu et al. 2016), speech recognition (Amodei et al. 2015), reinforcement learning (Mnih et al. 2013), and robotics (Levine et al. 2016). For a comprehensive review see also LeCun et al. (2015). The

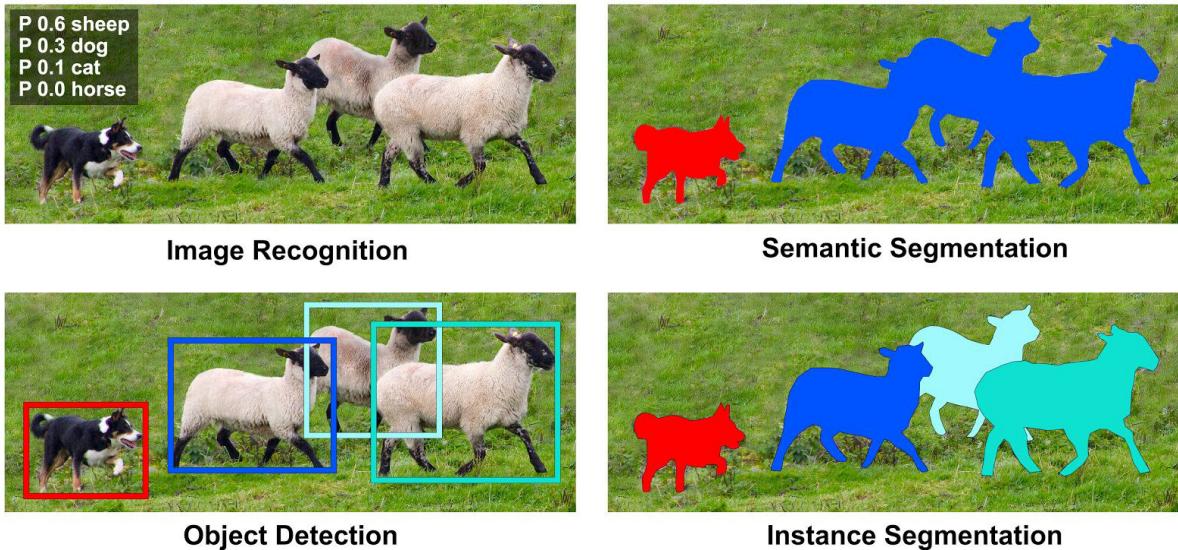
term “deep learning” was coined by the use of neural networks (non-linear statistical models) with many layers, so-called “deep” neural networks. By building up a hierarchical structure of more and more abstract representations of the input data, neural networks are able to make sense of complex datasets and use that knowledge to make predictions about similar data. Shallow neural networks have been around since the 1960s and have seen alternating phases of hype and disinterest by the research community. Key factors for the successful training of more powerful, deep neural networks were the increased computing power due to the use of graphics processing units (GPUs), the introduction of several key techniques (e.g. improved network parameter initialization) and a massive increase in the amount of available training data.



**Figure 4:** Feature hierarchy of image objects (example features edited from Jones 2014).

As mentioned above, the deep learning revolution began in 2012 with the introduction of deep convolutional neural networks (CNNs) (Krizhevsky et al. 2012) for the computer vision task of image recognition. CNNs are roughly inspired by the research of Hubel & Wiesel in 1959: After discovering that some neurons in the cat’s visual cortex get excited when the cat is shown images of edges of a particular orientation, they proposed a concept of information processing in the visual cortex (Hubel & Wiesel 1959). In their hypothesis, neurons that are low in the visual hierarchy scan for low-level image features (e.g. oriented edges) (fig. 4). These neurons feed to other cells that pick up on more complex features (e.g. the geometric shape of an eye by combining multiple edges), which in turn construct high-level features (partial or full image objects like a face). Convolutional neural networks, which also rely on feature hierarchies and a layered architecture, are a powerful tool for image recognition, the computer vision task of finding a single label for an image from a set of predefined categories. No manual feature engineering is required, the network is able to use raw image pixel data. Supervised training on large, annotated image datasets automatically tunes the CNNs to recognize and process relevant image features. In this way, the CNN maps the input image pixels to abstract feature representations and eventually to class probabilities.

The adaption and extension of the CNN architecture has also made deep learning the new standard for more complex computer vision tasks like object detection, semantic segmentation and instance segmentation (fig. 5).



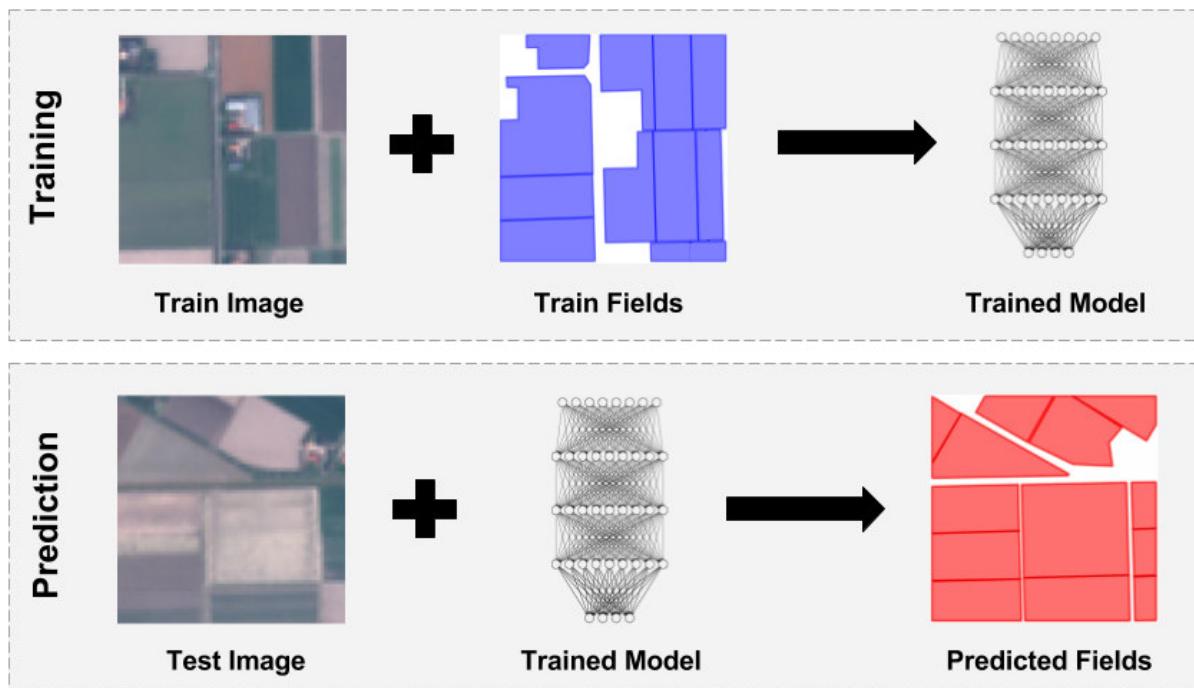
**Figure 5:** Four of the most important computer vision tasks.

The appropriate computer vision task for the segmentation and classification of agricultural fields via deep learning is **instance aware semantic segmentation** (IAIS) (or simply instance segmentation). IAIS delineates and classifies image objects while distinguishing between object instances of the same class. Instance segmentation lies at the intersection of object detection (predicting the bounding box and class of a variable number of image objects, but not segmenting them) and semantic segmentation (labeling each image pixel with a semantic category label, but not distinguishing between object instances of the same category). With semantic segmentation it would be possible to find single instances of agricultural fields if they were completely surrounded by areas of different landcover classes. However, agricultural field objects are often directly adjacent to each other and have touching boundaries. In this case, semantic segmentation would yield “clumped” polygon objects. Instance segmentation is able to distinguish between these connected or overlapping instances of the same class.

Deep learning has not reached mainstream adoption in the **remote sensing** community yet. However, considering the rapidly increasing volumes of earth observation imagery and geospatial data, it promises to be a great fit for many remote sensing applications. Recent studies have explored the potential of deep learning on satellite and aerial imagery for image recognition (Basu et al. 2015, Zhong et al. 2017), semantic segmentation (Längkvist et al.

2016, Marmanis et al. 2016, Saito et al 2016), and object detection (Vakalopoulou et al. 2015, Shenquan et al. 2016). To the author's knowledge, no existing publication has leveraged deep learning instance segmentation to segment and classify objects in remote sensing imagery.

### 1.5. Aims and thesis structure



**Figure 6:** High-level overview of Instance segmentation training and prediction.

This thesis aims to show the potential, advantages and challenges of deep learning instance segmentation for the automated delineation and classification of agricultural field parcels from medium resolution satellite imagery. The project uses a fully convolutional neural network architecture, adapted from Li et al. (2016). 159 042 digitized agricultural fields in Denmark and cloud free Sentinel-2 imagery (RGB bands) are used as training data. The thesis presents two experiments with increasing complexity: Firstly, the model is trained to segment all agricultural field instances. In the second experiment, the model is trained to simultaneously segment and classify the field instances based on their crop type. With the exception of some preprocessing to make the data compatible to the model, no significant pre- or post-processing steps are applied (i.e. multitemporal data or complex vector refinement, which are often employed in other approaches, are not considered in this thesis). Instead, the methodology focuses on a simple input-output mapping (fig. 6). Additional pre- or post-processing techniques could certainly add value later on, but would distract from the

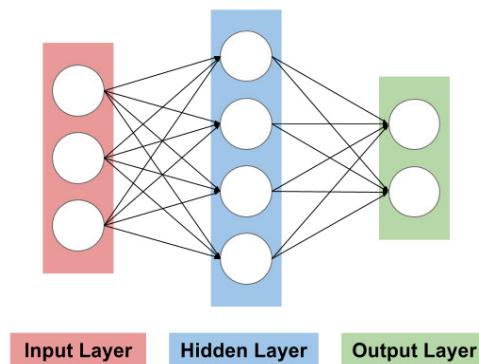
current goal of establishing baseline result. Also, when considering that the human brain is very successful at delineating field parcels even from simple RGB images, these additional layers of complexity should (in principle) not be necessary to yield promising results.

The thesis is structured as follows: Chapter 2 introduces the necessary deep learning basics to provide an understanding of the Instance segmentation task. It establishes the concepts of neural networks, image recognition via convolutional neural networks (CNN), semantic segmentation via fully convolutional neural networks (FCN) and object detection. Building on top of these concepts, the chapter then presents the functionality and state of the art of deep learning instance segmentation. Chapter 3 describes the Denmark study area, the training data and the methodology of preprocessing the satellite imagery and geospatial data. Also, the implementation and training characteristics of the FCIS model are introduced. Chapter 4 presents and discusses the results. Chapter 5 suggests possible additions and improvements over the current approach. Chapter 6 gives a conclusion.

## 2. Deep learning for instance segmentation

### 2.1. Neural network basics

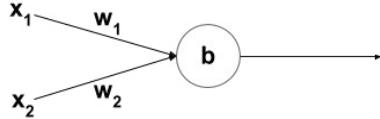
Artificial neural networks (ANNs) are information processing models, which are roughly based on the hierarchical and interconnected structure of biological neurons, the basic signal processing units in the brain. Usually, neural networks learn by example, which means supervised learning on training data configures the network to perform a specific task (e.g. image recognition). Neural networks can also be trained via unsupervised learning for data clustering applications, but these techniques are not considered here. ANNs don't require dedicated feature engineering, but can be trained on raw input data (e.g. image pixel values). These characteristics enable neural networks to perform complex prediction tasks, which are difficult to solve with simple rule-based techniques. The following section gives an introduction to the structure, functionality and training of neural networks; for a more comprehensive review see Nielsen (2015).



**Figure 7:** Simple neural network with input layer, one hidden layer and output layer.

Neurons are depicted by circles, neuron connections by arrows.

A neural network is often depicted as stacked layers of neurons (fig. 7). Layers between the input layer (which stores the data input, e.g. image pixel values), and output layer (storing the prediction result) are called hidden layers. When the network contains two or more hidden layers, it is usually called a deep neural network. Each neuron in a hidden layer is “fully connected” with all neurons in the adjacent layers via weighted connections. When performing a prediction (forward pass through the network), the network processes the input data through the network by carrying out a series of matrix operations. Each neuron receives the outgoing signals of all neurons in the previous layer, integrates and evaluates them, and passes on a new signal to the neurons in the next layer.



**Figure 8:** Perceptron.

The most basic form of a neural network is the perceptron (fig. 8), which has a single neuron. The perceptron's incoming signals  $x_1$  and  $x_2$  are multiplied with the respective connection weights  $w_1$  and  $w_2$ . Both products are summed and added to the neuron's bias  $b$ . The bias represents a threshold value of the neuron. The perceptron's output signal  $y$  is binary: 1 if the evaluated result of the neuron is positive, 0 if it is negative. More generally, the signal processing of a neuron in a multilayer neural network can be described by:

$$y = \sigma(\sum_i x_i w_i + b)$$

For each neuron in a hidden layer, the product of all incoming signals  $x_i$  and respective connection weights  $w_i$  is added to the neuron's bias  $b$ . Instead of a simple binary threshold, a nonlinear activation function  $\sigma$  is applied to the result (e.g. the sigmoid, tanh or ReLU function). A non-linear activation function is critical for the network to learn complex mappings between input and output. Using a linear activation function, the neural network can only learn linear mappings, no matter how many layers or neurons it has. After applying the activation function, the output signal  $y$  is passed on to the neurons in the next layer. By repeating this calculation for every neuron, the signals are processed throughout the layers of the network, which yields more and more abstract, internal representations of the original input data. The last network layer generates the prediction output (e. g. a distribution of class probabilities).

In general, all network parameters are initialized randomly. Therefore, the network is not able to make any meaningful predictions in the beginning. Training the network via supervised learning means to perform an initial prediction (forward pass), evaluate the prediction performance, tune the neuron's weights and biases (backward pass), and repeat the whole process iteratively until the network learns to predict the desired output for a specific input. As outlined, the network needs to evaluate how well the parameters have already been optimized during the training process. After each forward pass, the predicted output and the expected label are compared. This is done by a loss or cost function (e.g. mean squared error or cross entropy), which quantizes the quality of the result. Usually, the loss is averaged

over a number of samples in the training data set. The training then aims at minimizing the loss by updating the network parameters during each backward pass, via some form of gradient descent optimization. Recursively going back through the network layers, the gradient of the loss with respect to the network parameters (neuron weights and biases) is calculated using the computationally efficient backpropagation method. The gradient  $\partial L / \partial w$  is a vector in the high-dimensional parameter space of the network, pointing in the direction of the maximum increase of the loss. Consequently, in order to decrease the loss, the weights are updated by subtracting the gradient (i.e. moving the parameters into the opposite direction of the gradient in the parameter space):

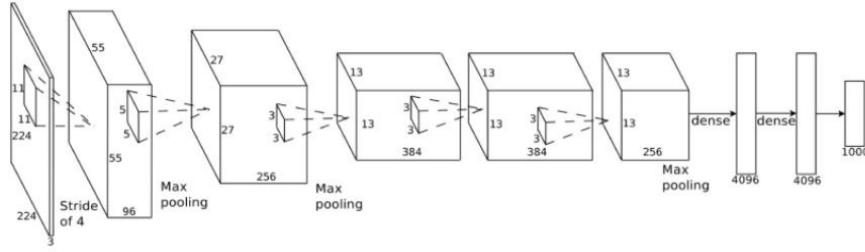
$$w' = w - \eta \frac{\partial L}{\partial w}$$

Here,  $\eta$  is the learning rate, a hyperparameter of the network that controls the magnitude of the parameter update. Besides this simple form of (stochastic) gradient descent, many modified optimization methods are used nowadays (e.g. RMSprop or Adam). Due to computational limitations, a forward-backward pass through the network does usually not contain all training samples, but is performed iteratively for mini batches of the dataset. Training for one epoch means that the the network has seen all training samples once.

## 2.2. Convolutional neural networks for image recognition

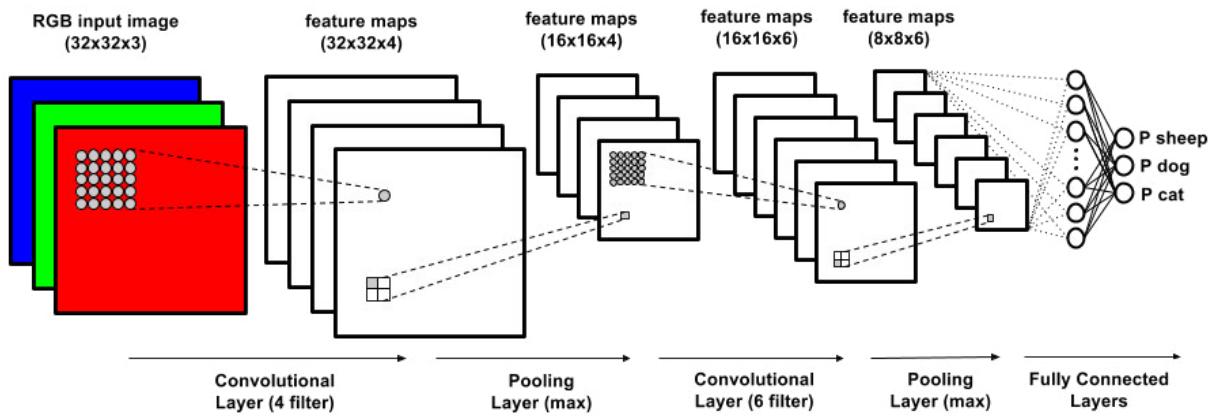
Although regular Neural Networks can perform simple image recognition tasks on input images of small dimensions, they are not efficient for processing image data at scale. Given a medium sized input image with dimensions of 256x256x3 (height, width, number of color channels), one fully connected neuron in the network's first hidden layer would already have 196608 weights. The huge number of parameters makes training the full network very inefficient and prone to overfitting.

A convolutional neural network (CNN or convnet) is a type of neural network that is more suitable for image data. Its architecture takes advantage of the spatial structure of the image, enabling more efficient training of the network's parameters. Initially, CNN architectures were mostly used for the task of image recognition. Although they natively lack the ability to perform spatial predictions, CNNs can be easily adapted and repurposed for more complex computer vision tasks like object detection, semantic segmentation or instance segmentation.



**Figure 9:** Architecture of the Alexnet network (Krizhevsky et al. 2012).

Pioneered by LeCun et al. (1988), the modern concept of CNNs was introduced and popularized by the Alexnet or Supervision model (Krizhevsky et al. 2012). Alexnet (fig. 9) appeared as the winning entry of the 2012 ImageNet Large-Scale Visual Recognition Challenge (ILSVRC). The challenge aims to predict the labels of 100,000 test images in 1000 categories, with 1.2 million images in the training dataset. Until 2012, the challenge was dominated by approaches leveraging Support Vector Machines. With an accuracy of 85%, Alexnet outperformed all other entries of 2012 and the previous years, beating the 2012 competition's second place by a huge margin of 11%. This prestigious win had an enormous impact on computer vision research and helped spark the deep learning revolution. Since 2012, each year's ILSVRC challenge winner has been a deep learning model.



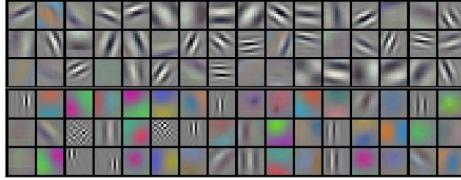
**Figure 10:** Schematic functionality of a convolutional neural network.

A typical CNN architecture for image recognition (fig. 10) takes an array of image pixels as input, extracts more and more abstract and computationally efficient image features, and finally outputs an array of class probabilities using a classifier on the most abstract feature vectors. Each feature extractor stage consists of a convolutional layer which applies multiple image filters to the data, each searching for a specific image pattern. The filters of the first convolutional layer look for these features in the input image, producing a filter activation

map or feature map. The filters of each following convolutional layer take as input the produced feature map of the previous convolutional layer. The filters of lower convolutional layers usually scan for simple, generic features (sharp edges of a specific orientation or color gradients). Combining the knowledge of the previous layers, later convolutional layers learn to register increasingly complex and abstract features (polygons, textures). With more and more abstract representations of the original input data, the network eventually builds up filters for whole image objects. The filter patterns are not fixed, but tune themselves during the network training process to recognize relevant image features. Before being passed on to the next convolutional layer, each feature map goes through an activation function (or non-linearity), and a pooling layer (reducing the parameter count) to make the prediction and training process more efficient. Based on the high-level features extracted by the last convolutional layer, the last stage of the network classifies the input image into one of the predefined categories. The various stages of the CNN (as depicted in fig. 10) are explained in more detail below.

**Convolution.** The convolutional operation can be described as a moving window that slides across the data. Its size is also called kernel size. Assuming a stride of 1 (i.e. the filter moves across the data in steps of one pixel) and zero-padding of 2 (i.e. the input image boundaries are extended by two rows and two columns of zeros), the moving window is applied to every possible spatial location of the RGB input image in the first network layer. Considering the CNN depicted in fig. 10, the image size is 32x32x3. The moving window filter of size [5x5x3] extends over the full depth of the input volume, i.e. it considers the information of all image bands at the same time. For each position in the input image, the moving window computes the filter activations via a dot product between the filter kernel and the raw image pixel values. The resulting 2-D feature map of size 32x32x1 is a representation of the original input image in terms of how strongly the filter was triggered by a particular image region. Closer matches with the filter pattern result in higher activations. A moving window stride bigger than 1, or no zero padding around the input volume yields a compressed feature map. This can result in a faster, but also less accurate CNN, because the feature maps contain less information. The convolutional operation can also be described via artificial neurons: The image filter is then represented by a grid of neurons, each connected to a small, localized pixel region in the input image, called the local receptive field. The filter pattern that each neuron is looking for (the convolutional kernel) is determined by the neuron's weights and biases. Every neuron in the grid shares the same weights and biases. This sharing of parameters enables the detection of a specific feature across the entire image. Each value in

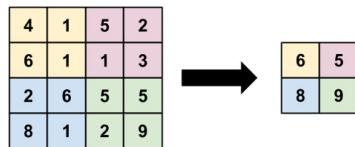
the feature map is calculated as the product between a neuron's weight and the RGB-pixel values in the image region that the neuron is connected to, plus the neuron's bias.



**Figure 11:** Visual representation of the 96 convolutional filter kernels learned by the first convolutional layer of the Alexnet CNN (Krizhevsky et al. 2012).

In practice, each convolutional layer uses multiple feature filters (fig. 11). Similar to the input image, each convolutional layer is arranged in 3 dimensions (width, height, depth). The width and height indicate the moving window dimensions or the arrangements of the neurons in each filter, the depth indicates the number of filters in the convolutional layer. Each filter takes into account the full depth of the 3D-input volume (original input image or stacked feature maps from a previous convolutional layer). The convolutional operation converts the 3D-input volume to a new 3D volume of activations. As an example, a convolutional layer with 4 filters is able to detect 4 different features across the whole input volume, producing an activation volume of  $32 \times 32 \times 4$ , no matter which depth the input volume has.

**Activation and pooling.** Before a feature map is processed by the next convolutional layer, it is usually passed through a non-linear activation function and downsampled via a pooling layer. The activation function, in this example a Rectified Linear Unit (ReLU) layer, sets all negative values in the feature map to 0. The dimensions of the feature maps array remain the same.



**Figure 12:** Max pooling operation (2x2, stride 2).

The most common operation performed by a pooling layer is max pooling (fig. 12). From each 2x2 square in a feature map, only the maximum value (the strongest activation) is kept. While the array depth remains the same, the array height and width is reduced to half (e.g. in the first feature extractor stage, the pooling layer decreases the array size from  $32 \times 32 \times 4$  to  $16 \times 16 \times 4$ ). The downsampling operation aggregates the activation information, while

discarding some detail of the spatial information. This makes the number of parameters and the computations in the network more manageable. After passing through multiple feature extraction stages, the original input image array is reduced into a much smaller volume.

**Classification.** In the classification stage of the CNN, the 3D-volume of high-level feature maps from the last convolutional layer (respectively last pooling layer) is transformed to a 1D volume of class scores, each representing the probability that the original input image belongs to that class. For that, the feature maps are vectorized, i.e. stretched out vertically and concatenated . By performing matrix multiplication of the neuron's weights and the connected filter activations, a number of fully-connected layers and a final softmax layer generate a dedicated score between 0 and 1 for each of these classes. The sum of all class scores is 1.

**Training.** Training of the CNN works similar as for a regular (feedforward) Neural Network (described in chapter 2.1). The CNN is trained via stochastic gradient descent. The aim is to reduce the error between the predicted and expected image classes. While the activation and pooling layers implement fixed functions that do not change during training, the weights and biases in the convolutional layers as well as the parameters in the classification stage are initialized randomly and tuned during the training process. After the initial forward pass through the network, the loss is calculated from the current prediction. Subsequently, in a backwards pass through the network, the gradients are calculated via backpropagation. Then, the network parameters are updated by subtracting the gradients. The parameters are adjusted iteratively, so that over time, the network is able to detect specific image patterns that are common in images of the desired classes. By tuning the weights of the fully connected layers, the network learns which of the extracted high-level features correlate most with which class.

Successfully training a CNN from scratch (with random initialization of all layer parameters) requires datasets of enormous size. Therefore, **Transfer Learning** is often employed to tune existing networks for new use cases. Here, the parameters of a different network, pretrained on a huge reference dataset within the same application category, are reused. Then, only the classification layers of the network are retrained, so that the network is able to predict the classes of the actual training dataset (equivalent to training a linear classifier on top of the extracted features). Depending on the dataset size and how much it differs from the training data for the initial network, it is also possible or required to finetune the network itself: **Finetuning** means to adjust the pretrained parameters of the last or even all convolutional

layers. Transfer learning and finetuning can speed up the training process enormously and enable learning on small datasets. They are especially successful in computer vision applications, because the image features recognized by the earlier feature extraction stages (e.g. edges or color gradients) are so generic that they fit almost any dataset.

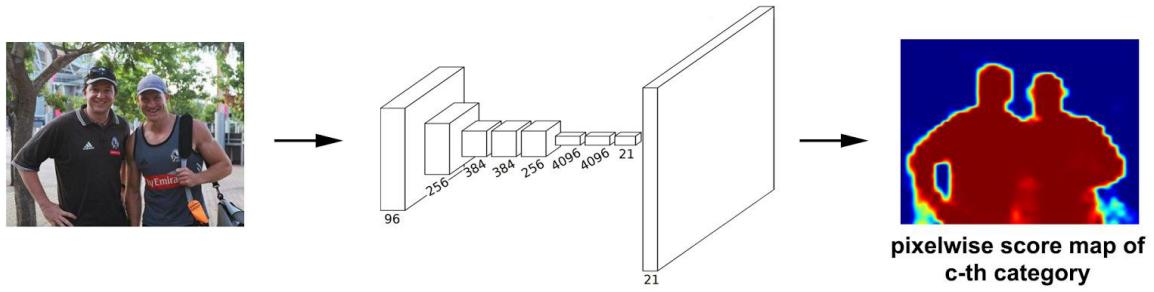
Following Alexnet (Krizhevsky et al. 2012), several popular, more complex CNN architectures emerged. **ZF Net** (Zeiler et al. 2013), winner of the 2013 ILSVRC challenge, is a modified version of the original Alexnet architecture with tweaked hyperparameters. **GoogLeNet** (Szegedy et al. 2014) won the 2014 ILSVRC challenge. It introduced average pooling and an Inception module, which enables the combination of pooling and multilevel convolutional operations at different filter sizes, while keeping the number of required model parameters comparably small. **VGGNet** (Simonyan et al. 2014) focused on a simple but powerful architecture (it only performs 3x3 convolution and 2x2 pooling operations throughout the entire network), with 13 convolutional layers and 3 fully connected layers. The model took second place in the 2014 ILSVRC competition. **Resnet** or Residual Network (He et al. 2015) introduced the concept of deep residual layers via “skip connections”. The connections between non-adjacent convolutional layers enable training this very deep network architecture. It won the ILSVRC 2015 challenge and is currently one of the most widely used models for transfer learning.

### 2.3. Fully convolutional neural networks for semantic segmentation

Image recognition models yield a single distribution of class probabilities for the entire image. In contrast, semantic segmentation requires a distribution of class probabilities for each image pixel (fig. 5). Conceptually, CNNs can be leveraged for pixelwise labeling in a patch-based approach via a moving window (CNNs require an input of fixed size). Then, the center pixel of each patch is assigned the predicted class scores for that image region. However, this approach is computationally expensive, as the CNN is applied on each image region independently.

The fully convolutional neural network (FCN) (Long et al. 2014) is based on the convolutional neural network architecture, but provides a more efficient solution for the prediction of pixelwise class labels (fig. 13). CNN architectures can be easily adapted to FCNs, and FCNs can still use transfer-learning with pre-trained CNN parameters. FCN introduces three major changes to the basic CNN architecture: 1) It replaces the CNN’s fully connected layers by convolutional layers. This enables the network to predict a class score map for each class for

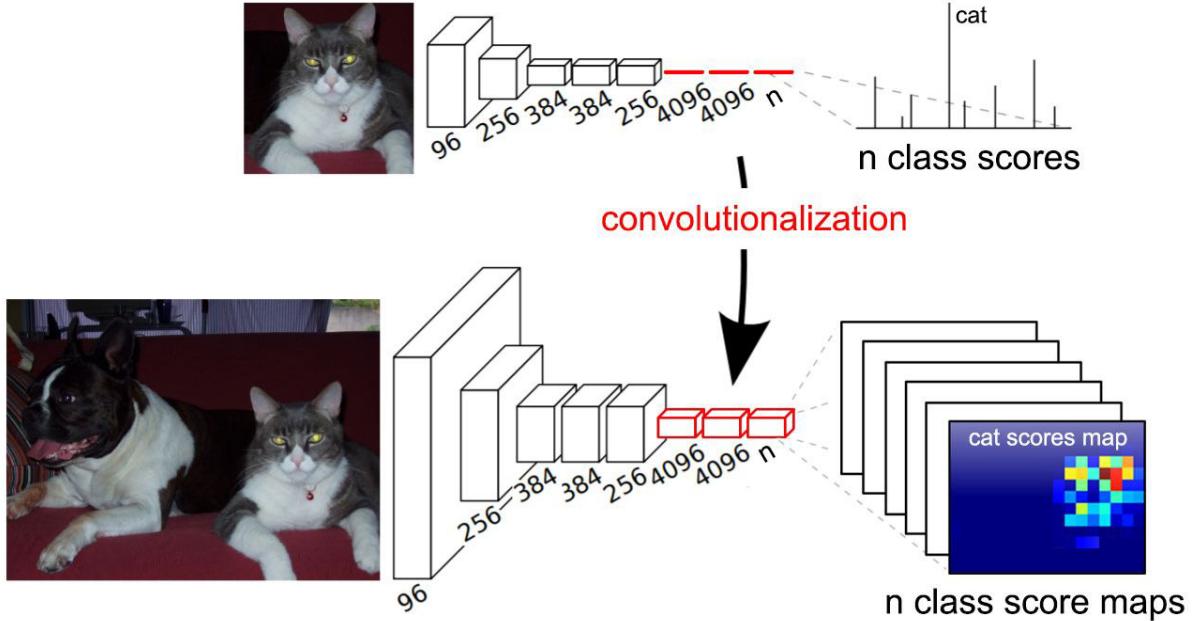
images of arbitrary size. The score maps of each class can then be merged into the desired “all-classes” score map by taking the maximum class probability of each pixel. Compared to the patch-based approach, the FCN method is a lot faster, because computations can be shared between overlapping patch regions. Due to the downsampling operations in the network architecture, the resulting output maps are coarse and require upsampling to the original image size. 2) Consequently, FCN introduces an in-network upsampling stage, the so-called deconvolutional layer. 3) The detail in the coarse class score map can be improved by adding “skip connections”, which incorporate less downsampled features from earlier convolutional layers into the final output. Some concepts of the FCN network are outlined in more detail below.



**Figure 13:** Class score map output and network architecture of FCN for semantic segmentation (edited from Long et al. 2014 and Dai et al. 2016).

**Convolutionalization.** Due to the fixed amount of neurons, a fully-connected layer requires an input of fixed size and outputs a 1D-array. Instead, the filter kernels of a convolutional layer can be applied to an input volume of arbitrary size, and produce a 3D-array of spatial maps. Ultimately, both a fully-connected and a convolutional layer compute weighted sums of their input values. This means that every fully connected layer can be replaced by a convolutional layer with a specific setup: The number of filters in the convolutional layer has to equal the number of neurons in the fully connected layer, and the receptive field has to be of the same size as the height and width dimensions of the input volume to the fully connected layer. E.g., take a fully connected layer with 4096 neurons, that processes feature maps of size  $8 \times 8 \times 256$  from the previous convolutional layer (i.e. 32 times downsampled from an input image of size  $256 \times 256$ ). This fully-connected layer can be replaced by a convolutional layer with 4096 filters (stride 1, padding 0) and a receptive field of size  $8 \times 8 \times 256$ . As the feature maps and filters have the same size, each filter will only be applied once, computing a single weighted sum. The 4096 resulting feature maps of size  $1 \times 1 \times 4096$  have the same number of network parameters and computational requirements as the fully

connected layer. For  $n$  classes, a final  $1 \times 1$  convolutional layer (trained via a pixelwise loss function) reduces the feature maps to  $n$  class scores.



**Figure 14:** Convolutionalization of fully-connected to convolutional layer. CNN outputting class scores (top) and CNN with convolutionalized layers outputting coarse spatial class score maps (bottom) (edited from Long et al. 2014).

By convolutionalizing the final, fully-connected layers in an image recognition CNN, the entire network can be treated as a series of filters. It can be applied to images of arbitrary size. This is the first main idea behind the FCN. When using the FCN on an image, the resulting output is not just a 1D-array containing one class score per class, but a 3D-array that contains, for each class, a spatial “heatmap” or 2D-array of pixelwise class scores (fig. 14). E.g. if the filter kernels of size  $8 \times 8 \times 256$  of the last convolutionalized layer look at feature maps of  $12 \times 12 \times 256$  (i.e. 32 times downsampled from an input image of size  $384 \times 384$ ), the filter kernels fit the feature maps five times in length and width, resulting in  $n$  spatial class score maps of size  $5 \times 5$ . All in all, the FCN maps the 3D input image to a 3D output of class score maps. This configuration is a lot faster than applying a CNN on every possible image location in a patch-based approach (Long et al. 2014), as the parameters of overlapping patch regions are shared in the network.

**Deconvolution.** Due to the downsampling operations in the network architecture, the class score maps are smaller than the input image size and exhibit limited spatial detail. They require upsampling to the original image resolution. Instead of fixed bilinear interpolation

(each upsampled pixel is the weighted average of the four diagonally neighbouring pixels in the input), the FCN network uses deconvolutional layers (also called transposed convolution, backward strided convolution, upconvolution or half strided convolution) for learnable, fast, in-network upsampling. The deconvolutional layer at the end of the network performs the upsampling of the class score maps by carrying out the inverse operation of a convolutional layer. An upsampling of factor  $f$  by the deconvolutional layer with stride  $f$  corresponds to a convolutional operation with stride  $1/f$ .

**Skip connections.** A more detailed or denser segmentation can be achieved by incorporating the higher resolution feature maps of earlier convolutional layers via “skip connections”. In addition to the feature maps of the last convolutional layer, the feature maps of the two previous convolutional layers (downsampled by factors of 8, 16 and 32 respectively) are directly forwarded to the final  $1 \times 1$  convolution classification stage. For each class, this yields three class score maps of different size. The score maps resulting from the feature maps of the earlier convolutional layers are of higher spatial resolution, but were generated from less abstract features, and thus contain “worse” semantic information. The class score maps are then upsampled (by their respective downsampling factor) to the original image resolution via deconvolutional layers. For each class, summing the respective three score maps yields a final output class score map with good semantic information and detail.

While most subsequent semantic segmentation models adopted the paradigm of fully convolutional networks, they introduced even more sophisticated approaches to improve the spatial detail in the class score maps (e.g. Badrinarayanan et al. 2015, Yu & Koltun 2015, Chen et al. 2017).

## 2.4. Object detection via Regional CNN

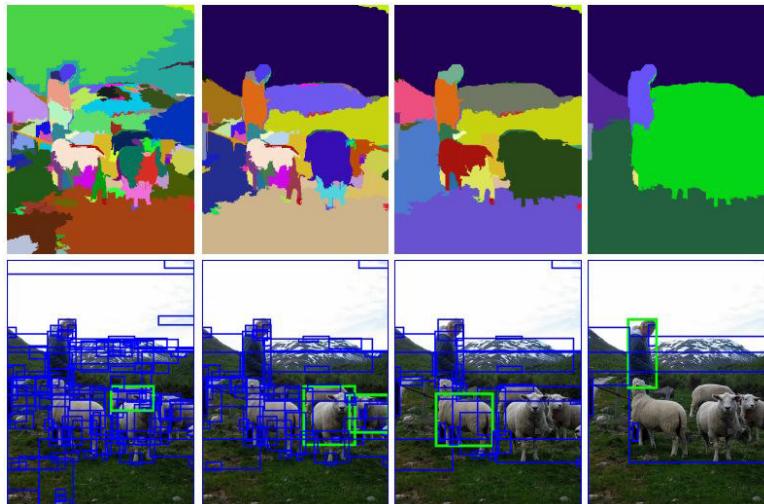
Object detection aims to find the bounding box locations and classes of image objects. The exact amount of objects is unknown. Conceptually, object detection is enabled by proposing a number of rectangular box regions in the input image to see if any of them correspond to actual image objects. This could be done by proposing boxes at every possible location and scale, and then separately applying a CNN on the image content of each box. However, this approach would be extremely performance-intensive.

R-CNN (Regional CNN) (Girshick et al. 2014) enables more efficient bounding box object detection with convolutional neural networks. Furthermore, the derivatives Fast-RCNN

(Girshick 2015) and Faster-RCNN (Ren et al. 2016) greatly enhance the model training and testing speed. At the time of writing, Faster-RCNN is the leading framework for object detection, and also serves as the basis for many instance segmentation models.

To evaluate models for object detection (as well as instance segmentation), there are two important metrics: The average precision (AP) is the proportion of correctly delineated bounding box or segment instances. The criterion for correct delineation is based on the Intersection-over-Union (IoU) ratio. This is the area of intersection between predicted and ground truth bounding box, divided by the area of union. See chapter 3.4 for a more detailed description of these evaluation metrics.

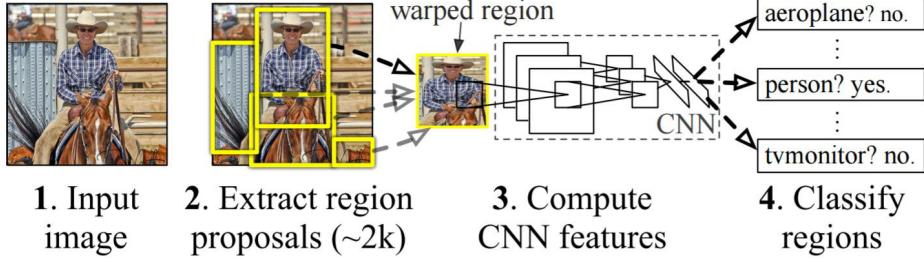
**R-CNN with Selective Search.** R-CNN first searches for a manageable number of class agnostic region proposals, i.e. a set of approximate image regions which are most likely to contain image objects. To find these regions, R-CNN uses the Selective Search algorithm (Uijlings et al. 2012), but is also compatible with other proposal methods. Selective Search greedily merges adjacent superpixels that share textures, colors, or intensities at multiple scales (fig. 15). R-CNN considers only the bounding boxes of these proposals.



**Figure 15:** Superpixels (top) and corresponding bounding boxes (bottom) of Selective Search applied at different scales. The green bounding boxes corresponding to entire image objects show the necessity of multiple scales (Uijlings et al. 2012).

The image content of each region proposal bounding box is then warped to a standard square size to fit the CNN specifications. The CNN is separately applied to each of these image snippets to extract their convolutional features (fig. 16). For each class (plus a dedicated background class), a linear SVM that is specifically trained on convolutional features of objects of that class is applied to the convolutional feature vector of each

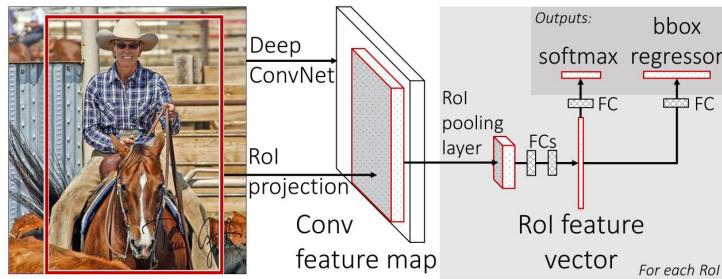
proposal. For each region proposal or region of interest (RoI), this yields a class score for each of the  $c+1$  predefined classes. The highest class score determines the proposal's class label and is also used as a detection score, indicating the confidence in the detection.



**Figure 16:** R-CNN stages (Girshick et al. 2014).

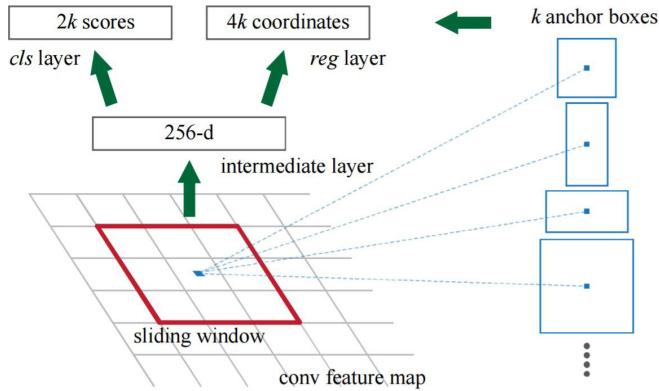
Redundant proposals, for which an overlapping, higher scoring proposal of the same class exists are then eliminated via non-maximum suppression (NMS). This method loops through every region proposal and compares it with all remaining region proposals of the same class. If region proposal A and another proposal B have a significant IoU overlap and the detection score of A is lower than that of B, region proposal A is rejected and removed from the loop. The remaining proposals are considered as the actual predicted image objects. Finally, to reduce localization errors and ensure tighter fitting bounding boxes, a class-specific, simple linear regression is applied on the remaining, classified region proposal. This step mildly corrects their bounding box coordinates.

**Fast R-CNN.** In the multi-stage pipeline of R-CNN, the CNN is applied separately on every region proposal. Because many region proposals are partially overlapping, R-CNN performs a lot of redundant calculations. The main achievement of Fast-RCNN (Girshick 2015) is the introduction of Region of Interest Pooling (RoIPool) for more efficient feature extraction. The classification stage of the CNN is removed, and the network is applied only once to the full image. Then, the feature vectors of the region proposals are extracted via RoIPooling directly at the corresponding locations on the output feature maps of the last convolutional layer (fig. 17). Subsequently, the network performs classification and bounding box regression in parallel. RoIPool is much quicker than applying the CNN multiple times and on partially overlapping regions. Instead of training and applying three different models separately (CNN feature extraction, SVM classification, bounding box regression), Fast R-CNN combines these tasks into one joint network with shared parameters, making the network much more efficient. The training uses a multi-task loss and updates all network layers jointly.



**Figure 17:** Fast-RCNN architecture (Girshick 2015).

**Faster R-CNN.** The final iteration of R-CNN, Faster R-CNN (Ren et al. 2016) focuses on accelerating the region proposal step, introducing the fully convolutional Region Proposal Network (RPN) in replacement of the fairly slow, external Selective Search algorithm. The class agnostic region proposals are instead generated directly on the convolutional feature maps, which is much more efficient and enables end-to-end training of the whole object detection model. While moving a sliding window over the feature maps of the whole image, each window position is checked with predefined anchor boxes (reference bounding boxes of common sizes and aspect ratios, by default 3 anchor boxes at 3 different scales) (fig. 18). For each anchor box, the RPN then regresses the proposal box coordinates (reg layer) and an “objectness” score via two parallel, convolutionalized fully-connected layers. The objectness score indicates how likely the box contains an image object or belongs to the background category. The region proposals are ranked by their objectness scores, and non-maximum suppression is applied with an IoU score of 0.7. From the remaining proposals, the top  $n$  ranked region proposals are used as the candidates for the object classification and bounding-box regression branches of Fast-RCNN. By using anchor boxes, fewer candidates need to be evaluated. Also, because the convolutional feature maps are shared between the RPN and Fast R-CNN, the proposals can be created at minimal computational cost.



**Figure 18:** Region Proposal Network (RPN) of Faster-RCNN (Ren et al. 2016).

## 2.5. Instance segmentation

### 2.5.1. Overview of task and models

Instance-aware semantic segmentation (or simply instance segmentation) means segmenting and classifying individual object instances. It combines elements of object detection (predicting bounding boxes and classes of objects without segmenting them) and semantic segmentation (labeling each image pixel with a semantic category label without distinguishing between objects of the same category).

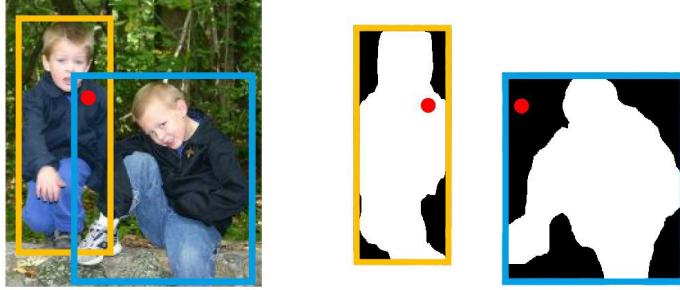


**Figure 19:** Mask-RCNN results on the COCO test dataset. Displayed are object instances, bounding boxes and object class labels (He et al. 2017).

Two of the most widely used reference datasets for instance segmentation are Microsoft Common Objects in Context (COCO) (80 classes, 300 k images) (Lin et al. 2014) and Pascal Visual Object Classes (VOC) (20 classes, 5 k images) (Everingham et al. 2010). The currently best performing model for instance segmentation with deep learning, Mask-RCNN (He et al. 2017), achieves 37.1 mean average precision (mAP) at [.5,.95] IoU (the average AP over all classes, averaged over IoU thresholds from 0.5 to 0.95 in 0.05 steps) on the COCO reference dataset (fig. 19). Second is FCIS (Li et al. 2017) with 29.2 mAP@[.5,.95] IoU. The FCIS model is used in the practical part of this thesis (and will be explained in further detail in chapter 2.5.4).

Instance Segmentation needs to operate on a region level, as the same image pixel can belong to different semantic categories depending on the focused image region or object instance. For example, a pixel can belong to the object mask of an image object (mask foreground), but the same pixel can be outside the mask foreground of another image object

(mask background) (fig. 20). This situation can not be distinguished by regular CNNs and FCNs when applied to the whole image, because convolution is translation invariant: A pixel will receive the same foreground/background scores, regardless of the relative position of the pixel in the focused context. Inspired by the success of region proposal techniques for object detection (see chapter 2.4), most instance segmentation models use per-RoI application via segment proposal.



**Figure 20:** Different semantic categories for the same pixel when focusing different instances. The same pixel (red point) belongs to the mask foreground (white) in the yellow instance, but the mask background (black) in the blue instance (edited from Li et al. 2016).

Earlier instance segmentation models (SDS: Hariharan et al. 2014, Hypercolumn: Hariharan et al. 2015, CFM: Dai et al. 2015a) relied on existing bottom-up segment proposal algorithms that extract and merge superpixels, e.g. Selective Search (Uijlings et al. 2012) or Multiscale Combinatorial Grouping (Arbelaez et al. 2014). Similar to an object detection workflow, the category-agnostic segment proposals are then classified with variants of R-CNN or Fast R-CNN (see chapter 2.3). The simple proposal techniques were later replaced by neural networks that learn to predict more accurate segment proposals, e.g. Deepmask (Pinheiro et al. 2015) or the fully-convolutional approach InstanceFCN (Dai et al. 2016). However, these segment proposals networks require a downstream classification network, which does not allow for parameter sharing. Multipathnet (Zagoruyko et al. 2016) uses proposals from Deepmask and utilizes a modified version of Fast-RCNN for more precise object localization. Multi-task Network Cascades (MNC) (Dai et al. 2015b) represents the first end-to-end trainable instance segmentation solution. It first predicts bounding box proposals, which are used subsequently to regress segment masks. Both bounding boxes and segment proposals are then used for classification via Faster-RCNN. Although MNC shares convolutional features between its subtasks, each task is still dependent on the output of the previous subnetwork, resulting in a complex network architecture and training.

Instead of successive segmentation and classification, recent models try to combine these tasks into parallel network architectures. FCIS (Li et al. 2017) is the first end-to-end trainable, fully convolutional network architecture for instance segmentation. It extends upon the ideas of position-sensitive class score maps for fully convolutional segment proposal prediction (Dai et al. 2016a) and object detection (Dai et al. 2016b). FCIS performs the object segmentation and classification simultaneously and efficiently, which makes the overall system very fast. Mask-RCNN (He et al. 2017), the current state of the art for instance segmentation, extends the object detection model Faster-RCNN (chapter 2.3) by adding a parallel, fully convolutional network branch for the prediction of RoI segmentation masks. Some of the most important models are described in more detail in the following chapters.

### 2.5.2. Deepmask and Multipathnet

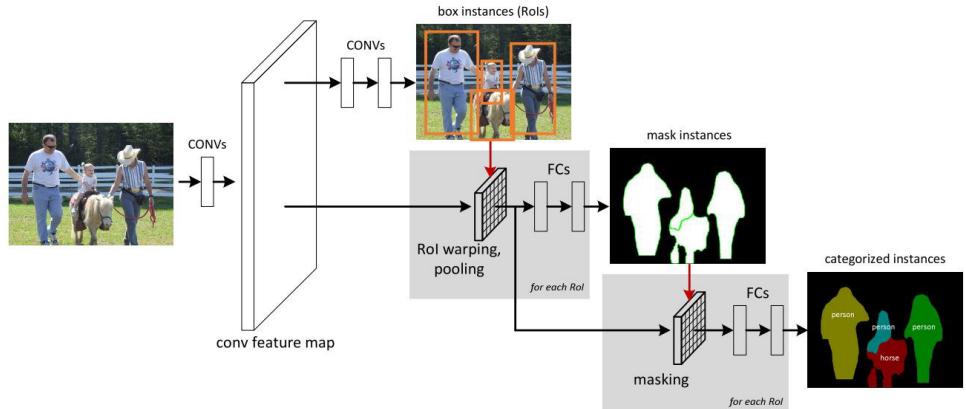
Deepmask (Pinheiro et al. 2015) introduced the first segment proposal method using convolutional neural networks. The architecture is based on the VGG network (Simonyan et al. 2014), a popular CNN for image recognition. The fully-connected layers of VGG are removed, and the convolutional features of the input image are extracted. The Deepmask architecture then splits into two network branches, which are applied directly to the convolutional feature maps of the last convolutional layer in a moving window approach at multiple scales. The first network branch predicts a class-agnostic, binary RoI segmentation mask. The second branch regresses a class-agnostic objectness score, indicating the likelihood that the RoI is centered on a full image object. Both tasks are trained jointly via stochastic gradient descent. The loss function is a sum of the logistic regression losses of both network branches. Because the segmentation mask is delineated from the downsampled convolutional feature maps, Deepmask bilinearly resamples the proposals to the original input image. The resulting proposals are relatively coarse and often suffer from imprecise alignment with the actual object boundaries. Sharpmask (Pinheiro et al. 2016) is an optional iteration of the Deepmask model, which aims to improve the accuracy and overall fit of the coarse segment proposals of Deepmask. It refines the proposal's alignment with the image object by incorporating information from less downsampled feature maps or earlier convolutional layers.

Multi-PathNet (Zagoruyko et al. 2016) identifies and classifies the segment proposals provided by Deepmask or Sharpmask. The architecture is based on the successful object detection model Fast R-CNN (Girshick et al 2015) with RoIPooling, but includes several modifications: While extracting the proposals' convolutional features via RoIPooling, the

network creates 4 “foveal regions” by cropping the corresponding feature map location at multiple scales. This enables the network to use the object context at multiple resolutions. Additionally, for each foveal region, it incorporates the more detailed information of earlier convolutional layers via skip connections. The network then classifies the segment proposals and applies bounding box regression to improve the localization of the object instances. Non-maximum suppression is applied to the proposals using the objectness scores provided by Deepmask to filter out overlapping predictions of the same object. Training is similar to Fast R-CNN, but the “integral loss” function used with Multipathnet applies the loss function at multiple IoU threshold values.

### 2.5.3. MNC

Multi-task Network Cascades (MNC) (Dai et al. 2015b), which won the COCO 2015 segmentation challenge, is a multi-stage structure of three consecutive VGG networks. The model is end-to-end trainable and shares parameters between its subtasks, but each task is still dependent on the output of the previous subnetwork.



**Figure 21:** Structure of the Multi-task Network Cascades (Dai et al. 2015b).

Firstly, the region proposal network of Faster-RCNN (chapter 2.3) predicts class agnostic bounding boxes and their respective objectness scores. Redundant proposals are filtered out via non-maximum suppression. The second stage takes the bounding box proposals and shared convolutional feature maps and extracts the feature vectors for the respective image regions via RoIPooling (fig. 21). For each bounding box, it predicts a class agnostic, pixel-level mask via binary logistic regression. Compared to Deepmask, which applies the segment proposal prediction at every image position, the preselection of RoIs via the bounding box prediction is computationally more efficient. The final classification stage takes the feature maps of the bounding box proposals and masks out the pixels inside the

bounding box that do not belong to the respective segment proposal (mask background; these values are set to zero). The resulting feature maps are limited to the foreground of the segment proposals, which are then used to predict the class scores. The entire workflow can be trained end-to-end via stochastic gradient descent, but the causal relation of the various outputs and subnetworks make the training via backpropagation non-trivial. As the loss of some network stages depends on other stages, the network is trained using a unified loss function.

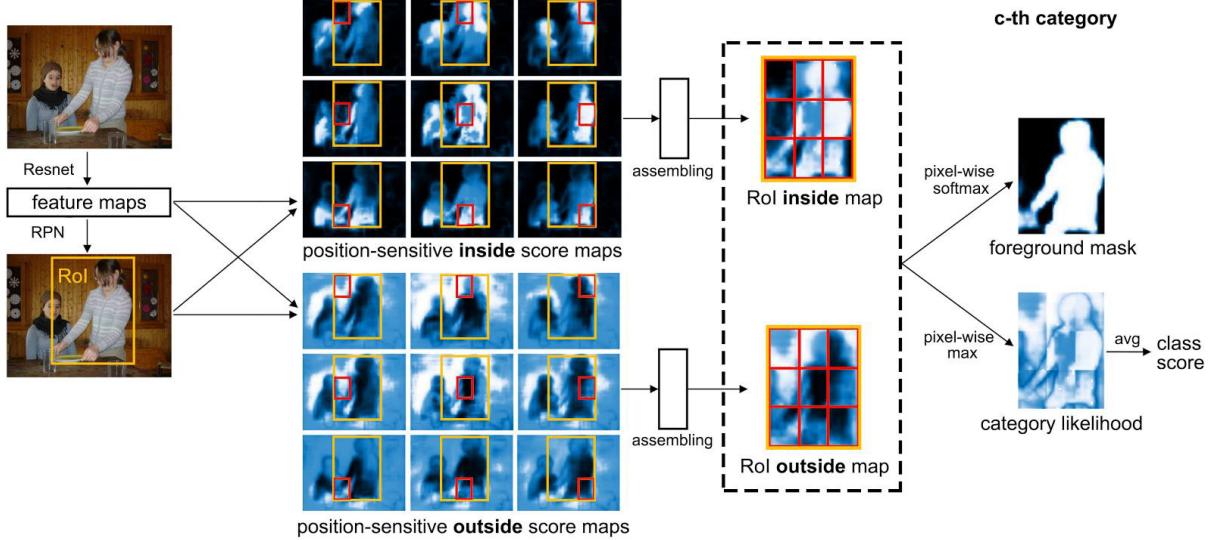
#### 2.5.4. FCIS

As described above, MNC shares convolutional features among its subnetworks. However, it does not share the per-RoI computation during the final classification step. Also, resizing the RoIs to a fixed size representation (which is required by the fully connected layers) can impair the spatial detail and segmentation accuracy. In contrast, fully convolutional neural networks (FCN, chapter 2.3) can be applied to images of arbitrary size and computations can be shared between overlapping RoIs. Therefore, they are much more suitable for spatial prediction. However, like CNNs, regular FCNs are translation-invariant.

A possible solution to this problem in instance segmentation is presented by the fully convolutional instance-aware semantic segmentation (FCIS) model (Li et al. 2017). It incorporates partially translation-variant, position-sensitive inside/outside class score maps. The model performs simultaneous detection and segmentation and shares parameters between sub-tasks as well as RoIs, which makes it a fast solution for instance segmentation. FCIS won the 2016 COCO segmentation challenge. It relies and extends upon the main ideas of InstanceFCN (Dai et al. 2016) for the fully convolutional prediction of segment proposals and R-FCN for fully convolutional object detection (Dai et al. 2016).

The FCIS model works in the following way: Firstly, the convolutional feature maps for the full image are extracted using the 100 convolutional layers of the ResNet-101 model (He et al. 2015) (the final pooling and fully connected layers are removed). As the downsampling factor of 32 in the last convolutional layer of ResNet produces too coarse feature maps for instance segmentation, the effective filter stride is reduced to 16 pixels. For each of the  $c+1$  classes (the additional class is dedicated to background), two sets of  $k^2$  position-sensitive inside/outside score maps are generated fully convolutionally. The  $2*k^2*(c+1)$  score maps represent the pixelwise probabilities that a pixel is at a specific relative position inside or outside an object instance of the particular category. E.g. in fig. 22, high pixel values in the top-middle inside score map indicate that these pixels have a high probability to be near the top-middle border area of image objects. A parallel RPN that shares the convolutional

features with FCIS proposes Rols. The Rol area (outer yellow rectangle, fig. 22) is split by a grid of  $k \times k$ , regularly distributed bins (inner red rectangles, fig. 22). For each Rol, one inside and one outside score map is assembled from the respective bins on the sets of inside/outside score maps of the full image. This makes the Rol score maps partially translation-variant: A pixel can receive different score values for different Rols if the pixel's relative position inside the Rols is different.



**Figure 22:** FCIS scheme (edited from Li et al. 2017).

For each class independently, the assembled Rol-specific inside and outside score map is then combined in two different ways: A pixel-wise softmax operation yields the class-specific foreground probability mask (the pixel has a high probability to belong to an object mask if it has a high Rol-inside and low Rol-outside score). A pixelwise maximum operation yields a class-specific category likelihood score map (high pixelwise probability to belong to the bounding box of an actual image object in case of combined low Rol-inside with high Rol-outside scores or low Rol-inside with low Rol-outside scores). The category likelihood map is then reduced to a single class score via average pooling. A softmax operation on all extracted class scores of the Rol results in a vector of  $c+1$  class probabilities. The highest class probability value indicates the class label of the Rol and acts as the confidence detection score. When a Rol bounding box only partially overlaps an image object, some of the position sensitive bins that are assembled for the Rol score maps are not activated, which leads to a low detection score. To filter out overlapping segment proposals, the Rols are ranked by their detection score, and non-maximum suppression with an IoU threshold of 0.3 is applied. For the remaining Rols, the foreground probability mask is selected by the class label and the final segment mask is determined by mask voting: For each remaining

RoI, the foreground masks of all RoIs with an IoU overlap of 0.5 or higher (including the ones removed during NMS) are selected. The masks are averaged pixelwise (weighted by the RoIs' detection scores) and then binarized.

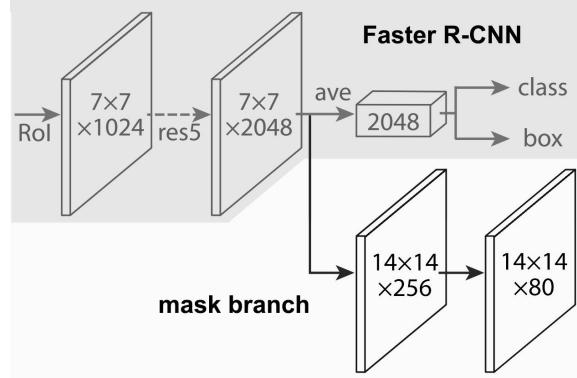
The introduction of partial translation-variance by the position sensitive score maps improved the AP@0.5 by 13.2 % (tested on the PASCAL VOC 2012 data) over using a translation-invariant baseline FCIS configuration with a single score map with 52.5 AP@0.5 (Li et al 2016). FCIS is trained via stochastic gradient descent with Online Hard Example Mining (OEHM) bootstrapping (Shrivastava et al. 2016). OEHM addresses the usual imbalance of many easy, but few hard examples (also due to the imbalance of positive and background RoIs). For each iteration, instead of randomly subsampling RoIs, the SGD modification automatically selects only high-loss examples, which significantly boosts the model accuracy with fewer iterations. The Resnet-101 convolutional layers are initialized with parameters pre-trained on ImageNet. All other learnable layers are initialized randomly. The shared score maps are generated fully convolutionally for the whole input image and correct aspect ratio. Each ROI is evaluated by the detection loss over  $c+1$  categories, the segmentation loss (sum of the pixelwise comparison of ROI foreground probabilities with the ground truth mask, normalized by ROI area) and the bounding box regression loss. All layers after the ROI and score map generation implement simple, fixed functions, enabling very efficient per-ROI computations.

### 2.5.5. Mask R-CNN

Mask R-CNN (He et al. 2017) is (at the time of writing) the state of the art instance segmentation model and outperforms all other solutions. It takes the object detection architecture of Faster-RCNN (including the region proposal network) and adds a fully convolutional network branch for the pixelwise prediction of class-agnostic, binary segmentation masks (fig. 23). The segmentation branch is parallel to the bounding box regression and classification branch of Faster R-CNN. All branches share the convolutional features of the input image and are applied to each ROI directly on the convolutional feature maps, which makes Mask-RCNN very efficient. The (multi-task) loss is the sum of the loss of each of the three network branches.

For each ROI, Mask-RCNN predicts a binary segmentation mask for each class fully convolutionally. The authors found that with conventional ROI Pooling of Faster R-CNN, the regions of interest in the convolutional feature maps do not completely align with the exact

location of the respective regions in the original image. For the classification branch of Mask R-CNN (and also for bounding box object detection in general), such small pixel misalignments are mostly irrelevant. However, instance segmentation requires pixel-level accuracy to delineate the exact object instances. Even small translations can have a big impact on the evaluation result. They are caused by rounding issues due to the different scales or pixel dimensions of the feature maps and input image. Therefore, Mask-RCNN replaces the RoIPool of Faster-RCNN with RoIAlign, which uses bilinear interpolation to exactly align the Rols with the corresponding regions in the original input image. RoIAlign greatly improves the mask accuracy by 10 % to 50 % depending on the evaluation metric's overlap threshold (He et al. 2017). To further improve the network accuracy, the convolutional feature extraction stage of Mask R-CNN uses the Feature Pyramid Network proposed by Lin et al. (2016), which enables efficient, in-network processing and exploitation of convolutional feature maps at multiple scales. Compared to FCIS, Mask-RCNN also increases the number of RPN anchors and evaluates the advanced Resnext-101 (Xie et al. 2016) instead of only Resnet-101.

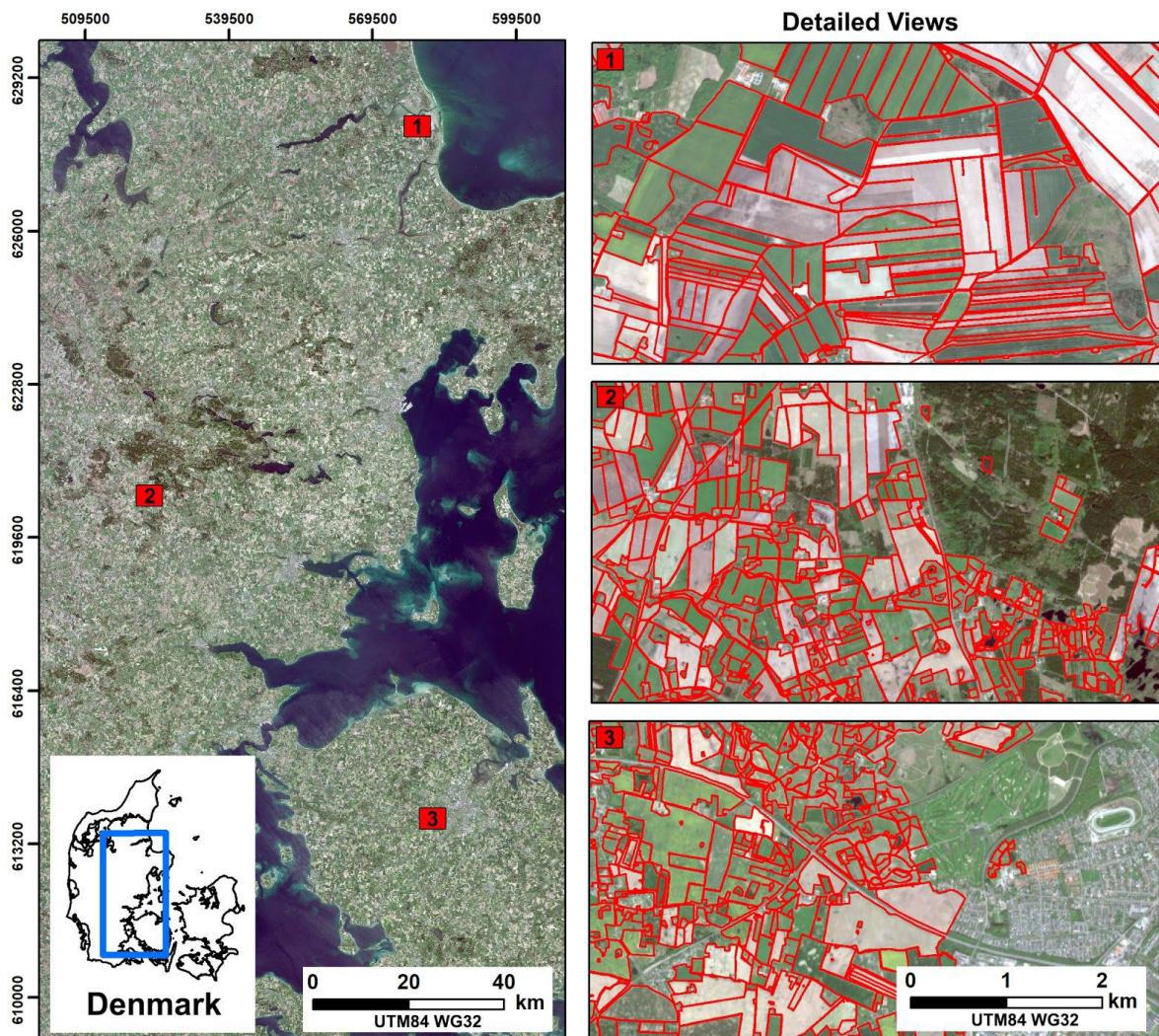


**Figure 23:** Mask-RCNN model with Faster-RCNN architecture (top) and parallel instance mask branch (bottom) (edited from He et al. 2017).

### 3. Methodology

The following chapter presents the training data, preprocessing, model training and evaluation methodology that was used for the automated delineation of field parcels. In short, medium resolution satellite images (Sentinel-2, RGB) and the corresponding, georeferenced field boundaries were preprocessed and cut into small image chips to fit the model requirements and available computational resources. Then, a fully convolutional neural network architecture, adapted from Li et al. (2016), was trained and tested in two configurations: 1) Segmentation of all agricultural field instances, 2) Segmentation of field instances and simultaneous classification of instance crop classes. Code will be made available at [github.com/chrisckri/InstanceSegmentation\\_Sentinel2](https://github.com/chrisckri/InstanceSegmentation_Sentinel2).

#### 3.1. Study area and datasets

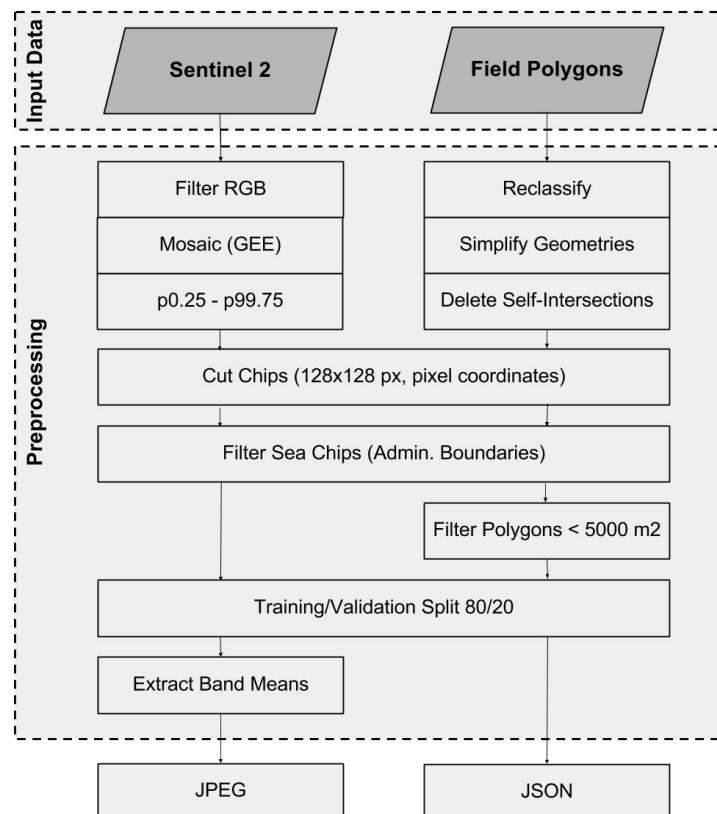


**Figure 24:** Denmark study area (left) and detailed views (right) showing the variety of agricultural parcel shape and geographical settings: (1) Large, regular field parcels in flat terrain, (2) Irregular shapes near forested areas (3), Urban setting with recreational areas.

The study area lies in central Denmark (fig. 24, left). It is defined by the extent of the Sentinel-2 satellite image mosaic and corresponding digitized agricultural parcels that were used as training and validation data. With an area of 20900 km<sup>2</sup> (190 km x 110 km), the region covers the eastern part of the Denmark mainland, adjoining parts of the Baltic Sea and the second biggest Danish island Funen. The area is characterized by dominant agricultural land use, temperate climate and mostly flat terrain.

The mosaic is made up of two cloud-free Top of Atmosphere (TOA) Sentinel-2 images of 10 m spatial resolution, captured at the beginning of the growing period in May 2016. The image tiles originate from the same Sentinel-2 image acquisition and overlap partially in north-south direction (full image metadata in table A1). The agricultural field polygons are from the 2016 Denmark ‘Marker’ dataset (MFLF 2017), which is part of the European Union Land Parcel Identification System (LPIS) initiative (Leo and Lemoine 2001). The dataset contains nearly 600 k field parcel polygons for all of Denmark, each attributed with a unique identification number, one of 293 crop type classes and the field area. 249298 parcels are contained in the study area, from which 159042 with a mean area of 52451 m<sup>2</sup> were used in this work (some classes are removed during preprocessing, see 3.2).

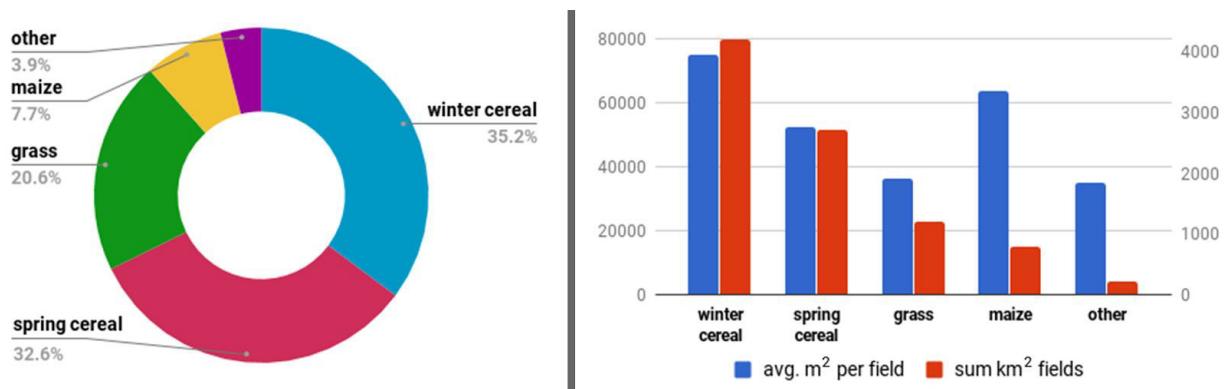
### 3.2. Preprocessing



**Figure 25:** Preprocessing workflow.

Fig. 25 shows a basic overview. The Sentinel-2 images were retrieved and mosaiced via Google Earth Engine (GEE), a cloud-based platform that offers access to geospatial datasets with analysis capabilities (Gorelick et al. 2017). All mosaic image bands except the red, green and blue channel (10 m spatial resolution) were discarded. The image values were clamped to their 0.25 to 97.75 percentile to preserve the dynamic range against outliers. Further preprocessing of image and polygon data was carried out with Python 2.7.8 (van Rossum 1995), mostly relying on the Python libraries rasterio, geopandas and shapely.

The parcel classes provided in the original dataset were reclassified to the five classes spring cereal, winter cereal, maize, grassfields and other fields. Parcels of classes not related to agricultural usage, e.g. natural or permanent environment, forest, recreational areas, or wasteland, were removed. The detailed reclassifications scheme can be found in table A2. The distribution of (reclassified) crop classes and parcel areas inside the study area is displayed in fig. 26 (detailed breakdown in table A3).



**Figure 26:** Distribution of (reclassified) parcel classes (left) and average area per class and sum of area per class (right).

After reclassification, the polygons were projected to the coordinate reference system of the image mosaic and intersected with the mosaic's bounding box. The remaining polygon geometries were simplified within a small distance tolerance of each point of the original geometry. The original topology was preserved. The original dataset contained polygons with geometry self-intersections, which were resolved by zero-distance buffering. The polygons were then transformed to pixel coordinates of the image mosaic.

The land area of the Sentinel-2 image mosaic was cut to 11320 chips of 128x128 pixels, each representing an area of 1.6384 km<sup>2</sup>. Chips outside of the high resolution Global Administrative Areas (GADM) Denmark boundaries (GADM 2016) were filtered out. The

agricultural fields dataset was intersected with the same grid and transformed to pixel coordinates of the 128x128 pixel chips (fig. 27).



**Figure 27:** Example 128x128 image chips with ground truth agricultural field polygons after preprocessing. Chips in the last row do not contain ground truth polygons.

Field polygons that overlap chip boundaries of the grid are split in at least two smaller polygons. As this can result in very small field instances, all polygons with an area below 5000 m<sup>2</sup> (corresponding to 50 image pixels) were filtered out. This resulted in 196704 field polygons with a mean area of 45583 m<sup>2</sup> (detailed breakdown in table A3). All image chips that contain at least one ground truth polygon were then randomly split into 80 % training data (8439 chips) and 20 % validation data (2110 chips). A secondary validation dataset contains the same validation chips in addition to all 771 “empty” image chips that do not contain any ground truth polygons.

To normalize the image chips to zero mean during training, the mean value of each image channel was calculated based on the training data. As all pixel values have the same range from 0 to 255 and similar relative scales, the standard deviation was not normalized. The image chips were then saved in RGB JPEG format without any compression. The polygon data was saved in the COCO JSON annotation format (Lin et al. 2014).

### **3.3. Model implementation details**

The instance segmentation task is carried out using the FCIS model (Li et al. 2016) (code release May 2016, <https://github.com/msracver/FCIS>). The architecture and functionality of FCIS is described in detail in chapter 2.5.4. Although Mask R-CNN (chapter 2.5.5) is the current state of the art instance segmentation model, at the time of writing the authors have not yet released their implementation (He et al. 2017).

The baseline FCIS model in combination with horizontal flip data augmentation was trained and evaluated on the training respectively validation satellite image chips (RGB channels, 128x128 pixel) and corresponding geodata of agricultural fields. Two different configurations were applied: 1) All fields treated as a single class (referred to as “single-class”), 2) each field annotated with one class label from the five reclassified crop classes (“multi-class”). The successfully trained algorithm takes in the validation satellite images and outputs instances of agricultural field parcels. The model implementation is based on MXNet (Chen et al. 2015), a lightweight, efficient deep learning framework. Training and evaluation were performed on a Tesla k80 graphics card of an Amazon Web Services (AWS) EC2 p2.xlarge instance. Following the considerations and recommendations for best practice deep learning instance segmentation in the literature, extensive testing of the internal FCIS-network parameters and training data properties (e.g. image chip dimensions, RGB vs. Near Infrared-R-G channels) was carried out. Eventually, most of the chosen model

hyperparameters follow the default settings of the original FCIS model implementation (Li et al. 2016) as they provide a good balance of prediction accuracy, training and inference time.

During model training, FCIS only considers image chips with existing ground truth instances, i.e. empty image chips are filtered out. The image chip values (range 0-255) are mean-centered by subtracting the per-channel mean values, which were extracted during the preprocessing. The order of the training data chips is randomly shuffled before each epoch, as a regular order could potentially bias the neuron parameters. Single-scale training is carried out, the shorter side of each image chip is resized to 600 pixel. Each image chip and the corresponding geodata is horizontally flipped during the training, the Roi specific results are averaged. This data augmentation helps making the model more accurate and robust to overfitting.

The model was trained end-to-end for 8 epochs, via stochastic gradient descent with OEHM bootstrapping (around 3 hours per epoch on the Tesla k80). Due to OHEM bootstrapping, each minibatch consists of a single image chip. The number of iterations is therefore equal to the number of training image chips times two (due to horizontal flipping). The SGD used a learning rate of 0.0005 and momentum of 0.9 (i.e. the parameter update takes into account previous update steps to avoid getting stuck in a local minimum). The learning rate was reduced to 0.00005 during the fifth epoch. Using transfer learning, the ResNet layers for the extraction of the convolutional features were initialized with the parameters of the ImageNet pre-trained ResNet-v1-101 model (He et al. 2015). Other learnable layers in the FCIS network were initialized with random values uniformly sampled between -0.1 and 0.1. During the forward pass, two sets of seven position-sensitive inside and outside score maps per category are generated. For shared convolutional features between FCIS and the RPN, a joint training scheme similar to Faster-RCNN (Ren et al. 2016) is used: The training alternates between tuning the region proposal and the detection/segmentation task. The FCIS loss is defined by the mask loss, detection loss and bounding box loss. The RPN uses twelve anchors (four scales at three aspect ratios). A Roi is considered a positive example when it has an IoU score greater 0.5 with the nearest ground truth object.

Inference time is about 0.25 s per image chip. Non-maximum suppression with an IoU threshold of 0.3 is applied to the Rols. For each of the remaining Rols, mask voting over all instances with an IoU overlap of 0.5 or higher (including the Rols removed during NMS) is applied. The pixelwise averaged masks (weighted by the Rols' detection scores) are binarized at a threshold of 0.5. Each final instance mask is encoded as a binary array in the dimensions of the input image chip. The binary masks are converted to polygon format via

Python opencv. The original topology is preserved. The resulting instance polygons and classes are exported to JSON format.

### 3.4. Evaluation metrics

The predicted field instances are evaluated by comparing them with the ground truth polygons in the validation dataset. The reported metrics include the average precision at different IoU thresholds (0.5, 0.75, [0.5:0.95]) and for three different object scales (S, M, L) as well as the corresponding recall values. These metrics are described in more detail below. The evaluation is carried out using the Python MS COCO API (Lin et al 2014, <https://github.com/pdollar/coco>).

For object detection and instance segmentation models, both the training and the evaluation stage require the metrics Intersection-over-Union (IoU) and average precision (AP): When comparing the position and appearance of a predicted object bounding box (object detection) or object segment mask (instance segmentation) with the respective ground truth instance, one needs to differentiate between a correctly predicted object (true positive or TP) and a falsely predicted object that does not exist in reality (false positive or FP). This differentiation is based on the IoU score of the two polygons. Commonly, an IoU score of 0.5 or above is considered a true positive. IoU, also called Jaccard Index, evaluates the similarity and dissimilarity of two polygons A and B. IoU is calculated by dividing the area of overlap between both objects by their area of union. For an IoU score of 1 the regions are exact matches, for a score of 0 they do not overlap at all. IoU is scale invariant.

$$IoU(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\text{Area Overlap}}{\text{Area Union}}$$

The instance segmentation method predicts a number of instances, each with a confidence value, the detection score. For each class independently, the instances are then ranked by their detection score. Each instance is compared with each ground truth object. Based on an IoU threshold (e.g. 0.5), each predicted instance is labeled as a TP or FP. For each ground-truth object, only the predicted instance with the highest detection score is labeled as the true positive: All other predicted instances that fulfill the IoU criterium but have lower

detection scores are labeled as false detections. If a true positive is found, both the ground truth and predicted object are removed from the loop. If a ground truth object is not matched by any predicted instance, it can be considered a false negative (FN). The number of false negatives results from the number of ground truth objects minus the number of true positives. From the TP/FP assignments, the precision and recall for each instance in the ranked list can be calculated. Precision (P) is the percentage of all predicted instances that match a ground truth object. Recall (R) is the percentage of ground truth objects that were correctly predicted. Each calculation takes into account the number of TP/FP/FN up to this point in the list, thus also includes the assignments from all prediction with higher detection scores. Per image chip, a maximum of 100 instances with the highest detection scores are considered. An example calculation of precision and recall for multiple predictions and ground truth objects can be found in table A4.

$$P_{class \ c} = \frac{\#TP(c)}{\#TP(c) + \#FP(c)}$$

$$R_{class \ c} = \frac{\#TP(c)}{\#TP(c) + \#FN(c)}$$

The precision values of all instances are averaged within 101 bins of 0.01 recall. When the resulting precision and respective recall values are plotted, the area under the resulting precision-recall curve gives the average precision (AP) for class c. To reach high AP values respectively have a high area under the curve, the method needs to have high recall and high precision. AP and AR are calculated for different IoU thresholds, e.g. 0.5 and 0.75. AP@[0.5:0.95], which is the standard COCO segmentation challenge metric, averages the AP at IoU thresholds from 0.5 to 0.95 in 0.05 steps. For the assessment of multi-class problems, the AP over all classes or mean average precision (mAP) is used. In addition to the overall AP and AR, the “per chip” AP and AR metrics are calculated over the predicted field instances of each chip independently. This enables a reasonable placement of the instance segmentation success for a particular image chip. Only instances with a detection score above the selected detection score threshold are taken into account. The per chip metrics can not be directly compared to the overall AP and AR metrics (as each chip contains a different amount of instances, the average of all per chip APs slightly deviates from the overall AP).

Analogous to the COCO dataset specifications and to make the results more comparable,

the evaluation and analysis focuses on chips with existing ground truth objects. For completeness, the results on a secondary validation dataset that includes all empty land area chips are briefly described as well. Given the spatial resolution of the Sentinel-2 images and the size of the agricultural field polygons, the area-specific COCO evaluation categories (Lin et al. 2014) do not correspond to the field delineation task. Instead, the S, M and L category area limits were selected by the distribution of field areas in the training data: S ( $< 16^2$  pixel), M ( $16^2 - 32^2$  pixel) and L ( $> 32^2$  pixel). This roughly corresponds to 43.1% of training field polygons  $< 0.025 \text{ km}^2$ , 46.5 % from  $0.025 \text{ km}^2 - 0.1 \text{ km}^2$  and 10.4 %  $> 0.1 \text{ km}^2$ .

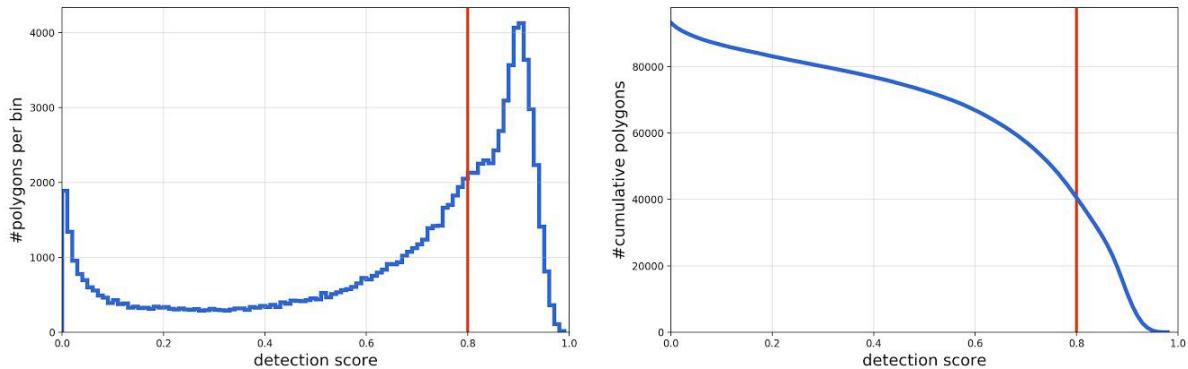
## 4. Results and evaluation

### 4.1. Experiment 1: Single-class

First, the results for instance segmentation in the single-class setting are reported. In this experiment, which is the main focus of this work, all agricultural fields were set to the same class index. Therefore, the sole task was to predict the polygons of all field instances in a given image patch, but not to classify them based on their crop type.

#### 4.1.1. Metrics

Overall, the validation dataset contains 2110 image chips. Running the model on all of these chips yields 100004 field proposals. On average, the model predicts 47.4 proposals per chip. This number is approximately normally distributed, with a minimum of 10 and a maximum of 76 polygons per chip. The confidence detection scores for these predictions range from 0.001 to 0.989. Fig. 28 visualizes the detection scores in more detail. The selected cutoff detection score of 0.8 (red line in fig. 28) leaves 42640 predicted field instances after filtering. In comparison, the validation dataset contains a similar number of 39160 ground truth polygons (an average of 18.6 polygons per chip).



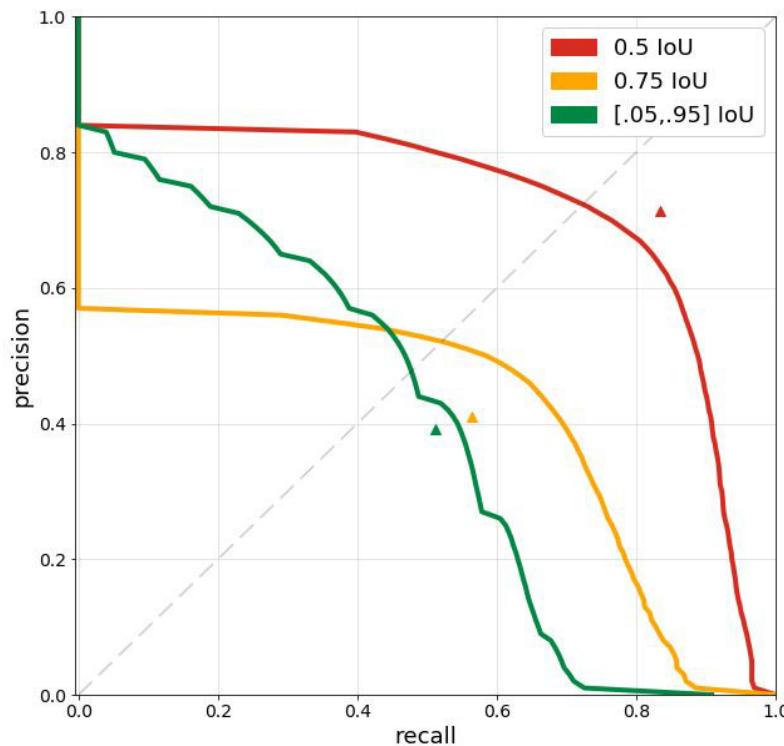
**Figure 28:** Histogram of detection scores vs. number of polygons (left) and number of polygons remaining at a specific detection score threshold (right). The red line is the detection cutoff score.

Table 1 summarizes the average precision (AP) and average recall (AR) at different IoU thresholds on the validation set. This dataset was not shown to the network during training. The metrics are explained in detail in chapter 3.4. As can be seen in the table, applying a low IoU threshold of 0.5 yields modest results: 71.4 % of the predicted field instances are correct (i.e. they coincide with actual field instances; precision), while 83.4 % of the ground truth

instances are successfully delineated (i.e. the model predicted a field instance at their location; recall). As expected, stricter IoU thresholds lead to disproportionately less successful delineations, because they require closer matches with the ground truth objects: Interestingly, AP and AR are very similar for IoU thresholds of 0.75 and [0.5:0.95], but substantially lower than for a threshold 0.5 (in relative numbers, AP and AR are reduced by around 40%).

**Table 1:** Average precision (AP) and average recall (AR) at different IoU thresholds for the single-class experiment (validation dataset). Metrics are calculated in two ways: By averaging over all field instances in all chips together (overall) and by calculating AP and AR for each image chip individually, then averaging these values over all image chips (per chip).

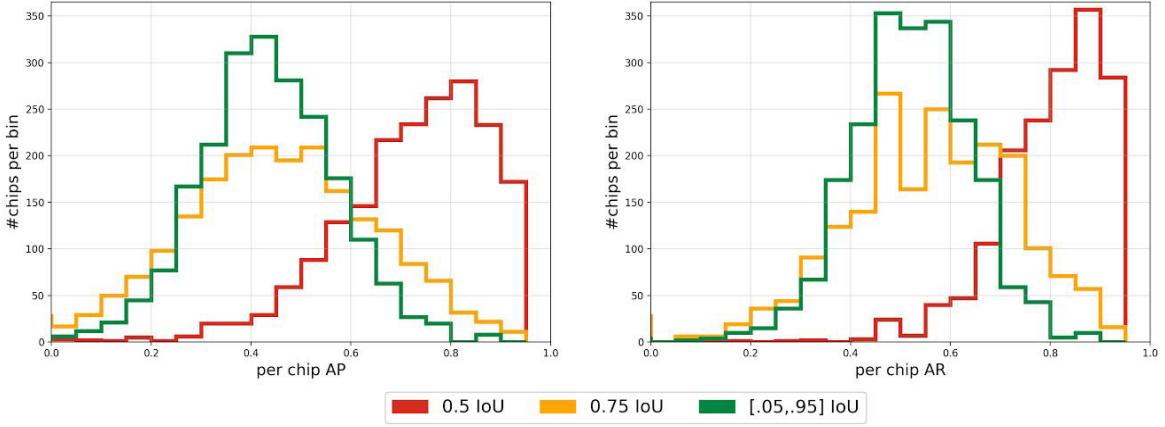
	AP@IoU			AR@IoU		
	0.5	0.75	0.5:0.95	0.5	0.75	0.5:0.95
<b>overall</b>	71.4	41.1	39.1	83.4	56.4	51.2
<b>per chip (mean)</b>	74.8	47.0	43.5	85.0	58.0	52.5
<b>per chip (standard deviation)</b>	16.3	20.9	13.8	12.5	19.1	12.1



**Figure 29:** Precision-recall curves at different IoU thresholds for the single-class experiment (validation dataset). Triangle symbols indicate the overall AP and AR values.

A general sense of the results can be gained by comparing them to the corresponding metrics on the COCO segmentation reference dataset: The AP of 39.1 at [0.5:0.95] IoU achieved in the single-class setting is in line with the performance on the COCO segmentation dataset evaluation by FCIS (29.2@[0.5:0.95]) and Mask-RCNN (37.1@[0.5:0.95]). However, note that the COCO dataset is much more complex and diverse than the field delineation task investigated here.

To gain a deeper understanding beyond the average metrics, the entire precision-recall curves (for the different IoU thresholds) are shown in fig. 29. The area under the curve corresponds to the AP of the predictions. To reach high AP values (i.e. high area under the curve), the method needs to have both high precision and high recall. The plot shows that the model requires a significant initial drop in precision to increase recall (for all IoU thresholds). In contrast, a perfect model would retain full precision until all ground truth objects are successfully delineated. After the initial drop off, the curves for AP at 0.5 and 0.75 IoU show a similar, balanced decrease of precision and recall. Although the APs at 0.75 and [0.5:0.95] have similar values, their precision-recall curves are significantly different: The averaging effects over the multiple IoU thresholds of AP at [0.5:0.95] IoU lead to a rippled curve pattern with a much stronger slope. Higher initial precision is balanced by lower recall.



**Figure 30:** Histograms of the number of chips with a specific per chip AP and per chip AR at different IoU thresholds (validation dataset).

The results presented so far were averaged over all field instances and image chips at the same time. In addition, the metrics (AP and AR) were calculated for each chip individually, in order to investigate how the model performs on the level of a single image chip. The mean and standard deviation of these per chip-values are also shown in table 1, they are consistently higher than the respective overall results. It should be noted that these values

are only calculated over the predictions of the specific image chip and can not be directly compared to the overall AP and AR metrics (see chapter 3.4.). Fig. 30 shows the distributions of AP and AR values on the chip level (i.e. how many image chips have a specific AP/AR value). The difference between low and high IoU thresholds can be observed clearly: For a low IoU threshold of 0.5, the distribution is right-skewed and limited to a certain range of relatively high per chip AP and AR values. In contrast, the distributions for IoU thresholds of 0.75 and [0.5:0.95] IoU approximately follow a normal distribution.

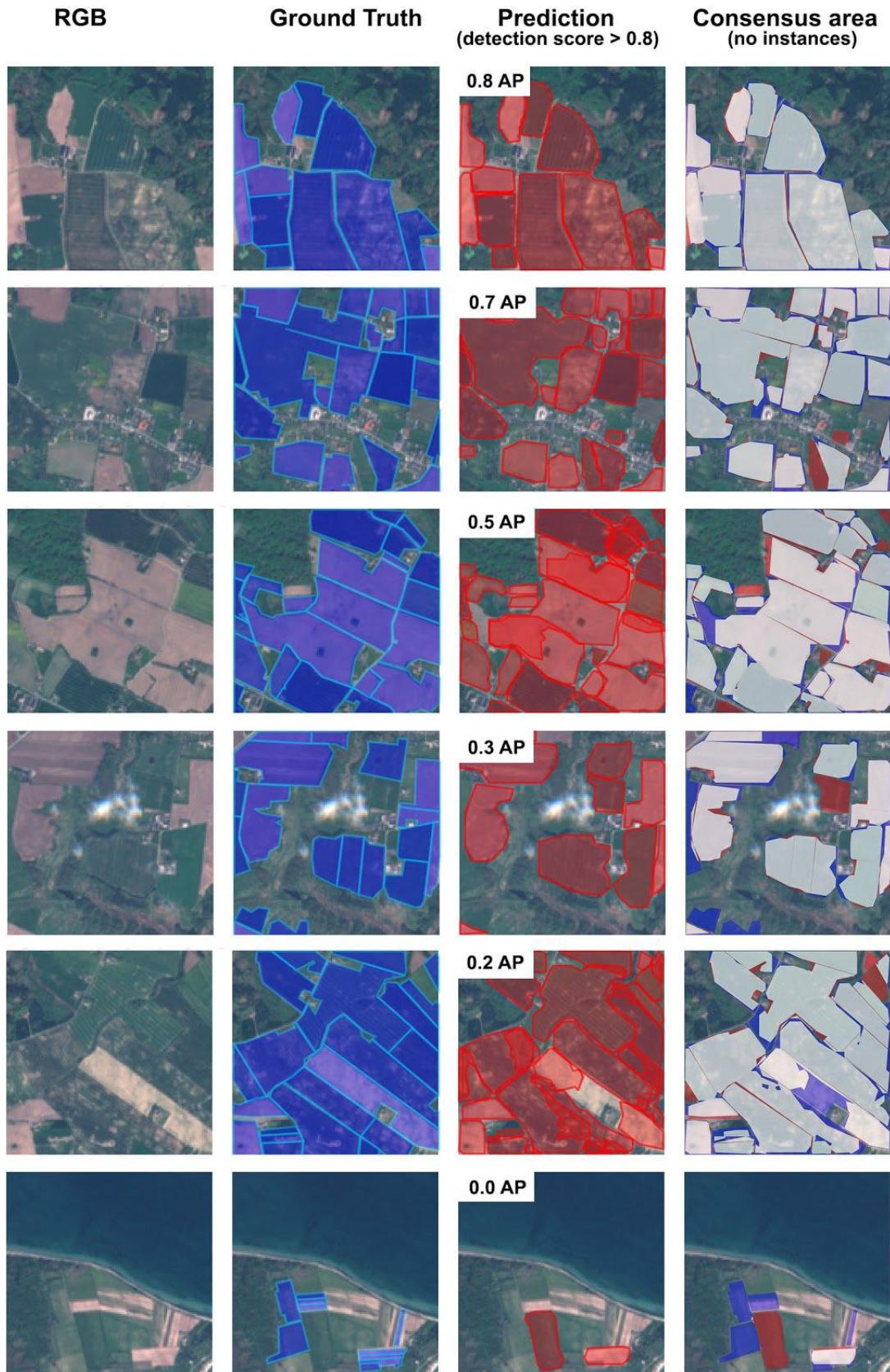
#### 4.1.2. Examples

Fig. 31 shows example image chips from the validation set. The ground truth field instances and the predictions by the model (with a detection score  $> 0.8$ ) are indicated. From top to bottom, the examples have decreasing per chip AP (all at IoU [0.5:0.95]).

For the examples with high per chip AP (top of fig. 31), the ground truth and the prediction results are in good accordance with each other in terms of the amount of field instances, their boundaries and the resulting overlap. Therefore, the consensus area (right column) closely resembles the actual field instances.

Complex shapes and borders are to some extent modelled by the predictions, even though small indentations and sharp edges often appear slightly rounded in the predictions. Also, straight boundaries of ground truth objects appear more irregular in the predictions. These effects hold for examples with high as well as low per chip AP. Heterogeneous texture inside single field parcels seem to be handled quite well. Sometimes, the predictions show some overlap with adjacent instances, i.e. the separation is not as strict as in the ground truth polygons.

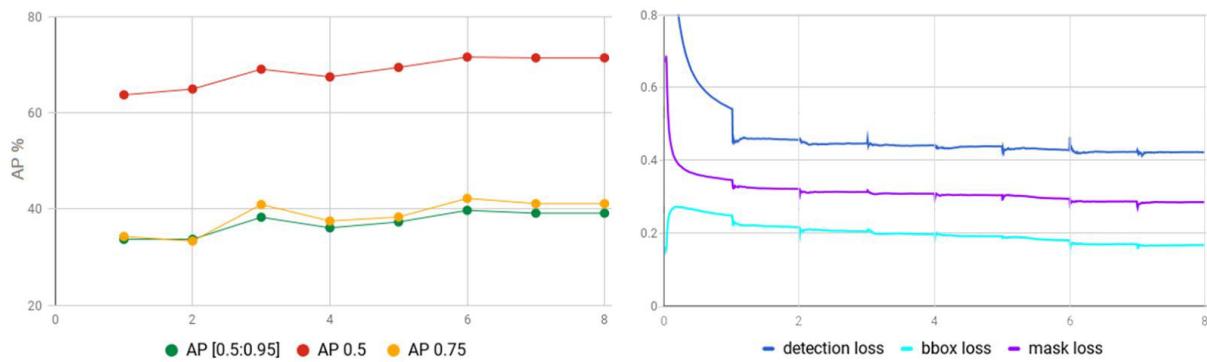
Image chips with low per chip AP (bottom of fig. 31) often show three main characteristics, which are discussed later on: 1) Predicted instances that cover an area which is actually covered by multiple small instances in the ground truth data, 2) full overlap of a smaller prediction by a bigger prediction, 3) false positives due to assumed field instances that are not existing in the ground truth data.



**Figure 31:** Example image chips from the validation dataset and predictions for the single-class setting (ranked by per chip AP@[0.5:0.95]). Only field instances with detection scores over 0.8 are displayed and considered for the per chip AP calculation. The consensus area (white) shows the overlap of ground truth and predictions.

#### 4.1.3. Training process

Fig. 32 displays the development of the AP metrics and loss functions during the training process. Due to transfer learning initialization, OHEM bootstrapping and the single image minibatch, the training losses (running average per epoch) quickly decrease within the first epoch. Afterwards, they slowly decreases further, until they plateau around epoch six to eight. The small peaks of the loss function at the beginning of each epoch are artifacts of the training process and insignificant for the performance of the model.



**Figure 32:** Development of the model over the training epochs: AP at different IoU thresholds (left) and loss functions (right).

#### 4.1.4. Field size

In manual segmentation, it is arguably harder to recognize and segment small agricultural fields. To investigate the effect of field size on the performance of the model, all field instances were assigned to one of three categories based on the size of their predicted polygon: S for small fields ( $< 0.025 \text{ km}^2$ ), M for medium-sized fields ( $0.025 \text{ km}^2 - 0.1 \text{ km}^2$ ) and L for large fields ( $> 0.1 \text{ km}^2$ ). Table 2 shows the AP and AR values for each category, as well as the distribution of fields across the categories.

As expected, the performance of the model increases with field size. Especially small field instances (category S) have a rather low AP and AR, while medium-sized and large fields show more similar results.. Also, with an amount of 55.9 % of the predicted field instances, unproportionally many small field objects are predicted - compared to only 43.1% of small objects in the ground truth training data.

A possible reason for these bad results on small field instances can be seen in fig. 31, in the example with 0.0 per chip AP: Apparently, the FCIS algorithm has problems to delineate field instances in congregations of small, directly adjacent fields, even if they show distinct

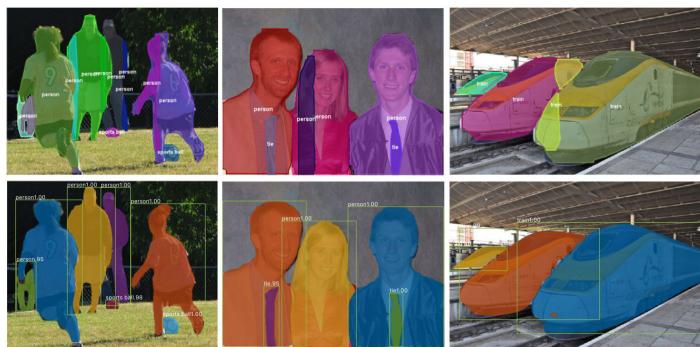
differences in texture and color . This may be related to the mask voting operation of FCIS in combination with the relatively low spatial resolution of the imagery. The segment masks of small, partially overlapping predictions could be merged unintentionally.

**Table 2:** Average precision (AP) and average recall (AR) by polygon area for the single-class experimental (validation dataset).

	AP@[0.5:0.95] IoU	AR@[0.5:0.95] IoU
<b>S</b> ( $< 0.025 \text{ km}^2$ , 55.9 %)	28.1	39.5
<b>M</b> ( $0.025 \text{ km}^2 - 0.1 \text{ km}^2$ , 32.1 %)	47.3	58.9
<b>L</b> ( $> 0.1 \text{ km}^2$ , 11.9 %)	51.0	68.6

#### 4.1.5. Overlap

The results also show some artifacts of strongly overlapping predictions. This situation is usually characterized by a small, redundant prediction that is completely (or to a large extent) overlapped by a bigger, often correctly delineated prediction. This seems to appear most frequently near the border of image objects. As the redundant predictions do usually not enclose a full image object, some of the position sensitive bins of the RoI-inside and outside score maps in the FCIS model should not be activated (see chapter 2.5.4). This would lead to a low detection score and the proposal being filtered out.



**Figure 33:** Systematic artifacts of redundant, overlapping predictions in FCIS results (top) compared to Mask R-CNN (bottom) (He et al. 2017).

A similar situation was observed by the authors of Mask R-CNN, when comparing prediction results from natural photos across models (fig. 33), “suggesting that it is challenged by the

fundamental difficulty of instance segmentation" (He et al. 2017). Neither the non-maximum suppression nor the mask voting operation of FCIS nor an additional NMS after the mask merging step are usually able to remove the redundant prediction, as two overlapping objects with distinctive size difference have a relatively low IoU score.

#### **4.1.6. Problematic training data**

The 2016 Denmark LPIS dataset, which was used as ground truth for the field polygons, generally matches the field structure in the satellite images that were used in this thesis. However, it was not specifically delineated from these images. Also, the main production period and exact production specifications of the dataset are not available. This causes some problems with model training and evaluation: When visually comparing the ground truth data and satellite images, definitive mismatches of declared or missing field instances can be spotted in some areas (e.g. fig. 31, last row). This is due to temporal changes before or after the reporting phase, missing or inaccurate field delineations or the inclusion of property boundaries in the dataset. Mismatches can also be caused by the challenging reclassification and selection of agricultural-related field classes due to the extensive original class legend and variability of some environmental classes (e.g. also containing a minority of field objects). These challenges influence both training as well as evaluation: Although the prediction might be reasonable for the actual image content, the evaluation of these instances leads to worse results or a high training loss.

#### **4.1.7. Empty image chips**

To guarantee the comparability of the results, in the validation dataset, empty image chips without any field instances were discarded (see chapter 3.4). This accounts for 771 land image chips in total. The scenes are mostly of urban, coastal and forest setting (fig. 27). Including these image chips in the evaluation slightly lowers the AP score by <0.01, but does not impair the performance of the algorithm in any significant way. For all empty image chips, only 47 additional instances with detection scores above 0.8 were predicted, distributed across 32 different image chips. These predictions count towards the overall number of false positives. However, when visualizing these false positives (fig. 34), it can be seen that most predictions are quite reasonable, and could potentially be confused by human annotators as well.



**Figure 34:** Example false positive predictions in image chips without ground truth instances.

## 4.2. Experiment 2: Multi-class

In the second experiment, the model was trained in a multi-class setting: Each field instance was assigned to one of five crop classes (compare 3.2). The task was now to predict the polygons of the field instances (as in experiment 1) and simultaneously classify their crop type.

### 4.2.1. Metrics

Table 3 summarizes the mean average precision (mAP) and mean average recall (mAR; averaged over all classes) at different IoU thresholds on the validation set. Table 4 reports the class-specific AP at [0.5:0.95] IoU for each of the five agricultural crop types.

Overall, the multi-class delineation was not as successful as the single-class setting: the (averaged) metrics are significantly lower than in the single-class experiment. Even with the low IoU threshold of 0.5, only 45.0 % of the predictions are correct (in contrast to 71.4 % for single-class). However, one has to note the different evaluation criteria in this experiment: False positives are not only field instances that can not be assigned to a ground truth object, but also instances that do not match the correct ground truth class label. Classification seems to be most successful for the winter and spring cereal classes, which are also the most frequent classes in the study area (each accounts for ~35 % of fields, fig. 26). As in the single-class experiment, applying the stricter IoU thresholds 0.75 and [0.5:0.95] leads to a significant drop of the mAP and mAR metrics (about 40 to 50 % in relative numbers). Overall, the network was apparently not able to handle the classification of field instances based on the available information in the RGB images.

**Table 3:** Mean average precision (mAP) and mean average recall (mAR) at different IoU thresholds for the multi-class experiment (validation dataset).

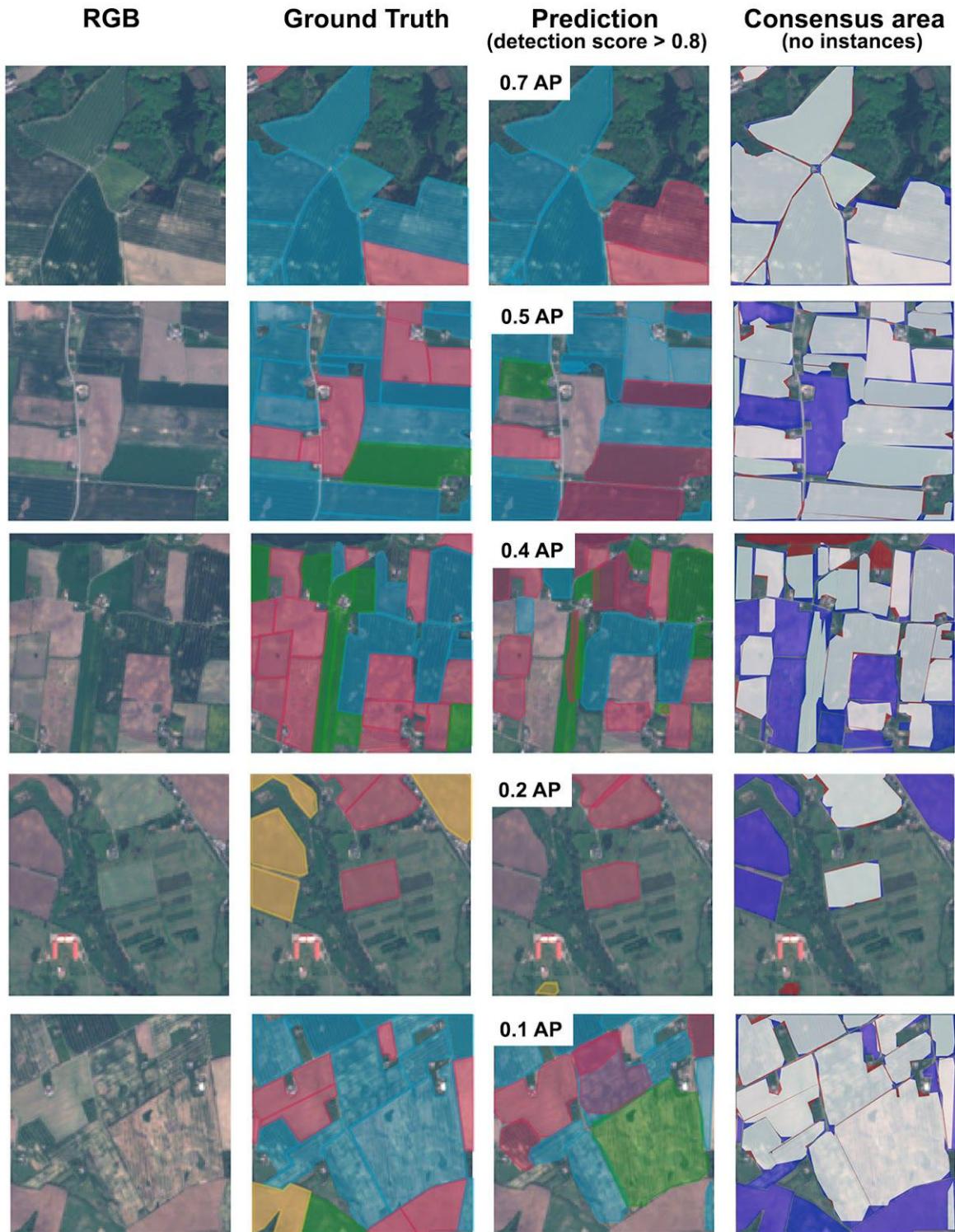
	AP@IoU			AR@IoU		
	0.50	0.75	0.5:0.95	0.50	0.75	0.5:0.95
overall	45.0	25.0	24.3	74.1	49.2	44.9

**Table 4:** Class-specific average precision at [0.5:0.95] IoU for the multi-class experimental (validation dataset).

	winter cereal	spring cereal	grassfields	maize	other
AP@[0.5:0.95]	44.0	32.2	22.4	19.1	3.6

#### 4.2.2. Examples

When comparing the results visually (fig. 36), the prediction of the polygons' shapes and layouts seems to match the performance of the single-class experiment. The predicted field instances closely resembles the ground-truth data, with some deviations particularly at low per chip AP. However, the crop classes are frequently misclassified.



**Figure 35:** Example image chips from the validation dataset and predictions for the multi-class setting (ranked by per chip AP@[0.5:0.95]). Only field instances with detection scores over 0.8 are displayed and considered for the per chip AP calculation. The consensus area (white) shows the overlap of ground truth and predictions.

## 5. Future work

The methodology presented in this thesis offers potential for many improvements. Additional research could greatly benefit the results achieved for agricultural fields segmentation, the applicability of the model to other scenarios and the general understanding of deep learning instance segmentation for the use with remote sensing imagery. The four main directions for future improvement are: 1) using more accurate and heterogeneous training data, 2) exploiting the satellite image information in a better way, 3) improving model training and 4) extending postprocessing.

**More accurate and heterogeneous training data.** The optimal training data for the field delineation task would represent the exact parcel situation at the time of the satellite image acquisition. Manual editing to make the data fit exactly to the selected satellite scenes would be time consuming, but could be viable to create general reference datasets. Also, in view of increasing public availability of high quality geo datasets, the use of additional and more heterogeneous ground truth and satellite image data (e.g. from multiple study areas in parallel) could improve the general model accuracy, robustness and transferability. Additional data augmentation (e.g. random scaling, distortions, color offset or jitter, introduction of image noise) might also help mitigate overfitting and increase the robustness to different geographic settings as well as atmospheric and lightning conditions. Although the model can generalize within certain limits (local generalization), its transferability is mainly limited by the variability of the training data. Furthermore, higher spatial resolution would certainly improve the exact delineation of the field boundaries, especially for small objects. At the time of writing, Sentinel-2 presents the highest resolution of freely available satellite imagery.

**Better exploitation of the satellite image information.** The exploitation of multitemporal satellite image information and the full spectrum of satellite image bands could certainly improve the results. The current FCIS implementation is limited to single-temporal three band imagery. The use of more image channels of different spectral wavelengths is mainly limited by computational resources. A workaround could also be the application of principal component analysis (PCA) on the full range of the available image bands and subsequent use of the first three dimensions in replacement of the RGB image channels. The proper integration of multitemporal information would instead require a different model architecture. The instance segmentation models that are currently available are designed for the use with natural photo images. Therefore, they do not consider a possible relation between several images. In addition to the raw image data, hand-crafted features could be integrated into the model of (e.g. band indices, edge detection preprocessing). This could possibly improve the

prediction accuracy, but somewhat opposes the deep learning premise of not requiring any feature engineering.

**Modified model training.** Various modifications of the training process and parameters could potentially improve the predictive power of the model. Often, this kind of improvement requires a tradeoff between hyperparameter tuning effort, training time and potential accuracy gain. Examples for future work are the use of additional RPN anchor scales and aspect ratios, additional training iterations, adjustments of the learning rate and multi-scale training. Model ensembles, which averages the results of multiple instances of the same model architecture applied on different subsets of the training data and with different initializations, can also increase the predictive performance. Higher computational costs by increased training complexity could be diminished by parallelizing the model on hardware with multiple GPU cores. Also, the pending implementation release of the current state of the art instance segmentation model Mask R-CNN (He et al. 2017) will enable access to an even more performant and simpler architecture.

**Improved postprocessing.** A multitude of external polygon postprocessing techniques could potentially improve the fit of the field boundaries (e.g. snake algorithms for directed line smoothing) or reduce the number of false positive predictions (e.g. removal of strongly overlapping polygons by modified NMS or overlap criteria). As the field parcel objects can not show any direct overlap in the satellite image, any partially overlapping subareas of the predicted polygons could be removed.

An important aspect of the current methodology is that its use in an operational context for regionwide delineation of field parcels is limited, as the use of instance segmentation with satellite images instead of natural photos causes a unique challenge: The satellite images are cut into smaller image chips because the method is limited by the GPU memory and requires multiple examples to iterate on. The chip-specific predictions would need to be stitched together. However, fields that are not completely enclosed within the boundaries of a single image chip are split into at least two polygon predictions. A simple combination of these parcel sections is not easily possible. One potential workaround could be to additionally use intermediate chips that correspond to a 25% area overlap of four adjacent image chips. The intelligent combination of the polygon predictions near the chip center could minimize the risk of predicting partial instances.

## 6. Conclusion

This thesis investigates the potential of deep learning instance segmentation for the automatic delineation of agricultural field parcels from satellite images. Considering the relatively basic input-output approach (without extensive postprocessing), the non-optimal training data, and the small amount of available precedent work that deals with instance segmentation of remote sensing images, the results look very promising. Major challenges for the existing approach are the correct prediction of small field parcels, miss detections of parcels considered as environmental areas as well as similar physical characteristics between parcels of different crop classes.

Considering the positive results obtained in this thesis, the application of deep learning instance segmentation to remote sensing and geospatial data shows a strong potential for the future: The models are able to use raw data and do not require manual feature engineering. Compared to many other image segmentation algorithms, applying the model at test time does not take much computation time, which hands itself well to the rapidly increasing volumes of satellite imagery and geodata. However, the predictive performance strongly relies on the accuracy, volume and variability of the training samples. Although the availability of open source high quality geospatial data has accelerated, it is still limited. For example, only Denmark and the Netherlands offer comprehensive and convenient access to their Land Parcel Identification System agricultural field datasets. The satellite image properties and image chip processing add to the already non-trivial training characteristics and resource requirements of deep learning models.

The presented methodology offers room for many future improvements, e.g. the integration of additional satellite image bands and multitemporal data, model architecture adjustments and additional post-processing.

## 7. Literature

Amodei, D. & Anubhai, R. & Battenberg, E. & Case, C. & Casper, J. & Catanzaro, B. & Chen, J. & Chrzanowski, M. & Coates, A. & Diamos, G. & Elsen, E. & Engel, J. & Fan, L. & Fougner, C. & Han, T. & Hannun, A. & Jun, B. & LeGresley, P. & Lin, L. & Narang, S. & Ng, A. & Ozair, S. & Prenger, R. & Raiman, J. & Satheesh, S. & Seetapun, D. & Sengupta, S. & Wang, Y. & Wang, Z. & Wang, C. & Xiao, B. & Yogatama, D. & Zhan, J. & Zhu, Z. (2015): Deep Speech 2: End-to-End Speech Recognition in English and Mandarin. <https://arxiv.org/abs/1512.02595>

Arbelaez, P. & Pont-Tuset, J. & Barron, J. & Marques, F. & Malik, J. (2014): Multiscale Combinatorial Grouping. Computer Vision and Pattern Recognition (CVPR), 2014. <https://arxiv.org/abs/1503.00848>.

Badrinarayanan, V. & Kendall, A. & Cipolla, R. (2015): SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation. <https://arxiv.org/abs/1511.00561>

Basu, S. & Ganguly, S. & Mukhopadhyay, S. & DiBiano, R. & Karki, M. & Nemani, R. (2015): DeepSat – A Learning framework for Satellite Imagery. <https://arxiv.org/abs/1509.03602>

Butenuth, M. & Straub, M. & Heipke, C. (2004): Automatic extraction of field boundaries from aerial imagery KDNet Symposium on Knowledge-Based Services for the Public Sector. p3-4. 2004.

Chen, T. & Li, M. & Li, Y. & Lin, M. & Wang, N. & Wang, M. & Xiao, T. & Xu, B. & Zhang, C. Zhang, Z. (2015): MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems. In Neural Information Processing Systems, Workshop on Machine Learning Systems, 2015. <https://arxiv.org/abs/1512.01274>

Chen, L. & Papandreou, G. & Schroff, F. & Adam, H. (2017): Rethinking Atrous Convolution for Semantic Image Segmentation. <https://arxiv.org/abs/1706.05587>

Da Costa, J. & Michelet, F. & Germain, C. & Lavialle, O. & Grenier, G. (2007): Delineation of Vine Parcels by Segmentation of High Resolution Remote Sensed Images. Precision Agriculture 8 (1 –2): 95 – 110. <https://doi.org/10.1007/s11119-007-9031-3>

Dai, J. & He, K. & Sun, J. (2015a): Convolutional feature masking for joint object and stuff segmentation. In CVPR, 2015. <https://arxiv.org/abs/1412.1283>

Dai, J. & He, K. & Sun, J. (2015b): Instance-aware Semantic Segmentation via Multi-task Network Cascades. <https://arxiv.org/abs/1512.04412>

Dai, J. & He, K. & Li, Y. & Ren, S. & Sun, J. (2016a): Instance-sensitive fully convolutional networks. In ECCV, 2016. <https://arxiv.org/abs/1603.08678>

Dai, J. & Li, Y. & He, K. & Sun. J. (2016b): R-FCN: Object detection via region-based fully convolutional networks. In NIPS, 2016. arXiv:1605.06409

Everingham, M. & Van Gool, L. & Williams, C. & Winn, J. & Zisserman, A. (2010): The PASCAL Visual Object Classes (VOC) Challenge. International Journal of Computer Vision, 88(2), 303-338, 2010.

GADM 2015: Global Administrative Areas v2.8. November 2015. <http://gadm.org/>

García-Pedrero, A. & Gonzalo-Martín, C. & Lillo-Saavedra, M. (2017): A machine learning approach for agricultural parcel delineation through agglomerative segmentation. International Journal of remote sensing, 2017. VOL. 38, NO. 7, 1809 – 1819. <http://dx.doi.org/10.1080/01431161.2016.1278312>

Girshick, R. & Donahue, J. & Darrell, T. & Malik, J. (2014): Rich feature hierarchies for accurate object detection and semantic segmentation. IEEE Conference on Computer Vision and Pattern Recognition (CVPR). <https://arxiv.org/abs/1311.2524>

Girshick, R. (2015): Fast R-CNN. <https://arxiv.org/abs/1504.08083>

Gorelick, N. & Hancher, M. & Dixon, M. & Ilyushchenko, S. & Thau, D. & Moore, R. (2017): Google Earth Engine: Planetary-scale geospatial analysis for everyone. Remote Sensing of Environment. Preprint. <https://doi.org/10.1016/j.rse.2017.06.031>

Hariharan, B. & Arbelaez, P. & Girshick, R. & Malik, J. (2014): Simultaneous detection and segmentation. In ECCV. 2014. <https://arxiv.org/abs/1407.1808>

Hariharan, B. & Arbelaez, P. & Girshick, R. & Malik, J. (2015): Hypercolumns for object segmentation and fine-grained localization. In CVPR, 2015. <https://arxiv.org/abs/1411.5752>

He, K. & Zhang, X. & Shaoqing, R. & Sun, J. (2015): Deep Residual Learning for Image Recognition. <https://arxiv.org/abs/1512.03385>

He, K. & Gkioxari, G. & Dollar, P. & Girshick, R. (2017): Mask R-CNN. <https://arxiv.org/abs/1703.06870>

Hinton, G. & Salakhutdinov, R. (2006): Reducing the dimensionality of data with neural networks. Science, Vol. 313. no. 5786, pp. 504 - 507.

Jia, Y. & Shelhamer, E. & Donahue, J. & Karayev, S. & Long, J. & Girshick, R. & Guadarrama, S. & Darrell, T. (2014): Caffe: Convolutional Architecture for Fast Feature Embedding. <https://arxiv.org/abs/1408.5093>

Längkvist, M. & Kiselev, A. & Alirezaie, M & Loutfi, A. (2016): Classification and Segmentation of Satellite Orthoimagery Using Convolutional Neural Networks. Remote Sens. 2016, 8, 329; doi:10.3390/rs8040329

LeCun, Y. & Bottou, L. & Bengio, Y. & Haffner, P. (1998): Gradient-Based Learning Applied to Document Recognition. Proc. of the IEEE, Nov. 1998.

LeCun, Y. & Bengio, Y. & Hinton, G. (2015): Deep Learning. Nature 521, 436–444. <https://doi.org/10.1038/nature14539>

Levine, S. & Pastor, P. & Krizhevsky, A. & Quillen, D. (2016): Learning Hand-Eye Coordination for Robotic Grasping with Deep Learning and Large-Scale Data Collection. <https://arxiv.org/abs/1603.02199>

Li, Y. & Qi, H. & Dai, J. & Ji, X. & Wei, Y. (2017): Fully Convolutional Instance-aware Semantic Segmentation. <https://arxiv.org/abs/1611.07709>

Lin, T. & Maire, M. & Belongie, S. & Bourdev, L. & Girshick, R. & Hays, J. & Perona, P. & Ramanan, D. & Zitnick, L. & Dollár, P. (2014): Microsoft COCO: Common Objects in Context. <https://arxiv.org/abs/1405.0312>

Lin, T. & Dollar, P. & Girshick, R. & He, K. & Hariharan, B. & Belongie, S. (2016): Feature pyramid networks for object detection. <https://arxiv.org/abs/1612.03144>

Long, J., & Shelhamer, E. & Darrell, T. (2014): Fully Convolutional Networks for Semantic Segmentation. Conference Paper.

Marmanis, D. & Schindler, K. & Wegner, J. & Galliani, S & Datcu, M. & Stilla, U. (2016): Classification with an edge: improving semantic image segmentation with boundary detection. Preprint. <https://arxiv.org/abs/1612.01337>

MFAF - Ministry for Food, Agriculture and Fisheries (2017): Marker 2016 dataset. [https://kortdata.fvm.dk/download/Index?page=Markblokke\\_Marker](https://kortdata.fvm.dk/download/Index?page=Markblokke_Marker)

Mnih, V. & Kavukcuoglu, K. & Silver, D. & Graves, A. & Antonoglou, I. & Wierstra, D. & Riedmiller, M. (2013): Playing Atari with Deep Reinforcement Learning. NIPS Deep Learning Workshop 2013. <https://arxiv.org/abs/1312.5602>

Mueller, M., & Segl, K. & Kaufmann, H. (2004): Edge-and Region-Based Segmentation Technique for the Extraction of Large, Man-Made Objects in High-Resolution Satellite Imagery. " Pattern Recognition 37 (8): 1619 – 1628. <https://doi.org/10.1016/j.patcog.2004.03.001>

Nielsen, M. (2015). Neural Networks and Deep Learning. Determination Press. <http://neuralnetworksanddeeplearning.com>

Pinheiro, P. & Collobert, R. & Dollar, P. (2015): Learning to Segment Object Candidates. <https://arxiv.org/abs/1506.06204>

Pinheiro, P. & Lin, T. & Collobert, R. & Dollar, P. (2016): Learning to Refine Object Segments. <https://arxiv.org/abs/1603.08695>

Ren, S. & He, K. & Girshick, R. & Sun, J. (2016): Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. <https://arxiv.org/abs/1506.01497>

Rossum, G. (1995): Python tutorial, Technical Report CS-R9526, Centrum voor Wiskunde en Informatica (CWI), Amsterdam, May 1995.

Rydberg, A. & Borgefors, G. (2001): Integrated method for boundary delineation of agricultural fields in multispectral satellite images. IEEE Transactions: Geoscience and Remote Sensing 39 2514–20.

Saito, S. & Yamashita, T. & Aoki, Y. (2016): Multiple Object Extraction from Aerial Imagery with Convolutional Neural Networks. Journal of Imaging Science and Technology R60(1). 2016.

Shenquan, Q. & Ying, W. & Gaofeng, M. & Chunhong, P. (2016): Vehicle Detection in Satellite Images by Incorporating Objectness and Convolutional Neural Network. Journal of Industrial and Intelligent Information Vol. 4, No. 2, March 2016.

Shrivastava, A & Gupta, A. & Girshick, R. (2016): Training Region-based Object Detectors with Online Hard Example Mining. <https://arxiv.org/pdf/1604.03540.pdf>

Simonyan, K. & Zisserman, A. (2014): Very Deep Convolutional Networks for Large-Scale Image Recognition. <https://arxiv.org/abs/1409.1556>

Szegedy, C. & Liu, W. & Jia, W. & Sermanet, P. & Reed, S. & Anguelov, D. & Erhan, D. & Vanhoucke, V. & Rabinovich, A. (2014): Going Deeper with Convolutions. <https://arxiv.org/abs/1409.4842>

Tiwari, P. & Pande, H. & Kumar, M. & Dadhwal, V. (2009) Potential of IRS P-6 LISS IV for Agriculture Field Boundary Delineation. Journal of Applied Remote Sensing 3 (1): 033528 – 033528. <https://doi.org/10.1117/1.3133306>

Torre, M. & Radeva, P. (2000): Agricultural Field Extraction from Aerial Images Using a Region Competition Algorithm, International Archives of Photogrammetry and Remote Sensing, Amsterdam, Vol. XXXIII, No. B2, pp. 889-896.

Turker, M. & Kok, H. (2013): Field-based sub-boundary extraction from remote sensing imagery using perceptual grouping. ISPRS Journal of Photogrammetry and Remote Sensing. Volume 79, May 2013, Pages 106-121.

Uijlings, J. & van de Sande, K. & Gevers, T. & Smeulders, A. (2012): Selective Search for Object Recognition.

Vakalopoulou, M. & Karantzalos, K. & Komodakis, N. & Paragios, N. (2015): Building detection in very high resolution multispectral data with Deep Learning features. IGARSS 2015 - 2015 IEEE International Geoscience and Remote Sensing Symposium.

Wu, Y. & Schuster, M. & Chen, Z. & Le, Q. & Norouzi, M. & Macherey, W. & Krikun, M. & Cao, Y. & Gao, Q. & Macherey, K. & Klingner, J. & Shah, A. & Johnson, M. & Liu, X & Kaiser, L. & Gouws, S. & Kato, Y. & Kudo, T. & Kazawa, H. & Stevens, K. & Kurian, G. & Patil, N. & Wang, W. & Young, C. & Smith, J. & Riesa, J. & Rudnick, A. & Vinyals, O. & Corrado, G. & Hughes, M. & Dean, J. (2016): Bridging the Gap between Human and Machine Translation. CoRR, vol. abs/1609.08144 (2016). <http://arxiv.org/abs/1609.08144>

Yalcin, H. & Günay, B. (2016): Delineation of parcels in agricultural lands in remote sensing. 2016 Fifth International Conference on Agro-Geoinformatics. Tianjin, China. <https://doi.org/10.1109/Agro-Geoinformatics.2016.7577601>

Yan, L. & Roy, D. (2014): Automated crop field extraction from multi-temporal Web Enabled Landsat Data. Remote Sensing of Environment 14442–64.

Yu, F. & Koltun, V. (2015): Multi-Scale Context Aggregation by Dilated Convolutions. Conference Paper ICLR 2016. <https://arxiv.org/abs/1511.07122>

Zagoruyko, S. & Lerer, A. & Lin, T. & Pinheiro, P. & Gross, S. & Chintala, S. & Dollar, P. (2016): A MultiPath Network for Object Detection. <https://arxiv.org/abs/1604.02135>

Zeiler, M. & Fergus, R. (2013): Visualizing and Understanding Convolutional Networks. <https://arxiv.org/abs/1311.2901>

Zhong, Y. & Fei, F. & Liu, Y. & Zhao, B & Jiao, H. & Zhang, L. (2017) SatCNN: satellite image dataset classification using agile convolutional neural networks. Remote Sensing Letters, 8:2, 136-145, DOI: <https://doi.org/10.1080/2150704X.2016.1235299>

## Figures

<b>Figure 1:</b> Polygons of agricultural field parcels overlaid on satellite imagery	1
<b>Figure 2:</b> Heterogenous geography of agricultural areas	2
<b>Figure 3:</b> The field of deep learning	4
<b>Figure 4:</b> Feature hierarchy of image objects	5
<b>Figure 5:</b> Four of the most important computer vision tasks	6
<b>Figure 6:</b> High-level overview of Instance segmentation training and prediction	7
<b>Figure 7:</b> Simple Neural network with input layer, one hidden layer and output layer	9
<b>Figure 8:</b> Perceptron	10
<b>Figure 9:</b> Architecture of the Alexnet network	12
<b>Figure 10:</b> Schematic functionality of a convolution neural network	12
<b>Figure 11:</b> Visual representation of the 96 convolutional filter kernels learned by the first convolutional layer of the Alexnet CNN	14
<b>Figure 12:</b> Max pooling operation	14
<b>Figure 13:</b> Class score map output and network architecture of FCN for semantic segmentation	17
<b>Figure 14:</b> Convolutionalization of fully-connected to convolutional layer	18
<b>Figure 15:</b> Superpixels and corresponding bounding boxes of Selective Search applied at different scales	20
<b>Figure 16:</b> R-CNN stages	21
<b>Figure 17:</b> Fast-RCNN architecture	22
<b>Figure 18:</b> Region Proposal Network of Faster-RCNN	23
<b>Figure 19:</b> Mask-RCNN results on the COCO test dataset	24
<b>Figure 20:</b> Different semantic categories for the same pixel when focusing different instances	24
<b>Figure 21:</b> Structure of the Multi-task Network Cascades	26
<b>Figure 22:</b> FCIS scheme	28

<b>Figure 23:</b> Mask-RCNN model with Faster-RCNN architecture and parallel instance mask branch	30
<b>Figure 24:</b> Denmark study area and detailed views showing the variety of agricultural parcel shape and geographical settings	31
<b>Figure 25:</b> Preprocessing workflow	32
<b>Figure 26:</b> Distribution of parcel classes and average area per class and sum of area per class	33
<b>Figure 27:</b> Example 128x128 image chips with ground truth agricultural field polygons after preprocessing	34
<b>Figure 28:</b> Histogram of detection scores vs. number of polygons and number of polygons remaining at a specific detection score threshold	40
<b>Figure 29:</b> Precision-recall curves at different IoU thresholds for the single-class experiment	
41	
<b>Figure 30:</b> Histograms of the number of chips with a specific per chip AP and per chip AR at different IoU thresholds	42
<b>Figure 31:</b> Example image chips from the validation dataset and predictions for the single-class setting	44
<b>Figure 32:</b> Development of the model over the training epochs	45
<b>Figure 33:</b> Systematic artifacts of redundant, overlapping predictions in FCIS results compared to Mask R-CNN	46
<b>Figure 34:</b> Example false positive predictions in image chips without ground truth instances	
48	
<b>Figure 36:</b> Example image chips from the validation dataset and predictions for the multi-class setting	50

## Tables

<b>Table 1:</b> Average precision and average recall at different IoU thresholds for the single-class experiment	41
<b>Table 2:</b> Average precision and average recall by polygon area for the single-class experimental	46
<b>Table 3:</b> Mean average precision and mean average recall at different IoU thresholds for the multi-class experiment	49
<b>Table 4:</b> Class-specific average precision at [0.5:0.95] IoU for the multi-class experiment	49

## 8. Appendix

**Table A1:** Sentinel-2 mosaic image metadata.

Product ID	Date	Sensing Orbit Direction	Cloudy Pixel Percentage	Cloud Cover Assessment
S2A_OPER_PRD_MSIL1C_PDMC_20160508T220612_R008_V20160508T104027_20160508T104027	2016-05-08	DESC	0	0.000673
S2A_OPER_PRD_MSIL1C_PDMC_20160508T220612_R008_V20160508T104027_20160508T104027	2016-05-08	DESC	0	0.000673

**Table A2:** Reclassification scheme.

Reclassified	Class IDs Denmark marker dataset
spring cereal	1,2,3,4,6,7,21, 55,56,210,211,212,213,214,215,224,230, 234,701,702,703,704,705
winter cereal	10,11,13,14,15,16,17,22,57,220,221,222,223,235
maize	5,216
grassfields	101,102,103,104,105,106,107,108,109,110,111,112,113,114,115, 116,117,118,120,121,122,123,125,126,162,170,171,172,173,174, 180,182,260,261,262,263,264,266,267,268,269,270,281,282,283, 284
other (fruits, vegetables)	23,24,25,30,31,32,35,36,40,42,51,52,53,54,55,56,57,124,160,161, 280,401,402,403,404,405,406,407,408,409,410,411,412,413,415, 416,417,418,420,421,422,423,424,429,430,431,432,434,440,448, 449,450,487,488,489,491,493,496,497,498,499,501,502,503,504, 505,507,509,512,513,514,515,516,517,518,519,520,521,522,523, 524,525,526,527,528,529,530,531,532,533,534,536,539,540,541, 542,543,544,545,547,548,549,550,551,552,553,560,561,563,570, 579
removed (environment/permanent, recreation, forest related, wasteland)	247,248,249,250,251,252,253,254,255,256,257,259,271,272,273, 274,276,277,278,279,285,286,287,305,308,309,310,311,312,313, 314,316,317,318,319,320,321,323,324,325,360,361,592,593,594, 596,597,602,603,604,605,900,903,905,907,908,920,921,997,580, 581,582,583,585,586,587,588,589,590,591

**Table A3:** Distribution of crop classes and parcel area.

	after reclass		after chips cut	
<b>chips</b>	-		10086	
<b>parcels</b>	159042		187602, 18.6 per chip (minimum area threshold 5000 sqm)	
<b>reclass_lcsub</b>	winter cereal	55973	winter cereal	73186
	spring cereal	51800	spring cereal	60930
	grassfields	32786	grassfields	32853
	maize	12242	maize	15005
	other	6241	other	4628
<b>area mean</b>	52451.78		45662.73	
<b>area in sqm reclass_lcsub</b>	winter cereal	75119.49	winter cereal	54486.44
	spring cereal	52388.08	spring cereal	41806.20
	grassfields	36355.37	grassfields	33593.77
	maize	63524.84	maize	48439.11
	other	34871.13	other	35721.10

**Table A4:** Example calculation of precision and recall (five predictions, three ground truth objects).

Detection score	Assigned	Current state	Precision	Recall
0.91	TP	1TP; 0FP; 2FN	1.0	0.33
0.73	FP (redundant TP)	1TP; 1FP; 2FN	0.5	0.33
0.61	FP (IoU<0.5)	1TP; 2FP; 2FN	0.33	0.33
0.53	TP	2TP; 2FP; 1FN	0.5	0.66
0.37	TP	3TP; 2FP; 0FN	0.6	1.0

## Acknowledgments

A big thank you goes to Adam Erickson from the Department of Biogeochemical Integration at the Max Planck Institute for Biogeochemistry Jena for the supervision of the thesis and the many helpful suggestions, ideas and motivating words.

Many thanks also to Prof. Dr. Christiane Schmüllius, head of the remote sensing department Jena, for the supervision of the thesis, her dedicated teaching and kind support throughout my studies in Jena.

I want to thank Guido Lemoine from the European Commission Joint Research Center for his presentation and advice on agricultural field datasets during the 2016 Google Earth Engine user summit, which sparked the initial idea for this thesis.

Also, I want to thank the people that contributed to my general interest in new technologies and transition into the field of programming, without which this thesis would have not been possible: Jonas Eberle from the remote sensing department Jena for the opportunity to contribute to many different projects and the interesting discussions about new geospatial trends. All members of the Google Earth Engine team, who not only created a powerful platform, but also an extremely engaging and helpful community. Finally, my brother Johannes for introducing me to the field of deep learning, the great programming advice throughout the last couple of years and generally for sharing many of my interests.

## **Statement of Authorship**

I hereby confirm that this thesis has been composed by me and is based on my own work, unless stated otherwise. No other person's work has been used without due acknowledgement. All references and verbatim extracts have been quoted, and all sources of information have been specifically acknowledged.

---

Munich, 13.09.2017