

Object Detection in Fine Art Photography

Bachelor thesis, spring semester, 2020

Fabian Meyer

Supervisor: Dr. Simone Lionetti

Lucerne University of Applied Sciences and Arts
Bachelor in Information Technology
May 30, 2020

Bachelorarbeit an der Hochschule Luzern - Informatik

Titel: Object Detection in Fine Art Photography

Student: Fabian Meyer

Studiengang: Bachelor Informatik

Abschlussjahr: 2020

Betreuungsperson: Dr. Simone Lionetti

Expertin / Experte: Roman Bachmann, Swisscom

Eidesstattliche Erklärung

Ich erkläre hiermit, dass ich die vorliegende Arbeit selbständig und ohne unerlaubte fremde Hilfe angefertigt habe, alle verwendeten Quellen, Literatur und andere Hilfsmittel angegeben haben, wörtlich oder inhaltlich entnommene Stellen als solche kenntlich gemacht haben, das Vertraulichkeitsinteresse des Auftraggebers wahre und die Urheberrechtsbestimmungen der Fachhochschule Zentralschweiz (siehe Markblatt „Studentische Arbeiten“ auf MyCampus) respektieren werde.

Ort / Datum, Unterschrift: _____

Abgabe der Arbeit auf der Portfolio Datenbank

Bestätigungsvisum Studentin / Student

Ich bestätige, dass ich die Bachelorarbeit korrekt gemäss Merkblatt auf der Portfolio Datenbank abgelegt habe. Die Verantwortlichkeit sowie die Berechtigungen habe ich abgegeben, so dass ich keine Änderungen mehr vornehmen kann oder weitere Dateien hochladen kann.

Ort / Datum, Unterschrift: _____

Verdankung

Danke liebe Freunde

Eingangsvisum (durch das Sekretariat auszufüllen):

Rotkreuz, den _____ Visum: _____

Hinweis: Die Bachelorarbeit wurde von keinem Dozierenden nachbearbeitet. Veröffentlichungen (auch auszugsweise) sind ohne das Einverständnis der Studiengangleitung der Hochschule Luzern Informatik nicht erlaubt.

Copyright © 2019 Hochschule Luzern - Informatik

Alle Rechte vorbehalten. Kein Teil dieser Arbeit darf ohne die schriftliche Genehmigung der Studiengangleitung der Hochschule Luzern - Informatik in irgendeiner Form reproduziert oder in eine von Maschinen verwendete Sprache übertragen werden.

Abstract

Contents

1	Introduction	1
1.1	Task description	1
1.2	Why fine art photography?	2
2	Current state	4
2.1	About convolutional neural networks	4
2.2	From CNNs to Mask R-CNN	5
3	Methods	8
3.1	Toolchain	8
3.2	Project management	8
3.3	Testing	9
4	Ideas and concepts	10
4.1	Basic idea	10
4.2	A more precise research question	10
4.3	Alternative research questions	10
4.4	Alternative models	10
5	Realisation	11
5.1	Data collection	11
5.2	Model selection	12
5.3	Framework selection	12
5.4	Choosing the programming language	12
5.5	Creating the tidied up image	12
5.6	Building the application	14
5.7	Deployment	15
5.8	Data labelling	15
5.9	Refining the model	16
5.10	Results	16
6	Evaluation und Validation	17
6.1	Vergleich mit Anforderungen	17
6.2	Technische Aspekte	17
6.3	Vorgehen	17

1 Introduction

Object detection is an important computer vision and machine learning task that consists in the localisation and classification of items within digital images. The business and industry applications of this technology are growing in number and relevance, especially because of the tremendous progress, largely driven by deep learning, that has been made in recent years.

Uses of object detection technologies include applications in security measures for example images from surveillance cameras, search engines, object tracking or counting in traffic, autonomous driving cars, robotics in general and the field of medical- and bioinformatics.

1.1 Task description

The aim of this work is to develop a web-application that is built around an object-detection model. This web-application has to take in a picture, looks for objects in it and has to return a "tidied up" image with the objects found in it. The "tidied up" image is inspired by "Kunst aufräumen", a series of fine-art photographies by Swiss artist Urs Wehrli, that consists of two images, a "messy" one and a "tidy" one. One such an example picture is shown here:

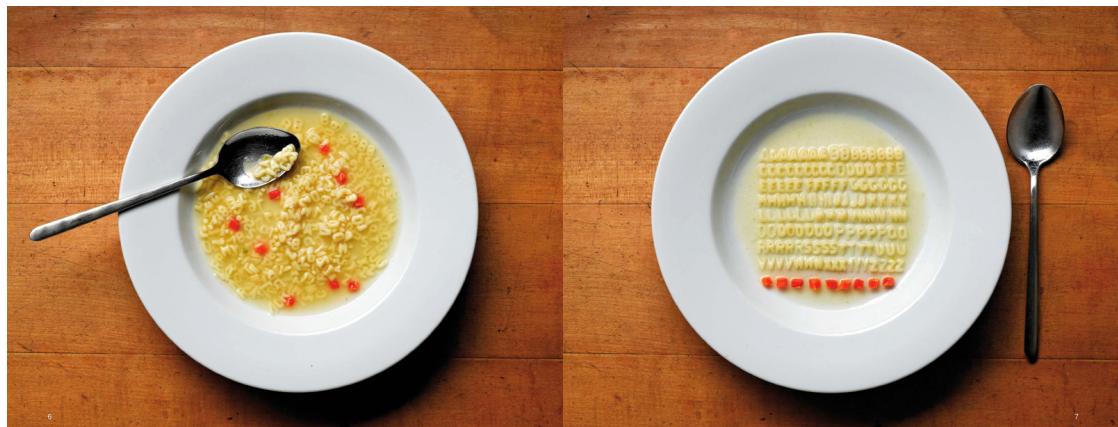


Figure 1.1: Sample image from "Kunst aufräumen" by Urs Wehrli

In a further step, an existing object detection model has to be refined by retraining it on a dataset with images from artistic photography and the performance of the two models have to be examined.

1.2 Why fine art photography?

Models for object detection are usually trained on datasets, containing pictures that show objects in a very clear manner. These images usually contain only a few objects and are shot in natural lighting with natural colours. Composition-wise they are simple as it is the goal to depict the object in a most clear way.



Figure 1.2: Sample image from COCO dataset

Fine art photography pictures on the other hand is in strong contrast to these images, as they are often depicting objects who do not belong together usually. These images are often shot in studios with artificial lighting and heavy post-processing or even digital manipulation. And they often contain a lot more objects in it which even can overlap each other.



Figure 1.3: Sample image from artist David LaChapelle

The question that now arises is, how do object detection models that are trained on traditional datasets perform on given images of artistic photography?

2 Current state

2.1 About convolutional neural networks

The model that is used for this project is from the family of convolutional neural networks or CNNs. These are models that are especially useful to perform operations on digital images. CNNs contain at least one convolutional layer that computes a mathematical operation that is called discrete convolution. Discrete convolution operation has been proved useful when applied on digital images to detect structural features like edges and corners. When applying discrete convolution to an image (a two dimensional array of values), a filter with a given two dimensional kernel is used to perform convolution on the image. With the given kernel, a segment that has the same size as the kernel is taken from the input image and an elementwise multiplication between these two matrices is performed. The sum of the results gets written in the output (also called feature map). As most kernels are much smaller than the input image, this inevitable results in smaller feature maps than input images, especially after applying convolution to the same image multiple times. To solve this problem, one can add a padding (increase the image size by adding specific values). An example of a two dimensional convolution operation with padding is shown here:

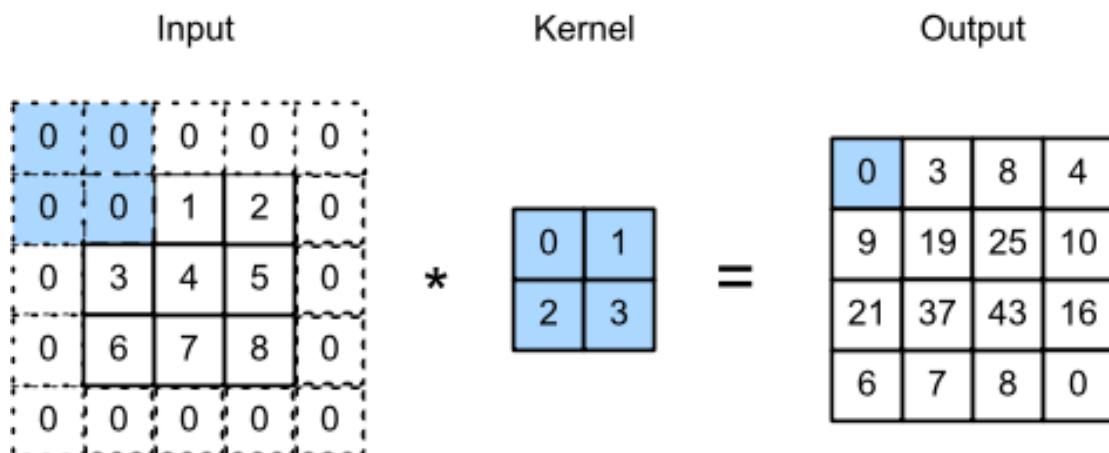


Figure 2.1: Example of a two dimensional convolution operation with padding

After the convolution is computed for a specific segment, the kernel shifts further some steps in the image array, called stride. There exists a lot of different kernels for different purposes. To detect edges and corners there exist specific kernels especially for that kind of task. It is also possible to learn the optimal kernel in the process of neural network training.

The advantage of CNNs compared to traditional multi layer perceptrons (MLPs) is that discrete convolutional operation is translational invariant. This means it is irrelevant where in a given image a specific object is located: The model will not learn the position of the specific object. Another huge advantage of CNNs over MLPs is that CNNs usually have sparse interactions. This means that it is possible to yield a good performance even if some adjacent layers are not fully connected with each other. Another way to speed up the training process is to use parameter sharing: Some weights will be used at multiple places at the same time. These two techniques, sparse interactions and parameter sharing does make it possible to train very deep convolutional neural networks in a feasible amount of time.

Another important kind of layer that is usually employed in CNNs are pooling layers. Also called subsampling or downsampling layers sometimes, they reduce the size of a given image, by computing a statistic metric to summarize multiple values. Some frequently used pooling layers are max- and average-pooling. The advantages of pooling layers are smaller input, parameter reduction and higher invariance to scaling and transformations.

The third important kind of layer is ReLU (rectified linear unit). The mere use of this layer is to preprocess the image before the next convolutional operation. This just adjusts the image brightness in a ways that all negative values (values that are darker than the middle grey of an image) got adjusted to middle grey.

The last layer of a CNN contains the number of classes that the network should predict. This layer is usually fully connected to the layer before.

2.2 From CNNs to Mask R-CNN

The very first CNN appeared in the year 1994, it was named LeNet5, by Yann LeCun. This fundamental work was the first network that used convolutional and pooling layers to process images. In a time long before consumer graphic processing units (GPUs), where even CPUs were slow, it was important to reduce the number of parameters to a bare minimum (accomplished with sparse connect layers). An image, showing LeNet5 with its different layers and operations in between is shown here:

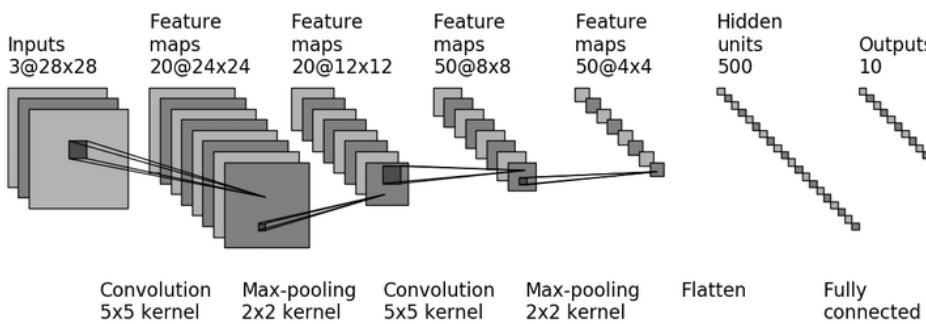


Figure 2.2: Overview of LeNet5 from the year 1994

Starting with AlexNet by Alex Krizhevsky, Ilya Sutskever and Geoffrey Hinton in 2012, CNNs gained significantly performance gains when applied on image classification tasks.

AlexNet took the idea of a CNN from LeNet5 and added more layers to it and was the first to include ReLU layers. It also used dropout techniques to avoid overfitting during the training process. AlexNet was still an image detection or image classification model, trained to label images. This means a single label (class) is given to a whole input image.

ResNet (residual neural network) in 2015 was the first CNN that included residual blocks or identity shortcut connections. These are connections in the network that skip one or more layers and just use the input value as an output (mathematical identity). This was another step to reduce computation cost and lead to even deeper networks. ResNet is still used today as part of the backbone in a lot of object detection and instance segmentation models.

With R-CNN (regions with CNN features) in 2013, the rise of object detection models began. These researcher asked themselves, how can one use the techniques from CNNs to not only classify an image but to classify objects in that image? These family of models are capable of labelling multiple objects depicted in the same image with each a bounding box and a class prediction. R-CNN models do that by proposing regions in that potential objects may lie. R-CNN is thus called a two stage model: The first stage scans the image and generates region proposals, whereas the second stage uses these regions to extract CNN features from it and to ultimately classify objects in it. For the classification step in the last layer, R-CNN uses a support vector machine (SVM) as a classifier. In the very last step, a regression is used to further tighten the coordinates of the bounding boxes of each object. [1]

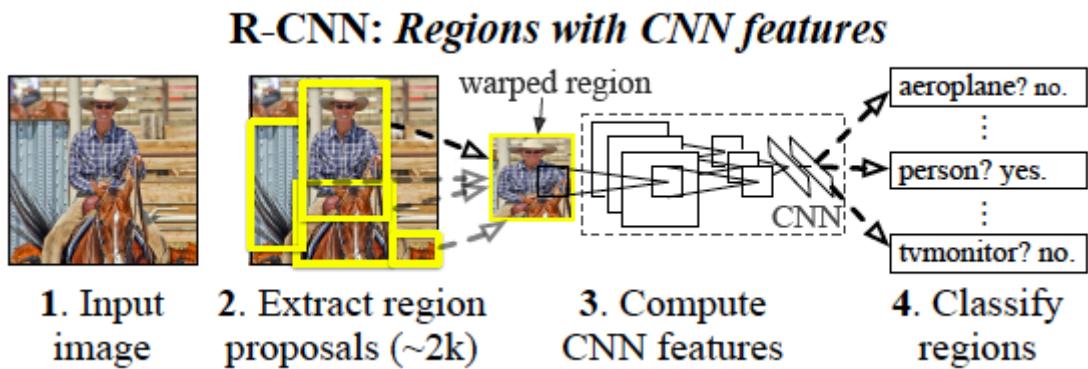


Figure 2.3: Architecture of R-CNN

In 2017, the team from Ross Girshick, one of the creators of R-CNN, delivered Fast R-CNN an improved, faster version of it. In R-CNN there are a lot of overlapping region proposals and every computation gets calculated again, even if the regions are very similar to each other. To circumvent this they invented region of interest pooling (RoIPool). RoIPool shares these computation across the regions of an image and can speed up computation time a lot. The other improvement was to put all computations in a single network (compared to R-CNN where for example classification ran in a single network).

Also in 2017, also by the team from Ross Girshik, the second iteration of R-CNN got re-

leased, called Faster R-CNN. The main improvement was to use just one CNN that produces a feature map for both region proposal and classification.

With Mask R-CNN from 2017 (also by Ross Girshik et. al) it is possible to not only predict the bounding box of an object but also to predict a mask that shows the exact shape of the object (called pixel level segmentation). These models are called instance segmentation models. This is accomplished by adding a branch to the model that computes a binary mask for a given object in the image. Mask R-CNN model is of the family of instance segmentation models. An overview of Mask R-CNN can be seen here:

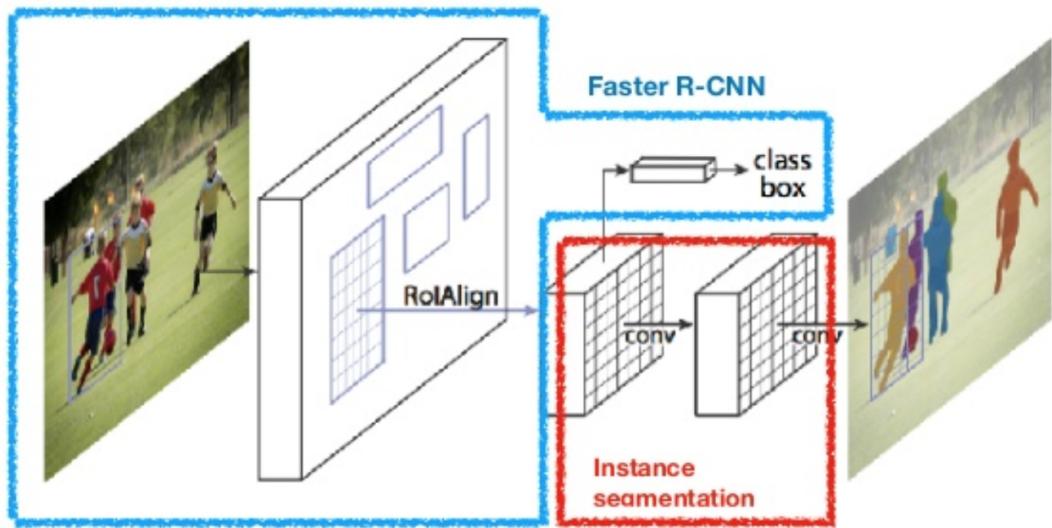


Figure 2.4: An overview of Mask R-CNN

An overview over the different tasks and approaches and their models in the field of object recognition can be seen here:

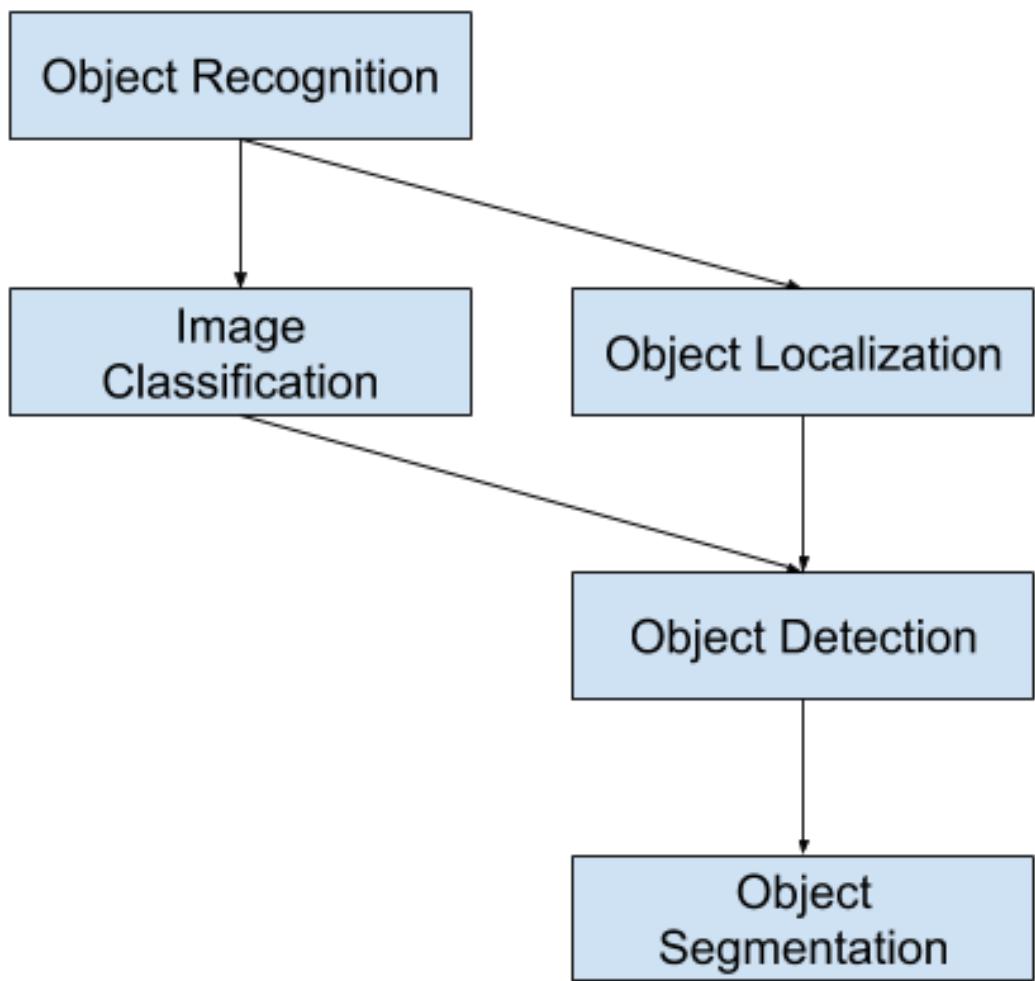


Figure 2.5: Tasks and approaches in object recognition

3 Methods

3.1 Toolchain

The whole project can be summarized in four single steps or tasks:

1. Data: Collect a dataset, for inference, for training and for testing
2. Model: Get to know at least one model, for inference and for retraining
3. Tidy: For a given input picture, return an image that displays a repertoire of all found objects in the input picture
4. App: Turn the former three tasks into an application that can be reached from a web browser.

In order to develop a minimal viable product, all these four tasks have to be solved first. The way of proceeding was to develop a minimal viable product first and then to start with refining the model and adjusting the app. This means going through the full circle first, before adjustments and proceeding into retraining is made.

3.2 Project management

In the beginning of the project, a project management plan had to be developed. The project management plan contains a timeline with all (then known) tasks and a list of all milestones. An image of the project management plan can be seen via this link: <https://trello.com/b/srqnMstX/object-detection-in-fine-art-photography>. A picture, showing the whole timeline is shown here:

To better structure the project, a number of milestones have been chosen:

1. Test pretrained models on Google Colab with new data
2. Build web application prototype
3. Finetune model with fine art photography data
4. Test finetuned model on new data
5. Finish web application with new model

Object Detection In Fine Art Photography

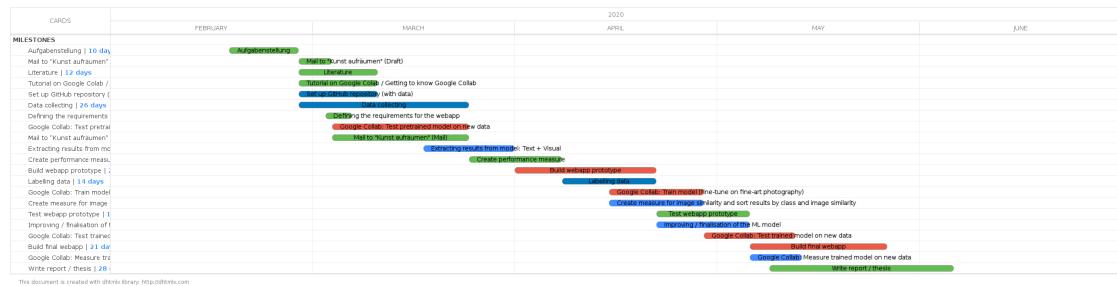


Figure 3.1: Timeline view of the project management plan

3.3 Testing

In favour of having more time available for the development cycle and retraining of the model, there was no testing strategy selected.

4 Ideas and concepts

4.1 Basic idea

4.2 A more precise research question

4.3 Alternative research questions

4.4 Alternative models

5 Realisation

5.1 Data collection

The first step in this project was to collect some data. Two pictures of fine art photography were already shown in the project description: An image by german photographer Andreas Gursky and one by Swiss artist Urs Wehrli. The main motivation behind data collection was to obtain images from different artists that each show a lot of different objects. This let us examine the performance of an object detection model applied to fine art photography, by using images that contain a lot of objects that optimally are also included in the COCO dataset.

The COCO (Common Objects in COntext) dataset is among the most popular image datasets created for object recognition tasks. It was lanced by Microsoft in 2014 and the most recent version was published in 2017. It uses 80 different classes for object detection tasks that can be seen here:



Figure 5.1: All 80 classes from COCO dataset 2017

The COCO dataset not only contains bounding boxes and binary masks for object detection and segmentation, it also includes informations for keypoint estimation, panoptic segmentation and image captioning tasks. It is one of the big baselines for object detection tasks for many years now.

In the end, 42 different images from the four artists Urs Wehrli, Andreas Gursky, David LaChapelle and Jeff Wall has been gathered. Urs Wehrli, the swiss artist behind "Kunst aufräumen" was asked to give his permission to use his images for this project. Thankfully he gave permission to do so. The dataset has then been examined with different object segmentation models on Google Colab.

Google Colab is a free to use infrastructure, powered by Google, that offers a zero-configuration Python and Jupyter environment. Running on Linux Ubuntu with the latest Nvidia GPUs this is a good option to get started with deep learning, as there are also a lot of notebooks about every single kind of neural network tasks available.

5.2 Model selection

With the gathered dataset, different object segmentation models from different frameworks have been tried out on the dataset in inference mode. All these models have been pretrained on the COCO-dataset. The tried out models were: Mask R-CNN, CenterNet, Detectron2 and ShapeMask. In general, Mask R-CNN outperformed the other models, when inference is run on the dataset. Detectron2 also delivered a good performance but it threw an error when running on many images in a loop on Google Colab.

5.3 Framework selection

At the same time different frameworks have been tested out. These were Tensorflow from Google, Detectron from Facebook and MMDetection, which is a part of the OpenMMLab project developed by Multimedia Laboratory by the Chinese University of Hong Kong. All these frameworks do at least contain one Mask R-CNN model. MMDetection has been chosen as the framework to develop the project in, because it worked out-of-the box when tried out on Google Colab. It returned results when running in inference mode on our dataset in about seven Minutes. MMDetection has a vast number of state-of-the-art models for computer vision tasks available and offers a high-level API that speeds up its usage. It is built on top of PyTorch (primarily developed by Facebook) and delivers a very good performance.

5.4 Choosing the programming language

Because most of the deep-learning frameworks are using Python and also of its ease of use, Python has been used as the sole programming language to develop this project. In addition Python does offer a lot of visual computing libraries that were used to create the tidied up image.

5.5 Creating the tidied up image

After the dataset has been collected, and the model and the framework have been chosen, the next step was to extract all single objects from the output of the model when running in inference mode on an image. Mask R-CNN model outputs for every input image among other things a list with each a list of bounding boxes, masks, confidence score and predicted class of all objects found in the image. These four results were saved in a list object. The output is per default limited to 100 objects per image. The confidence score was used to filter all objects to get the objects with the highest quality as predicted by the model. The binary mask of the found objects was then applied to the image, creating binary images of the object. An example can be seen here:

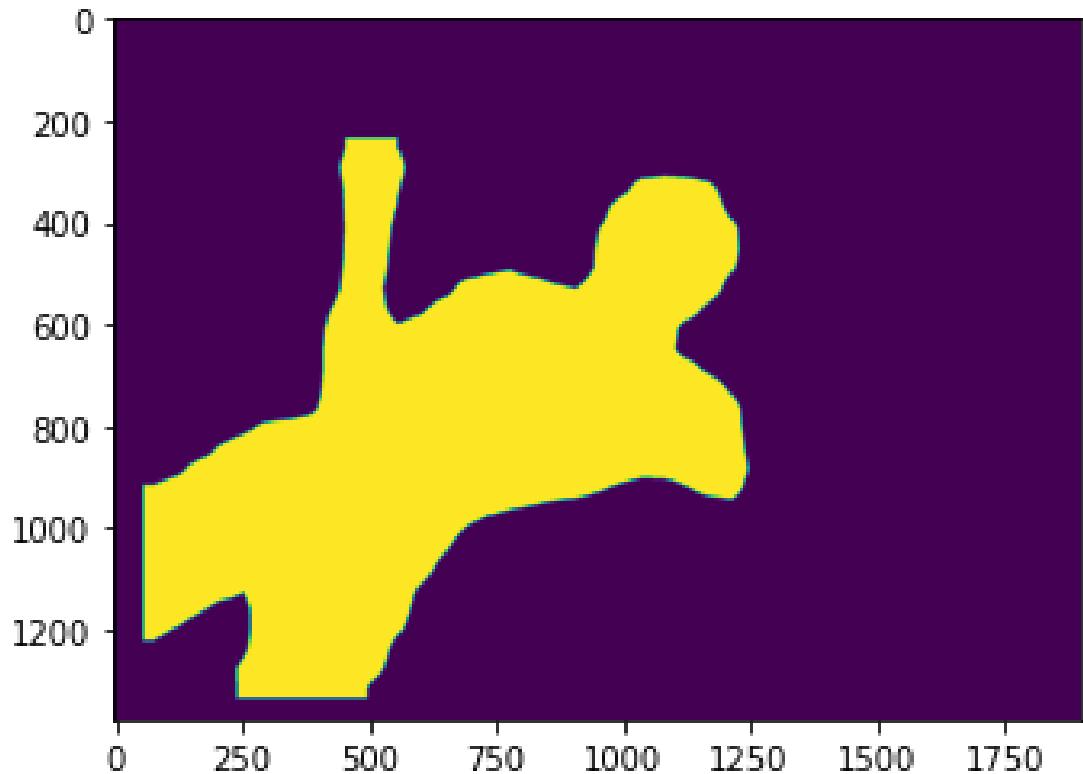


Figure 5.2: Sample binary mask image

With help of the binary masks, the background of the found object was coloured white whereas the foreground stayed as is. Subsequently the corresponding bounding boxes were taken to cut out the objects from the image. The last step is to create a grid with all the found objects. This was done with the help of Python's Matplotlib. The idea is to create a grid with the number of rows according to the number of classes and insert every object as a new column. The problem with using Matplotlib's `plt.subplots` is that it creates a grid where every image has the same size. To circumvent this, a more complicated approach was taken, by creating a grid manually with the help of the width and the height of all the found objects. A sample input of a tidied up image together with its output can be seen here:

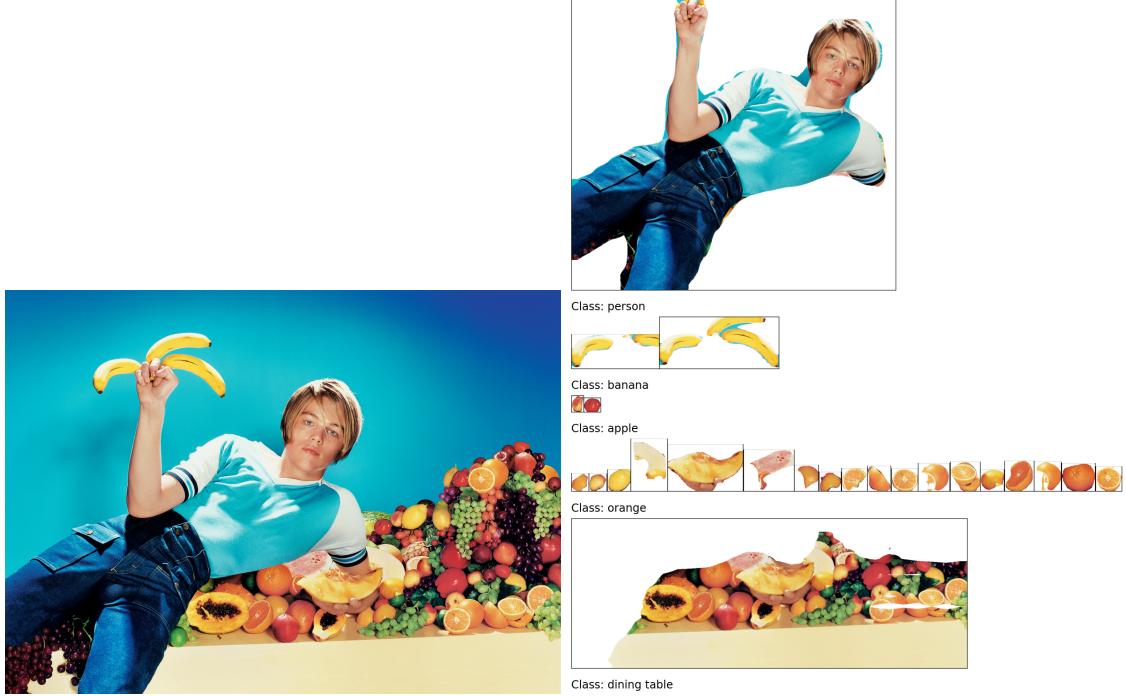


Figure 5.3: Sample input and output image

5.6 Building the application

After completing the tidied up image, the next step was to build an application that can load an image, run the model in inference mode and return the tidied up image from the input image. A difficulty was that MMDetection, like most deep learning frameworks, needs a Nvidia GPU to run even in inference mode. To solve this problem, a wrapper called SMD was used, that takes in MMDetection models and can run the inference on an arbitrary CPU. [2]

To convert the program into a full-blown web application, a framework, called Plotly-Dash has been used. Plotly-Dash is using Flask to create a webserver and is using HTML-, CSS- and JavaScript-technologies under the hood to create a running web application from a Python program. It offers out-of-the-box user interface (UI) components, that are useful to rapid prototype a web application. With the help of UI-elements like buttons and dropdown-lists, the user is given the possibility to upload and select images and models and to start the inference by himself. There are also two sliders built in: One to adjust the confidence threshold score of the model and one to adjust the input image size. Executed time was measured and gets written too to get an insight into performance of the model.

One of the goals of the web application was to build it as modular as possible. This was achieved with the possible selection of the model and three ways to select an image from (predefined list, upload and retrieve via URL). One difficulty with Plotly-Dash was that images when uploaded to the web server, got encoded with base64 encoding. It took some

time to find out, how to decode them for further use. A screenshot showing the web application and its UI can be seen here:

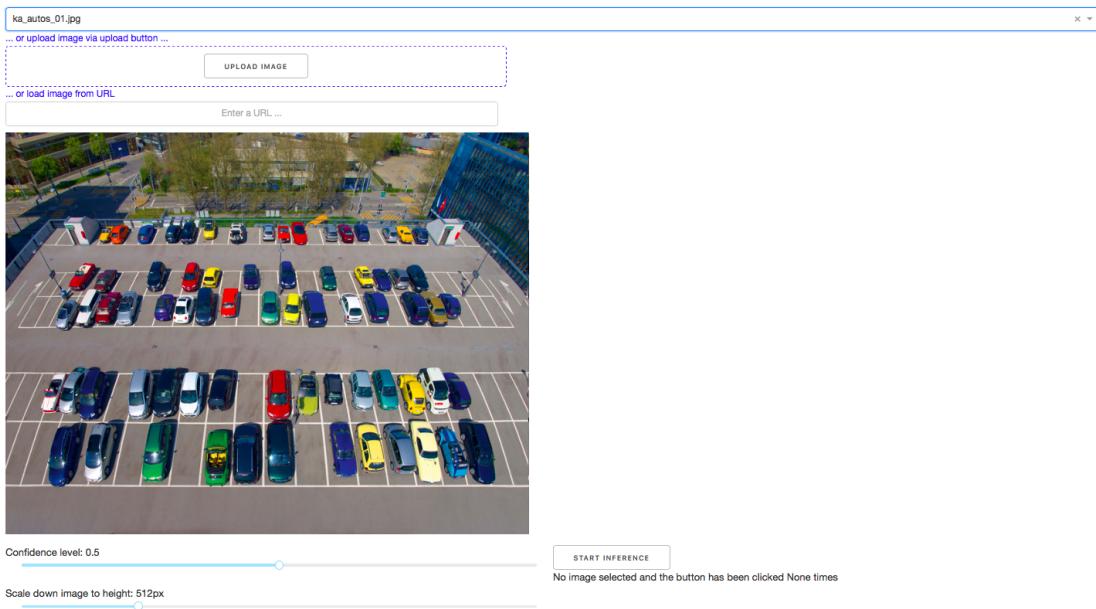


Figure 5.4: Screenshot of the webapplication

5.7 Deployment

To deploy the application, server space on the EnterpriseLab has been allocated. The web application can be accessed at: <http://bdaf20-iameyer.enterpriselab.ch>.

5.8 Data labelling

A dataset containing of three images from the series "Kunst aufräumen" by Urs Wehrli has been chosen to refine the model. As each of these three images contains a messy and a tidy version, it does make a good template to use the tidied up version as a training image and the messy version as a test image. To accomplish training with own images, one has to label these images first. When training an object segmentation model, it is necessary to train the classifier with images that contain masks in it. There are several tools that offer object masking with polygons or brushes. A total of six images (three messy ones and three tidy ones) have been labelled in order to use for refining the model. During labelling, a .xml-file gets generated, that after completion can be converted into the COCO-format as a .json-file.

5.9 Refining the model

Google Colab was used again as the platform to retrain the model. With just three training images it is quite extreme to train a classifier. This reflected in the poor performance in the beginning. After deciding to retrain a single model for every picture pair (messy and tidy one), performance got better. For a task like this, when training with a small dataset, data-augmentation is crucial: It lets the model use the same image over and over again after applying different kind of transformations and distortions to it.

When retraining with MMDetection, one has to adjust several things: A .json-file, containing all the classes, masks and bounding boxes in COCO-format. And a config.py-file that contains all the needed settings for the training process. During retraining, checkpoint-files get saved in a defined interval. These checkpoint-files can be used resume training on a later point of time. After one training cycle (between 1000 and 2000 epochs), the latest checkpoint-file was taken and examined on the test-image.

5.10 Results

Hello

6 Evaluation und Validation

6.1 Vergleich mit Anforderungen

6.2 Technische Aspekte

6.3 Vorgehen

List of Figures

1.1	Sample image from "Kunst aufräumen" by Urs Wehrli	1
1.2	Sample image from COCO dataset	2
1.3	Sample image from artist David LaChapelle	3
2.1	Example of a two dimensional convolution operation with padding	4
2.2	Overview of LeNet5 from the year 1994	5
2.3	Architecture of R-CNN	6
2.4	Tasks and approaches in object recognition	7
3.1	Timeline view of the project management plan	9
5.1	All 80 classes from COCO dataset 2017	11
5.2	Sample binary mask image	13
5.3	Sample input and output image	14
5.4	Screenshot of the webapplication	15

List of Tables

Formelverzeichnis

Bibliography

- [1] Girshick, R.B., Donahue, J., Darrell, T., & Malik, J. (2014). Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. 2014 IEEE Conference on Computer Vision and Pattern Recognition, 580-587.
- [2] Karaniewicz, A. (n.d.). SMD (simple mmdetection). Retrieved May 30, 2020, from <https://github.com/akarazniewicz/smd>