

Machbarkeitsstudie: Smart Warehouse

Echtzeit-Objektdetektoren im Vergleich

Studienarbeit

im Rahmen der Prüfung zum
Bachelor of Science (B.Sc.)

des Studienganges Angewandte Informatik
an der Dualen Hochschule Baden-Württemberg Karlsruhe

von

Felix Hausberger und Robin Kuck

Oktober 2019 - Mai 2020

-Sperrvermerk-

Abgabedatum: 18. Mai 2020
Bearbeitungszeitraum: 30.09.2019 - 18.05.2020
Matrikelnummer, Kurs: 2773463, 4409176, TINF17B2
Ausbildungsfirma: SAP SE
Dietmar-Hopp-Allee 16
69190 Walldorf, Deutschland
Gutachter an der DHBW: PD Dr.-Ing. Markus Reischl

Eidesstattliche Erklärung

Wir versichern hiermit, dass wir unsere Studienarbeit mit dem Thema:

Machbarkeitsstudie: Smart Warehouse

gemäß § 5 der „Studien- und Prüfungsordnung DHBW Technik“ vom 29. September 2017 selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt haben. Die Arbeit wurde bisher keiner anderen Prüfungsbehörde vorgelegt und auch nicht veröffentlicht.

Wir versichern zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

Karlsruhe, den 12. April 2020

Gez. Felix Hausberger und Robin Kuck
Hausberger, Felix und Kuck, Robin

Abstract

- *English* -

In this thesis the object detectors *You Only Look Once* and *Single Shot MultiBox Detector* are compared for precision, reactivity, training and inference behaviour and examined for their potential for industrial use. The background scenario of the *Smart Warehouse* offers live video data of a drone with goods in a warehouse, which are to be classified and localized in real time. In the future, this should make it possible to carry out inventories and inventory analyses of a warehouse in a time- and cost-efficient manner conserving resources.

The goal of this feasibility study is to find out whether the *Smart Warehouse* scenario is technically feasible. In addition, the focus is also on the object detectors themselves, their differences in architecture and behavior and whether they are generally suitable for industrial application scenarios.

Abstract

- *Deutsch* -

In dieser Arbeit werden die Objektdetektoren *You Only Look Once* und *Single Shot MultiBox Detector* nach Präzision, Reaktionsvermögen, Trainings- und Inferenzverhalten miteinander verglichen und auf deren Potential zum industriellen Einsatz untersucht. Das Hintergrundszenario des *Smart Warehouses* bietet dabei Live-Video Daten einer Drohne mit Warengegenständen in einem Warenhaus, die in Echtzeit klassifiziert und lokalisiert werden sollen. Dadurch sollen in Zukunft in der Industrie Inventuren und Bestandsanalysen eines Warenhauses zeit- und kostengünstig sowie ressourcenschonend ermöglicht werden können.

Diese Machbarkeitsstudie hat zum Ziel herauszufinden, ob das Szenario des *Smart Warehouse* technisch umsetzbar ist. Zusätzlich liegt der Fokus ebenso auf den Objektdetektoren selbst, deren Unterschiede hinsichtlich Architektur und Verhalten und ob sie allgemein für industrielle Anwendungsszenarien grundsätzlich geeignet sind.

Inhaltsverzeichnis

Abkürzungsverzeichnis	VI
Abbildungsverzeichnis	VIII
Formelverzeichnis	IX
Listenverzeichnis	X
0 Vorwort	1
1 Einführung	2
1.1 Forschungsumfeld	2
1.2 Problemstellung und Motivation	3
1.3 Vorgehensweise und Zielsetzung	3
2 Grundlagen und Forschungsstand	5
2.1 Deep Learning zur Bildverarbeitung	5
2.2 Neuronale Netze	6
2.3 Hyperparameter	12
2.4 Datensatzlehre	17
2.5 Grundlagen zu Objektdetektoren	19
2.6 Objektdetektoren	24
2.7 Datensatzformate	33
2.8 Cloud Infrastruktur	35
3 Konzeption	40
3.1 Erstellen eines Trainingsdatensatzes	41
3.2 Einführen von Bewertungskriterien	42
3.3 Auswahl der Objektdetektoren	44
3.4 Auswahl der Trainingsinfrastruktur	46
3.5 Auswahl einer Drohne	49
3.6 Spezifikation der Inventursoftware	50
4 Realisierung	52

4.1	Umsetzung der Objektdetektoren	52
4.2	Drohnen Anbindung	55
4.3	Dashboard Entwicklung	55
4.4	Zählalgorithmus	56
5	Ergebnisse	57
6	Bewertung	58
7	Zusammenfassung und Ausblick	59
	Literaturverzeichnis	XI
A	Anhang	XVII

Abkürzungsverzeichnis

AI	Artificial Intelligence
ANN	Artificial Neural Network
ASIC	Application-Specific Integrated Circuit
AWS	Amazon Web Services
CNN	Convolutional Neural Network
COCO	Common Objects in Context
CLI	Command Line Interface
CPU	Central Processing Unit
CUDA	Compute Unified Device Architecture
DBN	Deep Belief Network
DevOps	Development Operations
EASA	European Aviation Safety Agency
ELU	Exponential Linear Unit
FCN	Fully Convolutional Network
FLOPS	Floating Point Operations Per Second
FPGA	Field Programmable Gate Array
FPS	Frames Per Second
GCP	Google Cloud Platform
GPU	Graphics Processing Unit
IoU	Intersection over Union
LReLU	Leaky Rectified Linear Unit
MLP	Multi-Layer Perceptron
mAP	mean Average Precision

PReLU	Parametric Rectified Linear Unit
PascalVOC	Pascal Visual Object Classes
PaaS	Platform-as-a-Service
R-CNN	Regional Convolutional Neural Network
ReLU	Rectified Linear Unit
REST	Representational State Transfer
RoI	Region of Interest
RPN	Region Proposal Network
SaaS	Software-as-a-Service
SDK	Software Development Kit
SSD	Single Shot MultiBox Detector
TOPS	Tera Operations Per Second
TPU	Tensor Processing Unit
XML	Extensible Markup Language
YOLO	You Only Look Once

Abbildungsverzeichnis

2.1	Anwendungsgebiete von Deep Learning zur Bildverarbeitung im Überblick	6
2.2	Linear Threshold Unit	7
2.3	Das einschichtige Perzeptron	9
2.4	Gradientenverfahren	15
2.5	ReLU-Aktivierungsfunktionen	16
2.6	ELU-Aktivierungsfunktion	17
2.7	Convolutional Layer	20
2.8	Zero-Padding	20
2.9	Veranschaulichung von Feature Maps	22
2.10	Pooling Layer	23
2.11	Intersection over Union	24
2.12	Berechnung mAP	25
2.13	R-CNN Architektur	26
2.14	Faster R-CNN Architektur	27
2.15	SSD Bounding Box Proposals	28
2.16	Bounding Boxes	29
2.17	SSD Architektur	30
2.18	Vereinfachte Darstellung des YOLO Algorithmus	32
2.19	YOLO Architektur	33
3.1	Konzeptionelle Schritte	40
3.2	Die neun Datensatz-Kategorien	42
3.3	Smart Warehouse Regal	43
3.4	Vergleich SSD auf PascalVOC	44
3.5	Vergleich V100 - TPU Pod	49
3.6	SmartWarehouse User Interface	51
4.1	Detektion mehrerer Wasserflaschen	54
4.2	Detektion einer nahen Wasserflasche	54
4.3	Detektion einer entfernten Wasserflasche	54

Formelverzeichnis

2.1	Die Heaviside-Funktion	8
2.2	Die Softmax-Funktion	8
2.3	Die RMSE-Funktion	10
2.4	Neuberechnung der Gewichtungsmatrix durch partielle Differentiation . .	11
2.5	Superpositionsprinzip anhand der Varianz	12
2.6	Standardverteilung nach Xavier Initialisierung	13
2.7	Momentum Optimierung	14
2.8	Precision und Recall	23
2.9	Die Smooth L1 Funktion	29

Listenverzeichnis

2.1	PascalVOC Bildannotation	33
2.2	Konfigurationsdatei zum Trainingsjob auf FloydHub	37

0. Vorwort

Besonderen Dank ist an unseren Betreuer PD Dr. -Ing. Markus Reischl auszusprechen, ohne den die folgenden Forschungsergebnisse nicht zustande gekommen wären. Auch dem Informatik Labor unter Enrico Hühneborg der DHBW ist für die nötige finanzielle Unterstützung zum Erwerb der Drohne zu danken.

1. Einführung

1.1. Forschungsumfeld

Einen Teilbereich des maschinellen Lernens (engl.: machine learning) stellt das *Deep Learning* dar, welches auf künstlichen neuronalen Netzen (engl.: artificial neural networks) (ANNs) basiert [1]. Unter einer Vielzahl von Typen von ANNs wie Autoencodern, Deep Boltzmann Machines oder rekurrenten neuronale Netzen befindet sich ebenso die Klasse der *Convolutional Neural Networks* (CNNs), welche hauptsächlich zur Lösung von Klassifikationsproblemen in der Audio-, Text- und Bildverarbeitung genutzt werden [2].

Ein Forschungsfeld im *Deep Learning* stellen Objektdetektoren dar, welche basierend auf CNNs neben Bildklassifikationsproblemen ebenso in der Lage sind, Lokalisationsprobleme zu lösen. Solchen Objektdetektoren werden in der heutigen Zeit immer mehr Bedeutung zugesprochen angesichts neuer Herausforderungen wie autonomen Fahren, automatisierter industrieller Verarbeitung oder aber auch staatlicher Überwachung. Verschiedene Ansätze werden zur Realisierung von Objektdetektoren verwendet, unter anderem Netzarchitekturen wie *Regional Convolutional Neural Networks* (R-CNNs), *You Only Look Once* (YOLO) oder der *Single Shot MultiBox Detector* (SSD).

Gerade in Zeiten des industriellen Wandels in Richtung *Industrie 4.0* können solche Objektdetektoren ein großes Optimierungspotential für bestehende Industrieszenarien bieten, beispielsweise in der Lagerhaltung und Logistik. Kombiniert mit einer autonomen Drohne können Objektdetektoren es ermöglichen, ohne menschliche Hilfe Inventuren und Bestandsprüfungen in einem Lager- oder Warenhaus durchzuführen. Start-up Unternehmen wie *doks. innovation* werben bereits mit ähnlichen Lösungen, die 80% Zeiteinsparung und 90% Kostensenkung versprechen [3]. Lösungen wie *inventAIRyX* beschränken sich allerdings speziell auf Lagerhäuser, in denen die verpackten Waren mittels Sensoren identifiziert werden, was Großhändler mit Warenhäusern wie *Baumarkt* oder *Selgros* ausschließt. Statt Waren mittels RFID Chips oder Barcodes zu identifizieren, soll in dieser Arbeit der Einsatz von Objektdetektoren für dieses Szenario evaluiert werden.

Wie sich die unterschiedlichen Objektdetektoren unter Echtzeitvoraussetzungen im Be-

trieb verhalten, soll anhand des Industriebeispiels *Smart Warehouse* innerhalb dieser Arbeit untersucht werden.

1.2. Problemstellung und Motivation

Das *Smart Warehouse* beschreibt ein Warenhaus, welches unter Einsatz einer Drohne in der Lage sei soll, Inventuren und Bestandsprüfungen weitgehend ohne menschliche Hilfe durchzuführen. Das Live-Bild der Drohne soll von den Objektdetektoren dazu genutzt werden, Warengegenstände zu lokalisieren und zu klassifizieren.

Neben der Frage, ob ein solches Industrieszenario überhaupt umsetzbar ist, sollen die Objektdetektoren in diesem Anwendungsszenario nach verschiedenen Kriterien miteinander verglichen und beurteilt werden. Diese Kriterien lassen sich hauptsächlich in die Kategorien Präzision, Reaktionsvermögen, Trainings- und Inferenzverhalten untergliedern und werden genauer eingeführt. Dadurch lassen sich Aussagen darüber treffen, ob nach dem momentanen Forschungsstand um Objektdetektoren solche das Potential bieten, industriell eingesetzt zu werden.

Falls die Machbarkeitsstudie des *Smart Warehouse* glückt, so kann der Industrie ein kostengünstiges, zeitsparendes und ressourcenschonendes Modell zur Inventurverwaltung eines Warenhauses angeboten werden.

1.3. Vorgehensweise und Zielsetzung

Im Grundlagenkapitel 2 muss sich mit den theoretischen Grundlagen von CNNs und Objektdetektoren auseinander gesetzt werden. Hierzu ist zunächst eine Einführung in Deep Learning zur Bildverarbeitung und neuronale Netz erforderlich, darunter zu Perzepronen, dem Gradientenverfahren, dem Backpropagation Algorithmus und Hyperparametern zum Trainieren eines neuronalen Netzes (siehe Kapitel 2.1, 2.2 und 2.3).

Um weitere Grundlagen zum Training von neuronalen Netzen einzuführen, wird anschließend über die Anforderungen eines Datensatzes gesprochen (siehe Kapitel 2.4), bevor weitere Grundlagen zu Objektdetektoren eingeführt werden (siehe Kapitel 2.5).

Nachdem zu Beginn des Kapitels 2.5 kurz auf den Grundbaustein moderner Objektdetektoren eingegangen wird, den CNNs, können anschließend die Funktionsweisen und Architekturen der drei miteinander verglichenen Objektdetektoren der *R-CNN* Familie, *YOLO* und des *SSDs* erläutert werden. Bei *R-CNN* und *YOLO* ist zu bemerken, dass unterschiedliche Evolutionsstufen der Detektoren zu betrachten sind.

Um weitere Grundlagen zum Training von neuronalen Netzen einzuführen, wird anschließend über zwei wesentlichen Speicherformate eines Datensatzes gesprochen (siehe Kapitel 2.7), bevor verschiedene Cloud Anbieter für das Trainieren von *Deep Learning* Modellen aufgezeigt werden (siehe Kapitel 2.8).

In Kapitel 3 werden chronologisch Teilziele der Konzeption beschrieben, darunter das Erstellen eines Trainingsdatensatzes (siehe Kapitel 3.1), dem Einführen von Bewertungskriterien (siehe Kapitel 3.2), der Auswahl von Objektdetektoren, der Trainingsinfrastruktur und der Drohne (siehe Kapitel 3.3, 3.4, 3.4) und die Spezifikation der Inventursoftware (siehe Kapitel 3.6).

In der Realisierung werden die Herausforderungen zur Steuerung und Anbindung der Drohne betrachtet und zudem die Objektdetektoren auf die realen Datensätze trainiert. Auch die Entwicklung der Webapplikation zur Visualisierung des Live-Bildes und der erkannten Objekte wird Bestandteil dieses Kapitels sein. Die Ergebnisse der Realisierungsphase werden im folgenden Kapitel dargestellt.

Ziel der Arbeit ist es Aussagen über die Fähigkeit von Objektdetektoren zum Einsatz in der Industrie zu treffen, indem eine Bewertung der Verhaltensweisen der Objektdetektoren nach den eingeführten Bewertungskriterien durchgeführt wird. Dies soll mit Hilfe der Umsetzung des *SmartWarehouse* Szenarios bewiesen werden.

Zuletzt wird das Wesen der Arbeit nochmals kurz zusammengefasst und anschließend auf mögliche Verbesserungen und Ausblicke in die Zukunft aufmerksam gemacht.

2. Grundlagen und Forschungsstand

Neben einer Einführung in den Anwendungsbereich von *Deep Learning* zur Bildverarbeitung soll sich das folgende Kapitel speziell mit Architekturen unterschiedlicher Objektdetektoren auseinander setzen und herausstellen, wie sich diese voneinander abgrenzen. Davor wird allerdings zunächst grundlegendes Wissen über neuronale Netze und wie diese „lernen“ vermittelt sowie wie ein eigener Datensatz zu gestalten ist.

2.1. Deep Learning zur Bildverarbeitung

Ein klassisches Anwendungsgebiet von *Deep Learning* zu Bildverarbeitung oder auch allgemein von maschinellem Lernen beschreibt die *Klassifikation*. Hierbei werden bestimmte Kategorien, auch *Klassen* genannt, definiert, in die ein Bild eingeordnet werden soll. Die *Klassifikation* wird anhand von aus dem Bild extrahierten Merkmalen, auch *Features* genannt, getroffen. Die Merkmale werden zu einem *Merkmalsvektor* oder auch *Feature Map* zusammengefasst und von einem *neuronalen Netz* verarbeitet. Das Ergebnis der Verarbeitung durch das neuronale Netz ist die Einordnung in eine bestimmte Klasse.

Zusätzlich zur *Klassifikation* eines Bildes kann das auf dem Bild abgebildete Objekt ebenfalls lokalisiert werden. Es wird dann von sogenannter *Objektdetektion* gesprochen. Es können auch mehrere Objekte auf einem Bild detektiert werden. Ergebnis der Objektdetektion ist somit nicht nur eine Klasseneinordnung sondern ebenso eine klare Positionsangabe des Objektes auf dem Bild. Die Positionsangabe erfolgt durch Angabe einer sogenannten *Bounding Box*. Diese umrahmt das jeweils detektierte Objekt und wird durch ihren linken oberen Eckpunkt sowie ihre Höhe und Breite beschrieben.

Neben der klassischen *Klassifikation* und der *Objektdetektion* existiert ebenso ein drittes Anwendungsgebiet von *Deep Learning* zu Bildverarbeitung, die *Segmentierung*. Bei der *semantischen Segmentierung* wird versucht, jede einzelne Pixel eines Bildes einer Klasse zuzuordnen und dementsprechend farblich im Bild zu hinterlegen. *Instanzbasierte Segmentierung* hingegen zielt darauf ab, nicht nur jeden Pixel zu einer Klasse zuzuordnen, sondern ebenso eine Identität zu einem Objekt zuzuweisen [4]. Es setzt sich zusammen

aus *semantischer Segmentierung* und paralleler *Objektdetektion*.

Einen Überblick über die vorgestellten Anwendungsgebiete ist in Abbildung 2.1 zu sehen.

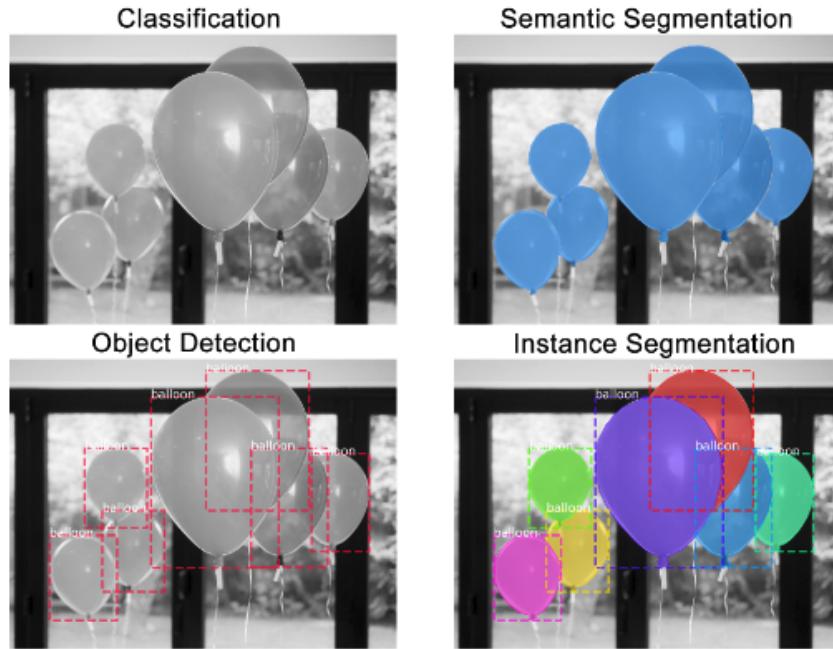


Abbildung 2.1.: Anwendungsgebiete von Deep Learning zur Bildverarbeitung im Überblick [5]

Für eine einfache *Klassifikation* eines Bildes können einfache sogenannte *Feed-Forward* Netze verwendet werden. Es kann hierbei aber auch auf konventionelle Methoden der Bildverarbeitung zurück gegriffen werden. Für die *Objektdetektion* stehen Architekturen wie *You Only Look Once* (YOLO), der *Single Shot MultiBox Detector* (SSD) oder neuronale Netze der *Regional Convolutional Neural Networks* (R-CNN) bereit. Aus dieser Familie entstammt ebenso das *Mask R-CNN* Netz, das zur *instanzbasierten Segmentierung* von Objekten verwendet wird.

2.2. Neuronale Netze

Ein neuronales Netz bildet die Grundlage des *Deep Learnings* [1]. Zunächst soll die einfachste Architektur eines neuronalen Netzes, das Perzeptron [1], exemplarisch erklärt

werden als auch der Lernprozess eines maschinellen Lernmodells an sich, um darauf basierend die Auswirkungen von Hyperparametern auf den Lernprozess des Modells zu erklären.

Das Perzepron

Der Aufbau eines typischen Perzeprons besteht aus einer oder mehreren Schichten so genannter *Linear Threshold Units* (LTU) wie in Abbildung 2.2 dargestellt.

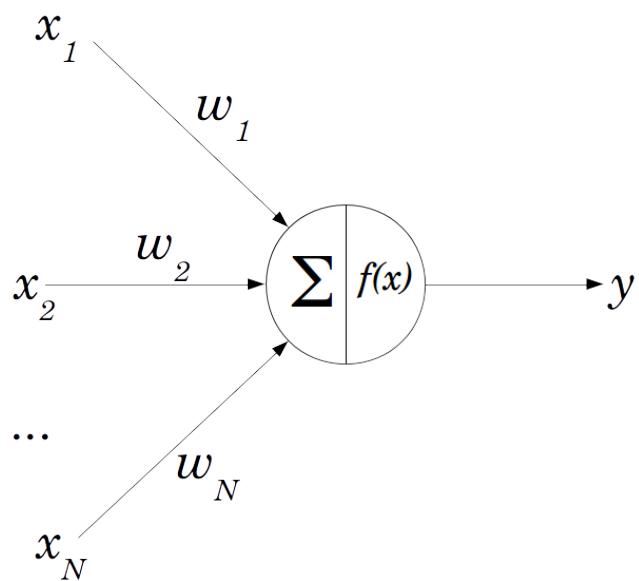


Abbildung 2.2.: Linear Threshold Unit [6]

Es besteht aus n Eingängen mit $x_i \in \mathbb{Q}$, die im Inputvektor \mathbf{x} zusammengefasst werden. Jeder Eingang wird mit einem Gewicht w_i aus dem Gewichtsvektor \mathbf{w} versehen [1]. Die LTU berechnet das Skalarprodukt $\mathbf{w}^T \circ \mathbf{x}$ aller Eingänge \mathbf{x} mit ihren Gewichten \mathbf{w} und wendet anschließend auf das Ergebnis z eine Aktivierungsfunktion an [1]. Das Ergebnis $h_w(x)$ kann anschließend als Eingabe für ein weiteres Perzepron dienen. Die einfachste

Aktivierungsfunktion für ANNs ist die *Heaviside-Funktion* [1]:

$$h_w(x) = s(\mathbf{w}^T \circ \mathbf{x}) = s(z) = \\ \left(w_1 \quad w_2 \quad \dots \quad w_n \right) \circ \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{cases} 1 & \text{wenn } z \geq 0 \\ 0 & \text{wenn } z < 0. \end{cases} \quad (2.1)$$

Falls eine Klassifizierung mit Wahrscheinlichkeiten vorliegen soll, so ist die letzte Schicht eines Perzeptrons meist mit der *Softmax-Funktion*

$$h_w(x) = \sigma(z)_j = \frac{e^{z_j}}{\sum_{i=0}^n e^{z_i}} \quad (2.2)$$

implementiert, die den Wert des j -ten LTUs einer Schicht mit allen anderen n Werten der LTUs derselben Schicht ins Verhältnis setzt [1]. Es gibt eine Vielzahl an möglichen Aktivierungsfunktionen, die im darauffolgenden Unterkapitel *Hyperparameter* betrachtet werden.

Die Aktivierung einer LTU hängt zusätzlich von einem Schwellwert θ ab, der durch einen sogenannten *Bias* festgelegt wird. Dies ist die Gewichtung des letzten Eingangs, der standardmäßig den Wert 1 liefert. Wird die Gewichtung negativ gewählt, so ist es schwieriger die LTU zu aktivieren, während eine positive Gewichtung die Aktivierung vereinfacht [1].

Nun bilden ein oder mehrere Schichten solcher LTUs ein Perzepron. Jede einzelne LTU ist dabei mit allen LTUs der vorherigen Schicht verbunden (siehe Abbildung 2.3) [1]. Hier wird auch von sogenannten vollständig verbundenen Schichten (engl.: *Fully-Connected Layer*) gesprochen. Die beiden LTUs zur Ausgabe können dabei Aussagen über eine Klassifikation von Daten anhand der Eingangsdaten treffen, während die LTUs im Input Layer wesentlich Daten weiter reichen. Die Verbindungen zur ersten Schicht des Hidden Layer sind stets mit Eins belegt. Existiert keine verborgene Schicht, so wird das ANN als einschichtiges Perzepron bezeichnet, ab einer oder mehr verborgenen Schichten wird

bereits von einem *Multi-Layer Perceptron* (MLP), einem mehrschichtigen Perceptron, gesprochen [1]. Ist das neuronale Netz optimal trainiert, so ist am Ende nur eines der LTUs zur Ausgabe aktiviert. Das folgende ANN ist zudem ein Beispiel für ein sogenanntes *Feed Forward Network*, bei dem die Auswertung der Daten von einer Schicht zur nächsten weitergereicht wird, ohne zu bereits besuchten Schichten zurückzukehren [1].

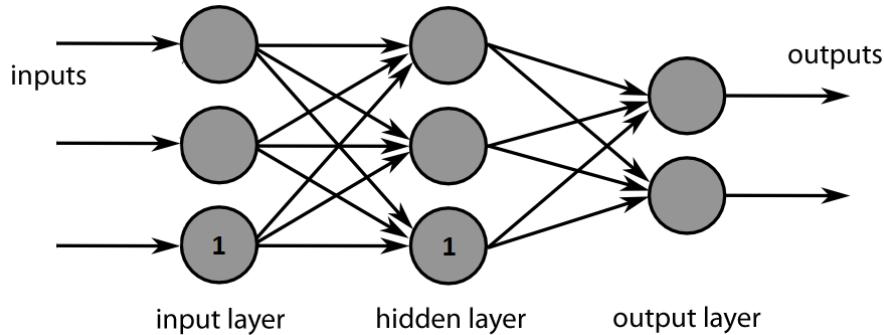


Abbildung 2.3.: Das einschichtige Perceptron [7]

Lernmethoden

Im Wesentlichen existieren vier Methoden, mit deren Hilfe neuronale Netze trainiert werden können. Beim *überwachten Lernen* werden den Trainingsdaten Lösungen, sogenannte *Labels*, hinzugefügt. Klassifikationsprobleme stellen eine typische Problemstellung für überwachte Lernverfahren dar. Auch die später eingeführten Objektdetektoren werden mittels überwachtem Lernen trainiert [1].

Beim *unüberwachten Lernen* werden den Trainingsdaten keinerlei Lösungen hinzugefügt, der Algorithmus muss selbstständig Klassifikationsaussagen treffen können. Ein Beispiel hierzu wäre *K-Means Clustering* [1].

Deep Belief Networks (DBNs) bestehend aus einzelnen *Restricted Boltzmann Machines* (RBMs) werden zunächst unüberwacht trainiert, bevor im *Fine Tuning* das Gesamtnetzwerk mit überwachten Lerntechniken fertiggestellt wird. Es wird hier von *halbüberwachtem Lernen* gesprochen.

Als letztes ist das sogenannte *Reinforcement Learning* zu nennen. Hierbei wird ein neuronales Netz durch das Erteilen von Belohnungen bzw. Bestrafungen so konditioniert, dass

es zukünftig basierend auf der wahrgenommenen Umwelt selbstsicher richtige Aktionen auswählen kann.

Gradientenverfahren und Backpropagation

Um zu verstehen, wie ein neuronales Netz durch überwachtes Lernen „lernt“, muss zunächst der Begriff der Kostenfunktion (engl.: *cost function*) eingeführt werden. Die Kostenfunktion ist ein Qualitätsmaß dafür, wie weit die Ausgabe einer LTU vom erwarteten Wert abweicht [1]. Angenommen dem neuronalen Netz wird ein Datensatz zur Klassifikation übergeben, so ist am Ende meist nicht nur eine LTU zu Ausgabe aktiviert, was auf eine eindeutige Klassifikation schließen würde, sondern meist mehrere zu einem frühen Stadium des neuronalen Netzes.

Eine oft genutzte Kostenfunktion ist die *Root Mean Squared Error Funktion* (RMSE) [1]:

$$E(\mathbf{z}, \mathbf{o}) = \sqrt{\frac{1}{n} \sum_{k=0}^n \|z_k - o_k\|^2} = \sqrt{\frac{1}{n} \sum_{k=0}^n (z_k - o_k)^2}. \quad (2.3)$$

Hierbei ist z der erwartete Ausgabevektor des Perzeptrons beim überwachten Lernen, während o die momentane Ausgabe darstellt. Den Fehler der Abweichung dieser beiden Werte gilt es nun schrittweise zu minimieren. Um dies zu erreichen können die drei Parameter

1. Gewichtung der Verbindungen zum Perzeptron
2. Bias zur Aktivierung der LTUs des Perzeptrons und
3. Stärke der Aktivierung des vorherigen Perzeptrons

angepasst werden [1]. Hierbei wird das sogenannte *Gradientenverfahren* eingesetzt. Es berechnet in einem iterativen Prozess über mehrere Testdaten das globale Minimum der Kostenfunktion nach den Gewichtungen der Verbindungen und damit auch nach den Bias Werten, die natürlich ebenso Gewichtungen darstellen. Ergebnis eines Durchlaufs im Gradientenverfahren (2.4) ist die Gewichtungsmatrix, die die Änderung der Gewichtung jeder einzelnen Verbindung eines Perzeptrons zu jeder LTU des Folgeperzeptrons angibt

[1]:

$$w_{ijt} = w_{ijt-1} - \eta \frac{\partial E}{\partial w_{ij}}. \quad [6] \quad (2.4)$$

Das Gradientenverfahren eignet sich allerdings nur für stetig differenzierbare Funktionen ohne Plateaus. Somit können beispielsweise bei der Heaviside-Funktion als Aktivierungsfunktion Probleme auftreten, da eine Ableitung der Kostenfunktion stets Null betragen würde, wohingegen bei der später eingeführten Sigmoid-Funktion im gesamten Definitionsbereich immer kleine Änderungen der Gewichtungen zu verzeichnen wären [1].

Nun stellt sich auch der Vorteil von MSE als Kostenfunktion gegenüber anderen, durchaus komplexeren Kostenfunktionen heraus. Während MSE genau ein Minimum, das zugleich das globale Minimum der Funktion darstellt, besitzt, haben andere Kostenfunktionen im Gradientenverfahren das Problem, dass anstelle des globalen Minimums auch nur lokale Minima erreicht werden können [1]. Dies hat zur Folge, dass mehrere iterative Durchläufe mit mehreren Testdatensätzen nötig werden, um durch unterschiedliche Startkonfigurationen die unterschiedlichen Minima miteinander vergleichen zu können und damit das globale Minimum herauszustellen.

Durch das Gradientenverfahren werden somit nur diejenigen Verbindungen verstärkt, die zum richtigen Ergebnis führen.

Nun bleibt nur noch die dritte Möglichkeit zur Minimierung der Kostenfunktion übrig, die Anpassung der Stärke der Aktivierung des vorherigen Perzeptrons. Zu diesem Problem veröffentlichten David E. Rumelhart, Geoffrey E. Hinton und Ronald J. Williams 1985 den sogenannten *Backpropagation-Algorithmus* [8]. Dieser berechnet mit Hilfe des Gradientenverfahrens welchen Anteil am Fehler der Ausgabe jede LTU des letzten Perzeptrons hat und anschließend welcher Anteil davon wiederum auf das vorherige Perzepron der vergorenen Schicht zurück zu führen ist. Das Gradientenverfahren wird solange wiederholt, bis die Eingangsschicht erreicht wurde, es berechnet also für jede LTU deren Anteil am Fehler des Ergebnisses [1].

Mit Hilfe des Gradientenverfahrens im Backpropagation Algorithmus wird nun also das neuronale Netz durch mehrere iterative Durchläufe trainiert, wobei das Training als Anpassung der Gewichtungen einzelner Verbindungen zu verstehen ist.

2.3. Hyperparameter

Hyperparameter sind die Parameter, die zur anfänglichen Konfiguration des neuronalen Netzes als auch zur Konfiguration des Lernprozesses herangezogen werden. Um im Laufe der Arbeit verstehen zu können, wie die Objektdetektoren auf Seiten der Netzarchitektur und des Lernverhaltens optimiert wurden, ist demnach ein kurzer Einblick in den Themenbereich der Hyperparameter von Nöten.

Anzahl der LTUs

Die Anzahl der LTUs im ANN ist dafür ausschlaggebend, wie hoch der Komplexitätsanspruch eines Klassifizierungsproblems sein darf, um noch vom ANN gelöst werden zu können. Die Anzahl der LTUs hängt hauptsächlich von den Eingangsdaten ab. Über die optimalste Anzahl an LTUs pro Schicht lässt sich allerdings nur schwer etwas vorhersagen. Generell gilt, dass bei gleicher Anzahl an LTUs tiefere Netze eine weitaus höheren Parametereffizienz aufweisen als breitere Netze, da diese schneller gegen den gewünschten Zustand konvergieren. Zudem lassen sie sich somit schneller und kostengünstiger trainieren. So müssten bei einem 2x32 Netz 1024 Gewichtungen angepasst werden, während es bei einem 32x2 Netz dies nur 128 sind [1].

Initialisierung der Gewichtungen

Auch stellt die Initialisierung der Gewichte eines ANNs zu Beginn des Trainingsprozesses eine berechtigte Frage dar. Falls keine bereits trainierten ANNs für ein Klassifikationsproblem vorliegen, so werden die Gewichtungen meist zufällig nach einer Normalverteilung gewählt [1].

Dies hat allerdings zur Folge, dass nach der Berechnung der gewichteten Summen aller LTUs die Gewichtungswerte der folgenden Schicht nicht mehr normalverteilt sind, da für die Varianz zweier unkorrelierter Zufallsvariablen das Superpositionsprinzip

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) \quad (2.5)$$

gilt.

Durch die größer werdende Standardabweichung können demnach Gewichtungswerte entstehen, die weit vom Mittelwert Null abweichen. Dies kann wiederum dazu führen, dass der Gradientenabstieg während des Backpropagation-Verfahrens nur langsam vollzogen werden kann, da der Gradient bei bestimmten Aktivierungsfunktionen wie der *Sigmoid-Funktion* gegen Null konvergiert [1].

Eine *Xavier Initialisierung* umgeht das Problem der sogenannten *schwindenden Gradiennten*, indem die Gewichte nach

$$W \sim U\left[-\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}\right] \quad (2.6)$$

gleichverteilt werden, wobei n_j die Anzahl an LTUs der j -ten Schicht sind [9].

Anzahl an Epochen

Die Anzahl der Epochen beschreibt die Durchläufe durch einen bestimmten Trainingsdatensatz während der Trainingsphase. Ist die Anzahl zu hoch gewählt, wird Gefahr gelaufen, sogenanntes *Overfitting* des ANNs zu erreichen. Dies bedeutet ein fehlendes Abstraktionsvermögen des ANNs durch Auswendiglernen der Trainingsdaten zu erreichen und damit alleinig eine richtige Erkennung der Trainingsdatensätze zu ermöglichen.

Lernrate

Die Lernrate η (2.4) gibt an, wie groß die Sprünge sein sollen und damit indirekt wie viele Iterationen benötigt werden, um das Minimum der Kostenfunktion zu erreichen. Ziel der Anpassung einer Lernrate ist es, mit möglichst wenig Iterationen und Testdaten die optimale Konstellation des neuronalen Netzes zu berechnen. Deshalb wird sie standardmäßig zu Beginn der Iterationen groß gewählt, um sich dem Minimum schnell zu nähern, während sie am Ende immer kleiner gewählt wird, um nicht über das globale Minimum hinaus zu gehen. Dieses Vorgehen wird als *Simulated Annealing* bezeichnet, während das Funktion zum Festlegen der Lernrate als *Learning Schedule* betitelt wird

[1].

Die Anzahl der Durchläufe wird zu Beginn des Verfahrens zunächst hoch angesetzt, das Verfahren wird aber genau dann gestoppt, sobald der Gradientenvektor unter eine gewisse Abbruchgrenze fällt. Zwar ist das globale Minimum zu diesem Zeitpunkt noch nicht erreicht, allerdings kann es auch nie vollkommen erreicht werden, da die für das Gradientenverfahren genutzten Aktivierungsfunktionen nie einen partiellen Ableitungswert gleich Null zulassen [1]. In diesem Sinne wird auch von *Toleranz* gesprochen.

Moment

Das Gradientenverfahren kann beschleunigt werden, indem während des Gradientenabstiegs frühere Gradienten Einfluss auf den nächsten Gradientenschritt nehmen. Es wird ein „Momentum“ aufgebaut. Damit das Momentum

$$m_x = \beta \cdot m_{x-1} + \eta \frac{\partial E}{\partial w_{ij}} \quad (2.7)$$

$$w_{ijt} = w_{ijt-1} - m$$

allerdings nicht zu groß wird, beschränkt der Hyperparameter $\beta \in [0, 1]$ die Größe des Momentums [1].

Die Momentum Optimierung kann dazu benutzt werden, das *stochastische Gradientenverfahren* bzw. *Mini-Batch* Verfahren zu beschleunigen und lokale Minima besser zu überwinden.

Auswahl des Gradientenverfahrens

Generell wird zwischen drei verschiedenen Arten unterschieden, das Gradientenverfahren durchzuführen (siehe Abbildung 2.4):

Beim *Batch* Verfahren werden in einem Trainingsdurchlauf, der *Epoche*, alle vorhandenen Daten des Trainingsdatensatzes herangezogen, um einen Gradientenabstieg zu vollziehen. Dies ist bei großen Trainingsdatensätzen auffällig langsam, dafür aber hinsichtlich der

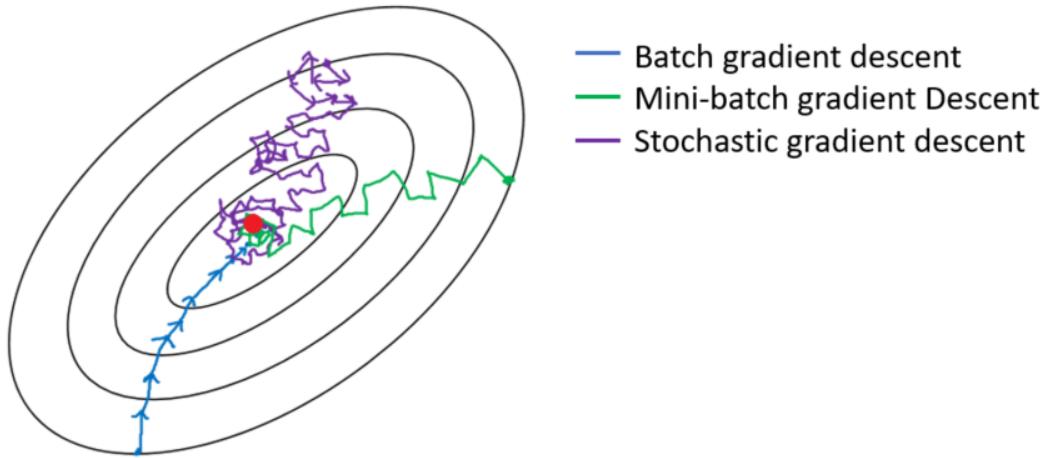


Abbildung 2.4.: Gradientenverfahren [10]

Erreichung des lokalen Minimums sehr zielstrebig [1].

Das *stochastische Gradientenverfahren* führt nach jedem einzelnen Dateneintrag im Trainingsdatensatz einen Gradientenabstieg durch. Da nur wenige Daten des ANNs verändert werden müssen, ist dieses Verfahren deutlich schneller, dafür aber unregelmäßiger hinsichtlich der Erreichung des Minimums. Oft wird das stochastische Gradientenverfahren verwendet, wenn nicht der komplette Trainingsdatensatz in den Hauptspeicher oder Grafikspeicher geladen werden kann. Diese Fähigkeit wird oft als *Out-of-Core* Fähigkeit bezeichnet. Es hat auch den Vorteil, besser das globale Minimum der Kostenfunktion aufzufinden, da bei lokalen Minima die Chance besteht, durch den unregelmäßigen Gradientenabstieg das lokale Minimum wieder zu überwinden [1].

Ein Kompromiss der beiden Verfahren bietet das *Mini-Batch* Verfahren, bei dem wiederholt Teilmengen des gesamten Datensatzes für einen Gradientenabstieg verwendet werden. Genauso wie das *Batch* Verfahren bietet das *Mini-Batch* Verfahren den Vorteil, die partiellen Ableitungen als Matrizenoperationen auf die Grafikkarten auszulagern, um die Performanz durch Parallelisierung zu steigern [1].

Aktivierungsfunktionen

Zwei bekannte und ähnliche Aktivierungsfunktionen sind die *Sigmoid-Funktion* und die *Tangens Hyperbolicus* Funktion. Da diese allerdings durch ihr schnelles Konvergieren gegen den Grenzwert anfällig für das Problem *schwindender Gradienten* sind [1], wird die *Rectified Linear Unit* (ReLU) bzw. *Parametric/Leaky Rectified Linear Unit* (PReLU/LReLU) Aktivierungsfunktion bevorzugt (siehe Abbildung 2.5).

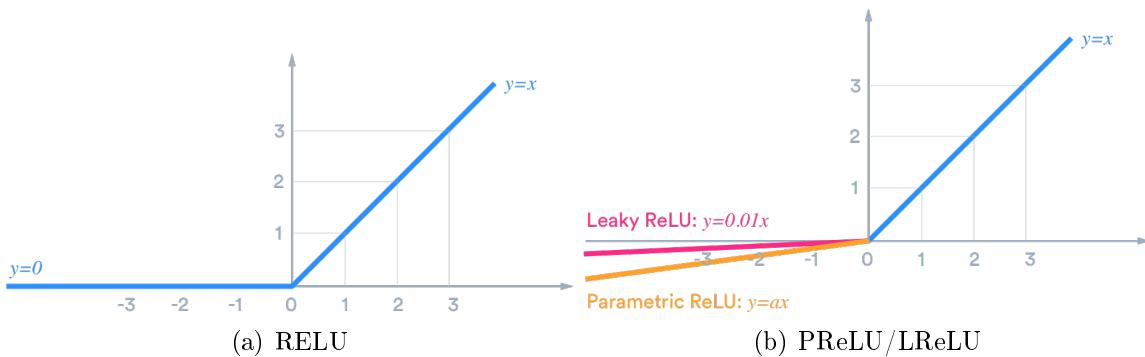


Abbildung 2.5.: ReLU-Aktivierungsfunktionen [11]

Bei ReLU kann es während des Trainingsprozesses dazu kommen, dass LTUs nach dem Gradientenabstieg einen negativen Wert aufweisen, weshalb sie nicht weiter aktiviert werden und für den Rest der Trainingsdauer „tot“ sind. Um dies zu verhindern, wurde *LReLU* dazu genutzt, um eine Reaktivierung zu ermöglichen, da auch für negative LTU Werte ein Gradient der Aktivierungsfunktion bestimmt werden kann. Bei *LReLU* ist die Steigung der Funktion im zweiten Quadranten statisch gewählt, während sie bei *PReLU* dynamisch von neuronalen Netz während des Trainingsprozesses selbst gelernt werden kann [1].

Eine letzte Variante der Aktivierungsfunktionen beschreibt die *ELU* Funktion (siehe Abbildung 2.6).

Sie besitzt nicht nur die Eigenschaft schwindende Gradienten und damit nicht anpassbare, sogenannte „tote“ LTUs zu verhindern, sondern ist im gesamten Definitionsbereich ebenso eine stetig differenzierbare Funktion, was das Gradientenverfahren beschleunigt. Als Standardwert für den Streckungsfaktor α der niederen Funktion wird oft Eins verwendet. Nachteil der *ELU* Funktion ist der erhöhte Rechenaufwand, was aber durch die schnellere Konvergenz kompensiert wird [1].

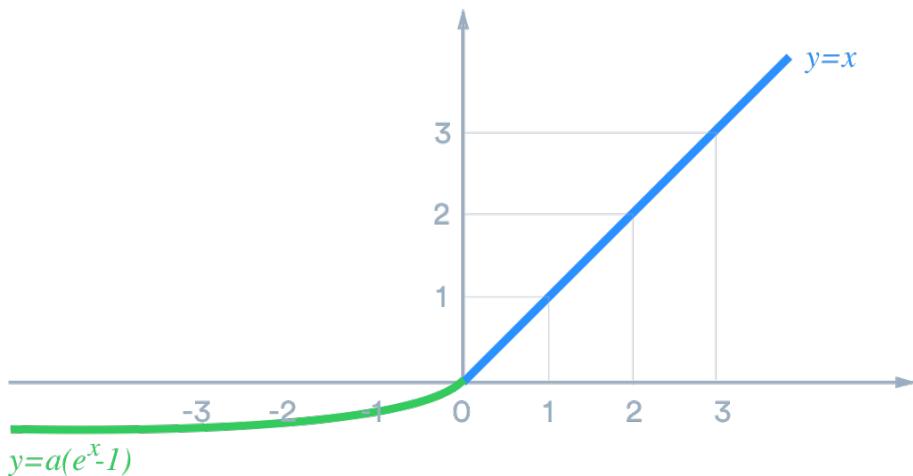


Abbildung 2.6.: ELU-Aktivierungsfunktion [11]

2.4. Datensatzlehre

Datensatzzusammensetzung

Zum Erstellen und Auswählen eines *Deep Learning* Modells wird der Datenbestand in der Regel in drei Kategorien unterteilt. Ein Datensatz wird für das Training des Modells verwendet. Durch das anschließende Anwenden des Modells auf zuvor ungesehene Daten, den Testdaten, wird der *Verallgemeinerungsfehler* gemessen, der möglichst niedrig ausfallen sollte. Fällt der allgemeine Fehler während des Trainings niedrig aus, der *Verallgemeinerungsfehler* während des Testdurchlaufs allerdings hoch, so liegt klassisches *Overfitting* vor, die Trainingsdaten wurden auswendig gelernt [1].

Anschließend werden in mehreren Durchläufen die Hyperparameter des Trainingsprozesses angepasst, sodass letztendlich der *Verallgemeinerungsfehler* für die Testdaten niedrig ausfällt. Kommt es anschließend zum Einsatz des Modells in der Produktivumgebung, so können trotz allem unerwartete Ergebnisse bezüglich des Abstraktionsvermögens des Modells auftreten, was daran liegt, dass das Modell allein auf die Testdaten hin optimiert wurde. Um dies zu vermeiden, wird ein dritter Datensatz, der Validierungsdatensatz, eingeführt. Mehrere Modelle werden dabei durch den Validierungsdatensatz getestet und das am besten abschneidende Modell mit dessen Hyperparametern ausgewählt. Der eigentliche Testdatendatatz wird anschließend nur noch zur Abschätzung des *Verallgemei-*

nerungsfehlers verwendet [1].

Oft wird der Trainingsdatensatz mit dem Validierungsdatensatz zum sogenannten *Trainval* Datensatz zusammengeführt. Dies steht im Kontext des sogenannten *K-Kreuzvalidierungsverfahrens*. Dabei wird der *Trainval* Datensatz in K gleich große, komplementäre Untermengen unterteilt. Eine dieser Untermengen dient anschließend als Validierungsdatensatz. Für jedes zu trainierende Modell mit unterschiedlichen Hyperparametern wird eine andere Untermenge als Validierungsdatensatz ausgewählt. Hierdurch steigt die Aussagekraft des Abstraktionsvermögens nach der Validierung und zudem müssen keine Trainingsdaten dauerhaft für die Validierung zurück gelegt werden. In der Regel werden 80% der Gesamtdaten als *Trainval* Datensatz verwendet [1].

Qualität und Quantität der Daten

Um ein funktionsfähiges Modell zu trainieren, muss der Datensatz einem gewissen Standard nachkommen. Demnach müssen die zu klassifizierenden Objekte vollständig im Bild enthalten und gut erkennbar sein. Zwar gibt es gerade im *PascalVOC* Datensatzformat ebenso die Möglichkeit, Objekte als „schwierig erkennbar“ zu markieren, dennoch sollen solche Objekte nicht die Mehrheit im gesamten Datensatz ausmachen. Auch die Aufnahme von Objekten in unterschiedlichen Umgebungen, Verdeckungsgraden, Entfernung und Blicklagen fördert langfristig das Abstraktionsvermögen des Modells.

Ebenso muss ein ausreichend großer Datensatz vorliegen, um das gewünschte Abstraktionsvermögen des Modells zu erreichen. Die Ergebnisse aus Abbildung 3.4 wurden beispielsweise durch Kombination der *Trainval* Datensätze von PascalVOC 2007 und 2012 erzielt und umfasst 16.551 Bilder im Trainingsverfahren [12, 13, 14].

Unter Hinzunahme des COCO *trainval135k* Datensatzes erreicht der *SSD* sogar das beste Ergebnis aus der ursprünglichen Veröffentlichung mit einer durchschnittlichen Präzision von 81.6% [12].

Techniken zum Trainieren bei geringen Datenmengen

Bei Betrachtung der obigen Ergebnisse wird schnell deutlich, dass für ein komplexeres *Deep Learning* Modell ein umfangreicher Datensatz von Nöten ist. Allerdings gibt es zwei

bekannte Techniken, wie auch mit kleineren Datenbeständen ein sehenswertes Ergebnis erzielt werden kann.

Beim sogenannten *Transfer Learning* können von einem bereits für ein ähnliches Problem trainiertes Modell die ersten Schichten des neuronalen Netzes für das neue Modell wiederverwendet werden. Die übernommenen Gewichtungen werden nicht mit trainiert. Neben einer kleineren Datenmenge zum Trainieren hat das *Transfer Learning* ebenso den Vorteil das Training selbst zu beschleunigen [1].

Eine weitere Technik beschreibt das künstliche Vergrößern des Datensatzes durch affine Transformationen wie Translation, Rotation oder Skalierung und wird *Data Augmentation* genannt [1].

2.5. Grundlagen zu Objektdetectoren

Ein Anwendungsgebiet des *Deep Learnings* beschreiben die Objektdetectoren, die im *Smart Warehouse* Szenario zur Lokalisierung und Klassifizierung von Bestandsobjekten genutzt werden. Im folgenden Kapitel soll demnach der Grundbaustein von Objektdetectoren, das *Convolutional Neural Network*, zunächst genauer betrachtet werden, bevor auf die gängigsten Objektdetectoren, die der *Regional Convolutional Neural Networks* (R-CNN), der *Single Shot MultiBox Detector* und der *You Only Look Once* Ansatz im darauf folgenden Kapitel eingegangen wird. Auch wird die gängiste Metrik zum Vergleich von Objektdetectoren, die *mean Average Precision* eingeführt.

Convolutional Neural Networks

Ein CNN besteht größtenteils aus drei grundlegenden Bausteinen, den sogenannten *Convolutional Layern*, *Pooling Layern* und den bereits bekannten *Fully-Connected Layern*.

Ein Convolutional Layer zeichnen sich unter anderem dadurch aus, dass jede LTU dieser Schicht nicht mit allen vorherigen LTUs der vorgegangenen Schicht verbunden ist, sondern nur mit einer festen, beschränkten Anzahl. Es ist also kein vollständig verbundenes neuronales Netz. Dieser „lokale Wahrnehmungsbereich“ macht es möglich, dass örtliche Informationen und Merkmale im Bild erhalten bleiben. Auch können große Bilder klas-

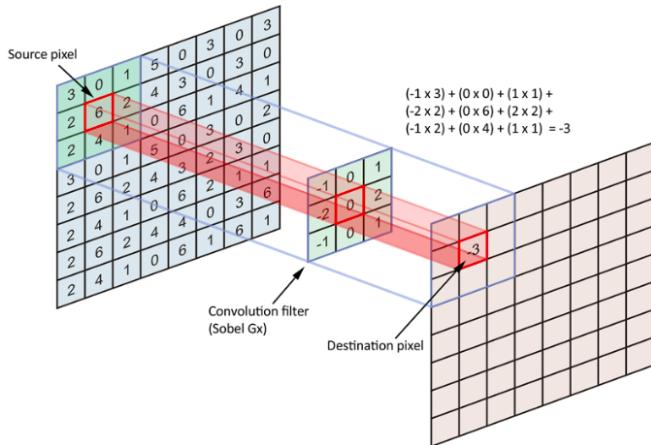


Abbildung 2.7.: Convolutional Layer [15]

0	0	0	0	0	0
0	35	19	25	6	0
0	13	22	16	53	0
0	4	3	7	10	0
0	9	8	1	3	0
0	0	0	0	0	0

Abbildung 2.8.: Zero-Padding [16]

sifiziert werden, ohne dass die Anzahl an nötigen Verbindungen im ANN unüberschaubar groß wächst. Die folgenden LTUs der zweiten Schicht sind ebenfalls wiederum nur mit einem Ausschnitt vorangegangener Neuronen verbunden und fassen die erkannten, kleinteiligen Merkmale der ersten Schicht zu übergeordneten, zusammengesetzten und komplexeren Merkmalen zusammen [1].

Um allerdings Convolutional Layer genauer zu verstehen, ist anstelle einer eindimensionalen Darstellung eines Layers eine dreidimensionale Darstellung besser geeignet.

Zunächst kann ein zweidimensionale Bild als Matrix dargestellt werden, bei der jedes Element der Matrix den Grauwert eines Pixels zwischen 0 und 255 trägt. Die dadurch entstandene zweidimensionale Schicht bildet den Input-Layer mit einer LTU pro Pixel. Anschließend werden auf diese Schicht nacheinander mehrere Filter angewandt, die die Gewichte des CNNs tragen und Muster aus dem Bild extrahieren. Die Stellen, die dem Muster ähnlich sind sollen verstärkt werden, während Stellen, die nicht dem Muster entsprechen durch eine Nullgewichtung ausgelöscht werden sollen. Der Lernprozess bei der Bilderkennung beruht also darauf, die bestmöglichen Filter für die gegebene Aufgabe zu finden und diese für eine Erkennung komplexer Mustern zusammen setzen zu können. In Abbildung 2.7 ist ein 3x3 Pixel Filter mit dessen Anwendung dargestellt. Ein Filter besitzt die Größe des künstlichen Wahrnehmungsbereiches einer LTU [1].

Ein Filter wird dazu verwendet, jeden Pixel der Eingabeschicht auf die folgende Schicht abzubilden. Um keine Informationen zu verlieren und den Filter ebenso auf Randbereich anwendbar zu machen, wird oft ein sogenanntes *Zero-Padding* auf eine Schicht

angewandt, bei dem die Randbereiche mit LTUs des Wertes 0 aufgefüllt werden (siehe Abbildung 2.8) [1].

Falls eine gleich große folgende Schicht gewünscht ist, wird eine Schrittweite (engl.: *stride*) von 1 gewählt. Dies dient vor allem dazu kleinere Strukturen noch zu erkennen. Der Filter wird von einem Pixel zum direkt benachbarten Pixel bewegt und angewandt. In tieferen, fortgeschritteneren Schichten kann die Schrittweite auch größer als 1 gewählt werden, da hier bereits nach dem Anwenden mehrerer Filter feinere Muster erkannt wurden und diese nun zu größeren zusammengesetzt werden. Dabei verkleinert sich die resultierende Schicht [1].

Das Ergebnis der Anwendung eines Filters wird als *Feature-Map* bezeichnet. Da mehrere Filter auf die gleiche Schicht angewandt werden, entstehen ebenso mehrere Feature Maps der Schicht. Werden diese Feature Maps übereinander gelagert vorgestellt, so entsteht der dreidimensionale, „faltungsbedingte“ (engl.: convolutional) Charakter eines Convolutional Layers. Eine Schicht eines Convolutional Layers ist mit den entsprechenden Wahrnehmungsbereichen aller vorhergehenden Feature Maps des vorhergehenden Convolutional Layers verbunden (siehe Abbildung 2.9) [1].

Falls zusätzlich eine Farberkennung gewünscht ist, besitzt der Input Layer für jeden der drei Farbkanäle des RGB-Schemas eine Schicht, die Werte zwischen 0 und 255 in ihren LTUs tragen und den Stärken des Rot-, Grün- und Blaukanals entsprechen [1].

Der zweite Grundbaustein eines CNN sind Pooling Layer. Ähnlich zu den Convolutional Layern ist auch hier jede LTU nur mit einer begrenzten Anzahl an LTUs des vorhergegangenen Layers verbunden, also nur mit dem lokalen Wahrnehmungsbereich. Der Hauptunterschied liegt aber darin, dass keine Filter existieren, die die Werte vorhergehender LTUs unterschiedlich gewichten und dabei Muster erkennen. Statt den Filtern werden Aggregatfunktionen wie *MAX()* oder *MEAN()* dazu verwendet, um die Eingaben in nachfolgende Schichten zu verkleinern. So wird beispielsweise bei einem MAX-Pooling Layer mit Schrittweite größer als 1 der jeweils größte Wert des lokalen Wahrnehmungsbereiches weitergereicht und damit die Eingabe in nachfolgende Schichten verkleinert (siehe Abbildung 2.10), was mit einem Informationsverlust verbunden ist. Diese Verkleinerung des Bildes ist ein wesentlicher Schritt, um bei der Mustererkennung weiter Informationen und Merkmale abstrahieren zu können [1].

Daneben ist ein Pooling über die Tiefe der Feature Maps möglich. Hier bleibt die Grö-

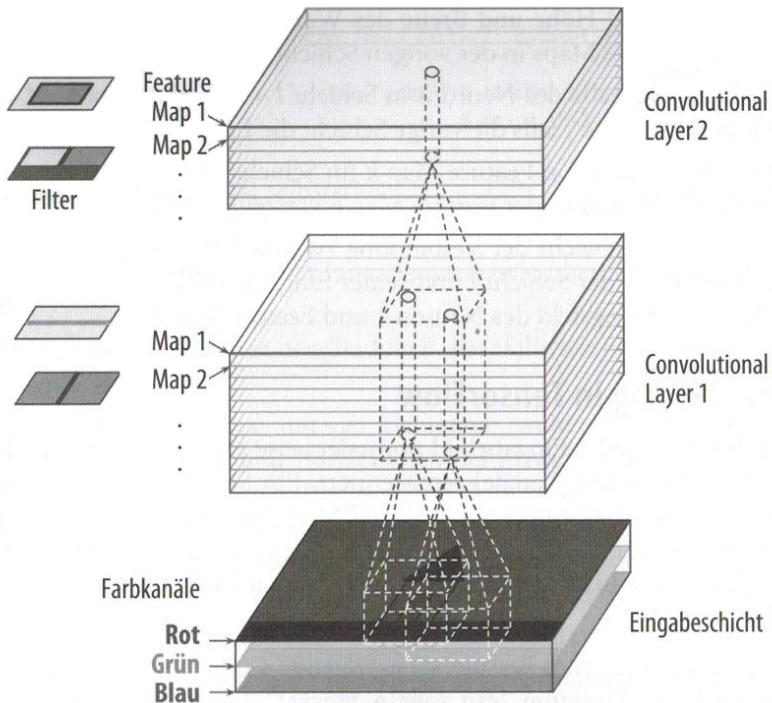


Abbildung 2.9.: Veranschaulichung von Feature Maps [1]

ße der resultierenden Feature Maps gleich, die Anzahl verringert sich allerdings. Die unterschiedlichen Farbkanäle werden somit nach und nach abstrahiert [1].

Nachdem nun beide Grundbausteine eines CNNs genauer erläutert wurden, lassen sich diese nun kombinieren um ein vollständiges CNN zu bauen. Hierbei gibt es unterschiedlichste Architekturen, größtenteils äußerst komplexe. Im Rahmen dieser Arbeit genügt es allerdings, die grundlegende Architektur zu erläutern.

Diese beginnt mit einigen Convolutional Layern, die aufeinander folgen und am Ende durch eine ReLU-Funktion nochmals gefiltert und durch ein Pooling Layer abgeschlossen werden. Dies wird je nach Komplexität der zu erkennenden Muster und der Größe der Bilder einige Male wiederholt. Das ursprüngliche Bild wird durch die Pooling Layer zwar immer kleiner, allerdings auch durch die Convolutional Layer immer tiefer. Das CNN schließt mit einem normalen *Feed-Forward ANN* mit *Fully-Connected Layern* ab, generiert dabei einen *Feature Vektor* und trifft durch eine Softmax-Funktion eine Klassifizierungsaussage des Bildes [1].

Diese Architektur ermöglicht ebenso die Wiederverwendbarkeit einzelner Schichten und

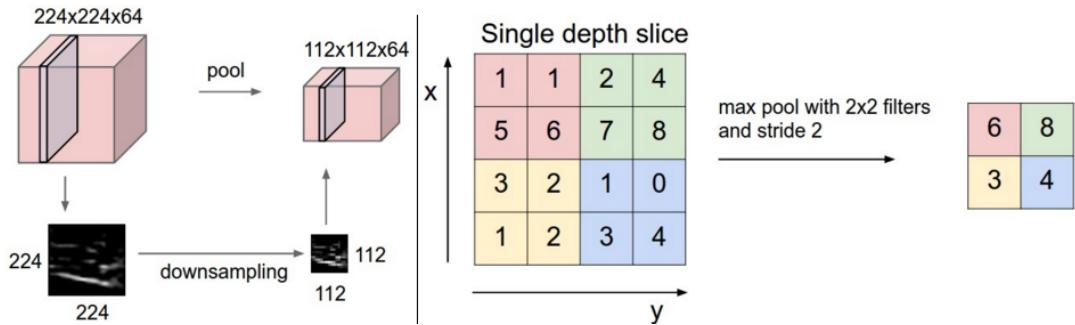


Abbildung 2.10.: Pooling Layer [17]

Gewichtungen für ähnliche Klassifikationsprobleme, bei denen gleiche Muster vorzufinden sind [1].

Mean Average Precision

Um die Genauigkeit von Objektdetektoren zu messen, wird oft die Metrik *mean Average Precision* (mAP) gewählt. Diese setzt sich aus zwei grundlegenden Größen zusammen [18]):

$$\text{Precision} = \frac{\text{TruePositive}}{\text{TruePositive} + \text{FalsePositive}} \quad (2.8)$$

$$\text{Recall} = \frac{\text{TruePositive}}{\text{TruePositive} + \text{FalseNegative}}$$

Precision sagt also etwas über die Verlässlichkeit einer Klassifikation aus, während *Recall* Aussagen über die Erkennungsfähigkeit eines Objektdetektors trifft. Wichtig ist es hierbei anzumerken, dass mehrfach detektierte Objekte nur einmal als positiver Befund aufgefasst werden, die restlichen Detektionen gehen als *False Positives* [19].

Die Klassifikation, ob eine Bounding Box das gewünschte Objekt enthält und demnach ein positiver Fall vorliegt, wird anhand eines sogenannten *Intersection over Union* (IoU) Schwellwertes bestimmt (siehe Abbildung 2.11). Er beschreibt ein Maß der Überdeckung

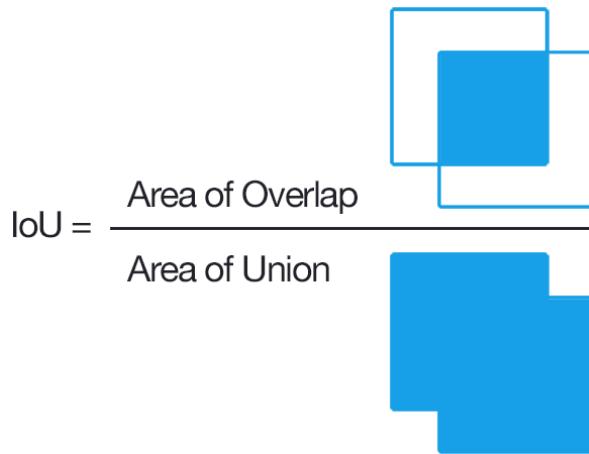


Abbildung 2.11.: Intersection over Union [20]

der detektierten Bounding Box zur wahren Bounding Box und wird auch als *confidence score* bezeichnet. Für den kompletten Datensatz werden nun für unterschiedliche *confidence scores* jeweils *Precision* und *Recall* bestimmt und anschließend in einem Graphen aufgetragen. Meistens werden die *confidence scores* so gewählt, sodass sich eine äquidistante Abstufung in den *Recall* Werten ergibt [19].

Im Graphen ist meist ein klassisches „Zick-Zack“ Muster zu erkennen (siehe Abbildung 2.12). Dieses Muster wird geglättet, indem nach jedem Einbruch für jeden *Recall* Wert der maximale *Precision* Wert rechts des aktuellen *Recalls* übernommen wird. Wird anschließend das diskrete Integral über alle *Recall* Werte gebildet, so ergibt sich der *Average Precision* Wert für eine zu klassifizierende Kategorie. Der Mittelwert der *Average Precisions* über alle Klassifikationskategorien hinweg ergibt letztendlich den *mAP* Wert [18].

2.6. Objektdetektoren

Ein Teilziel der Machbarkeitsstudie ist es, anhand vorbestimmter Kriterien eine ausgewählte Menge von Objektdetektoren zu vergleichen. Um im Laufe der Arbeit zu verstehen, wie diese Auswahl zu Stande kommt und wie sich bestimmte Ergebnisse im Vergleich begründen lassen, ist eine Einführung in die unterschiedlichen Architekturen der Objektdetektoren unumgänglich.

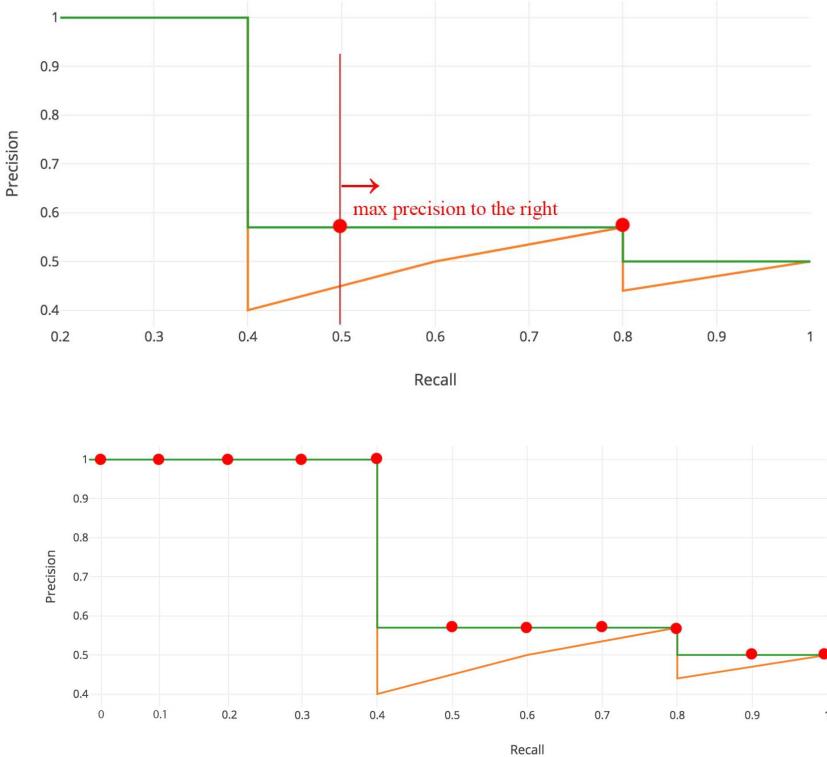


Abbildung 2.12.: Berechnung mAP [18]

Regional Convolutional Neural Networks

Regional Convolutional Neural Networks (R-CNNs) vertreten den Ansatz, für ein Bild mehrere Lokationsvorschläge für mögliche Objekte zu liefern, sogenannte *Regions of Interest* (RoIs), und diese anschließend zu klassifizieren.

Bei dem klassischen *R-CNN* Detektor werden durch den *Selective Search* Algorithmus 2000 solcher *RoIs* vorgeschlagen. Zur Merkmalsextraktion wird für jede *RoI* anschließend ein CNN eingesetzt. Der resultierende *Feature Vektor* wird zur Klassifikation eines Objektes einer *Support Vector Machine* unterzogen. Um zusätzlich die Bounding Boxen akkurat zu bestimmen, wird der *Feature Vektor* zudem einem Bounding Box Regressor unterzogen (siehe Abbildung 2.13) [21].

Da der sogenannte *Region Proposal* Schritt durch den *Selective Search* Algorithmus allerdings viel Zeit in Anspruch nimmt, entstand eine Weiterentwicklung des *R-CNN* Netzes, das *Fast R-CNN* Netz. Dieses tauscht den Schritt des *Selective Search* Algorithmus mit

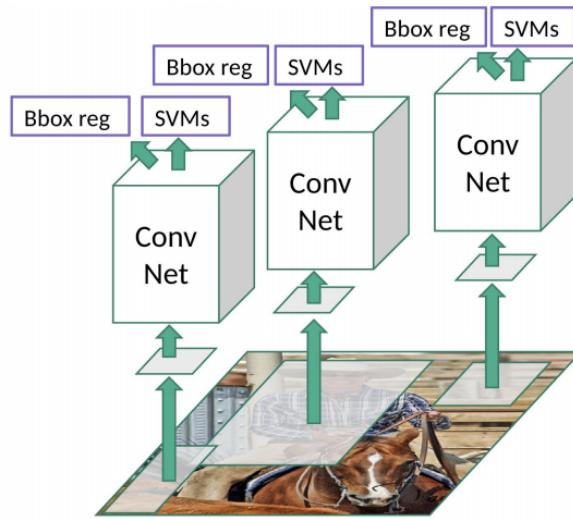


Abbildung 2.13.: R-CNN Architektur [21]

dem Einsatz des CNNs. Außerdem wird das klassische CNN leicht angepasst. Bei *Fast R-CNN* wird ein Bild zunächst einem CNN unterworfen. Bevor eine *Feature Map* durch *Fully-Connected Layer* zu einem einzigen *Feature Vektor* vereinfacht wird, werden aus der *Feature Map* die verschiedenen *RoIs* extrahiert. Dies geschieht wiederum mit dem *Selective Search* Algorithmus, mit dem Unterschied, dass dieser nun nur auf der *Feature Map* operiert und nicht auf dem gesamten Bild. Durch *RoI Pooling Layer* werden die einzelnen entstandenen Regionen in eine feste Größe transformiert und einzeln einer Klassifikation durch *Fully-Connected Layer* und einer Softmax-Funktion unterworfen. Die Komponente mit dem *Bounding Box Regressor* bleibt gleich. Durch den Tausch des CNNs mit dem *Selective Search* Algorithmus werden die mathematischen Faltungsoperationen nur einmal statt 2000 Mal pro Bild ausgeführt, was die Performanz des Detektors gegenüber eines klassischen *R-CNNs* enorm steigert [21].

Die letzte Optimierung der R-CNN Familie entstand durch das *Faster R-CNN* Netz. Dieses ersetzt den statischen *Selective Search* Algorithmus des *Fast R-CNN* Detektors durch ein eigenes lernfähiges, sogenanntes *Region Proposal Network* (RPN) (siehe Abbildung 2.14) [21].

Neben dem Einsatz von RCNN Detektoren zur *Objektdetektion* existiert ebenso ein Ansatz zur *instanzbasierten Segmentierung*, das *Mask R-CNN* Netz. Es nimmt zwei wichtige Anpassungen an der Architektur des *Faster R-CNN* Netzes vor. Da bei Segmentierungs-

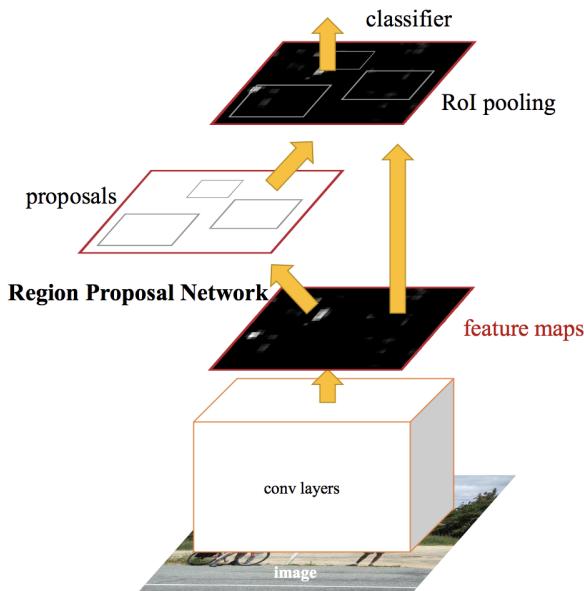


Abbildung 2.14.: Faster R-CNN Architektur [21]

problemen eine genauere Abgrenzung von Objekt und Hintergrund notwendig ist, wird das *RoI Pooling Layer* durch ein *RoI Align Layer* ausgetauscht. Hierbei wird das Rundungsproblem beim Pooling behoben. Angenommen eine *RoI* von 16x16 Pixeln wird mit einem MEAN-Pooling Layer der Schrittweite Drei verarbeitet, so ergibt sich pro Pooling Schritt ein Einzugebiet von 5.33 Pixeln. Dieses wurde abgerundet auf 5 Pixel. Bei *RoI Align Layer* wird durch bilineare Interpolation der Wert des 5.33ten Pixels ermittelt und in das Pooling mit einbezogen. Dies ermöglicht eine genauere Segmentierung an den Grenzen eines Objektes [22].

Außerdem wird parallel zum RPN ein sogenanntes *Fully Convolutional Network* (FCN) eingesetzt, einem Netz, dass rein aus *Convolutional Layern* besteht. Es dient, um für jede existierende Klasse eine pixelbasierte binäre Maske auszugeben, die für jeden Pixel die Zugehörigkeit zu einer Klasse bestimmt. Basierend auf dieser Maske werden die detektierten Objekte anschließend farblich hervorgehoben [23].

Single Shot MultiBox Detector

Zwar liefern die oben genannten Objektdetektoren akkurate Ergebnisse, allerdings sind sie als zu rechenintensiv und langsam einzuordnen, als dass sie für Echtzeit Applikationen eingesetzt werden könnten. Der *Single Shot MultiBox Detector* (SSD) unterscheidet sich von vorhergehenden Modellen, wie beispielsweise den R-CNN Detektoren, dahingehend, dass er bewusst auf den Schritt der Generierung von Bounding Box Vorschlägen und des *Poolings* verzichtet, um wesentlich schneller ablaufen zu können als andere Objektdetektoren. Die Präzision der Klassifikationen bleibt hierbei erhalten, selbst Bilder niedriger Auflösung können weiterhin verarbeitet werden. Dem *SSD* genügt also ein einziges tiefes neuronales Netz zum Lokalisieren und Klassifizieren von Objekten. Wie der *SSD* aufgebaut ist und welche Ansätze er verfolgt, soll in diesem Unterkapitel erläutert werden [12].

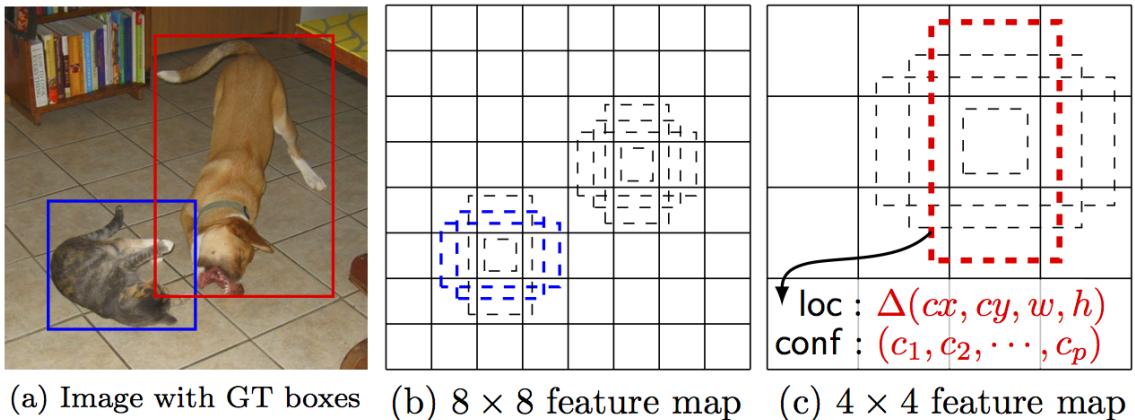


Abbildung 2.15.: SSD Bounding Box Proposals [12]

Die Architektur des *SSD* zielt darauf ab, durch unterschiedlich große *Convolutional Layer Feature Maps* unterschiedlicher Skalierung in die Klassifikation mit einfließen zu lassen. Anschaulich kann es sich vorgestellt werden, als werde das Bild in mehrere unterschiedlich große Gitterstrukturen unterteilt und die resultierenden Zellen jeweils einzeln klassifiziert. Dadurch ist es möglich, Objekte unterschiedlicher Größe zu erkennen. Für jede Zelle im Gitter wird eine gleiche Anzahl vordefinierter Bounding Boxen, die unterschiedliche Seitenverhältnisse aufweisen, definiert. Daher entstammt der Name „MultiBox“. Abbildung 2.15 zeigt beispielweise, wie eine Katze (in blau) und ein im Vergleich zur Katze größerer Hund (in rot) durch unterschiedlich große Gittereinteilungen und Boun-

ding Box Seitenverhältnisse detektiert werden. Durch die *MultiBox* Eigenschaft wird ebenso sichergestellt, dass sowohl horizontal als auch vertikal ausgeprägte Objekte in der selben Zelle gleichzeitig erkannt werden können (siehe Abbildung 2.16) [12].

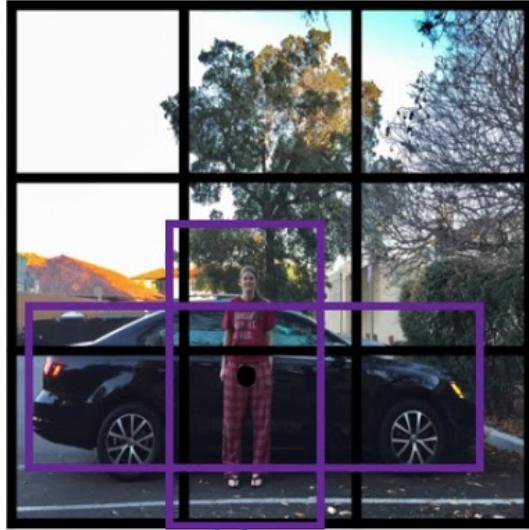


Abbildung 2.16.: Bounding Boxes [24]

Für jede dieser Bounding Boxen bestimmt der *SSD* Wahrscheinlichkeiten für Klassenzugehörigkeiten als auch Verschiebungen der vordefinierten Bounding Box zur wahren Bounding Box des Objekts für jede Klasse. Die Kostenfunktion ist durch die gewichtete Summe des Lokalisationsverlustes und des Klassifikationsverlustes bestimmt. Während der Klassifikationverlust durch eine Softmax-Funktion (2.9) bestimmt werden kann, wird der Lokalisationsverlust über die *Smooth L1* Funktion (2.9) bestimmt. Der Parameter l beschreibt die vorhergesagte Bounding Box, der Parameter g die originale Bounding Box nach den Trainingsdaten [12].

$$\begin{aligned}
 L_{loc}(l^j, g) &= \sum_{j \in Pos}^n \sum_{i \in (x,y,w,h)}^m SM_{L1}(l_i^j - g_i) \\
 SM_{L1}(l_i^j - g_i) &= \begin{cases} 0.5x^2 & \text{wenn } |x| < 1 \\ |x| - 0.5 & \text{sonst} \end{cases} \quad (2.9)
 \end{aligned}$$

Technisch basiert der *SSD* auf der Idee eines *Feed-Forward Convolutional Networks* (sie-

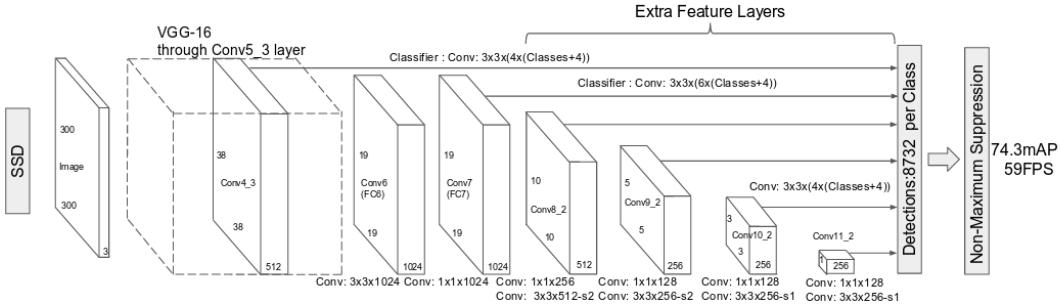


Abbildung 2.17.: SSD Architektur [12]

he Abbildung 2.17). Er benutzt ein *VGG-16* Basis Netzwerk¹, dessen *Fully-Connected Layer* am Ende entfernt wurden. Die resultierende *Feature Map* wird nun einer Reihe von ständig kleiner werdenden *Convolutional Layern* unterzogen. Jedes *Convolutional Layer* kann eine feste Anzahl an Detektionen bestimmen. Eine Detektion wird durch eine Klassenangabe und die Lage einer vorhergesagten Bounding Box bestimmt. Eine Bounding Box wird wie bereits erläutert durch den linken oberen Eckpunkt $P(x, y)$ und eine Höhe und Breite bestimmt. Bei c Klassen hat der *Feature Vektor* einer Detektion demnach die Größe $c + 4$. Bei einer *Feature Map* Größe von $m \cdot n$ und k verschiedenen vordefinierten Bounding Boxen ergeben sich also $m \cdot n \cdot k \cdot (c + 4)$ verschiedene *Feature Vektoren* für eine Feature Map [12]. Diese *Feature Vektoren* werden nun an das Ende des Netzes zur Klassifikation weitergeleitet.

Dieser Vorgang wird für alle *Feature Maps* für alle *Convolutional Layer* durchgeführt. Die daraus folgende Menge an Detektionen wird durch ein *Non Maximum Suppression Layer* in ihrer Größe reduziert. Als Maß zur Filterung wird die *IoU* der detektierten Bounding Box zur wahren Bounding Box verwendet. Überschreitet diese einen Wert von 0.5, so ist diese der originalen Bounding Box zugeordnet. Demnach ist es auch möglich, dass eine originale Bounding Box mehreren vordefinierten Bounding Boxen zugeordnet werden kann [12].

Während des Trainingsprozesses des *SSD300*² wurde eine Lernrate von $\eta = 10^{-3}$ für

¹ *VGG-16* ist ein auf dem Datensatz von *ImageNet* basierendes neuronales Netz, das bis zu 1000 unterschiedliche Kategorien klassifizieren kann [25].

² *SSD300* verwendet Bilder der Auflösung 300x300 Pixel. Alternativ existiert ebenso *SSD512* für Bilder der Auflösung 512x512 Pixel. Die Bilder können jedoch auch kleiner als die vorgegebene Auflösung gewählt werden.

das Mini-Batch Verfahren mit Batchgröße 32 und Moment $\beta = 0.9$ verwendet. Die Gewichtungen wurden *Xavier* initialisiert. Nach 40.000 Iterationen wurde die Lernrate für 10.000 Iterationen auf $\eta = 10^{-4}$ reduziert und schließlich auf $\eta = 10^{-5}$ [12]. Auf Basis der PASCAL VOC Datensätze aus 2007 und 2012 wurde mit dieser Konfiguration eine *mAP* von 74.3% für SSD300 respektive 76.8% für SSD512 erreicht.

You Only Look Once

Der Algorithmus *You Only Look Once* (YOLO) ist ein weiterer Objekterkennungsalgorithmus und betrachtet statt separaten Bildregionen das komplette Bild. Er benutzt nur ein neuronales Netz, um Bounding Boxen und Wahrscheinlichkeiten für bestimmte Klassen vorherzusagen.

Hierzu wird ein $S \times S$ Gitter über das Bild gelegt. Für jedes Feld im Gitter werden B Bounding Boxen erzeugt. Jede Box besitzt neben den zum Gitterfeld relativen Positionswerten einen Wert, der die Vorhersage der jeweiligen Klasse und die Präzision der Box repräsentiert. Dieser Wert wird als *confidence score* bezeichnet und wird durch die Multiplikation der Wahrscheinlichkeit für eine Klasse mit der *IoU*, also die Präzision der berechneten Box im Verhältnis zu der Box aus den vortrainierten Testdaten festgelegt [26].

Aus der Menge an Bounding Boxen werden schließlich mit Hilfe eines festgelegten Schwellwertes die Boxen mit lokalisierten Objekten bestimmt (siehe Abbildung 2.18).

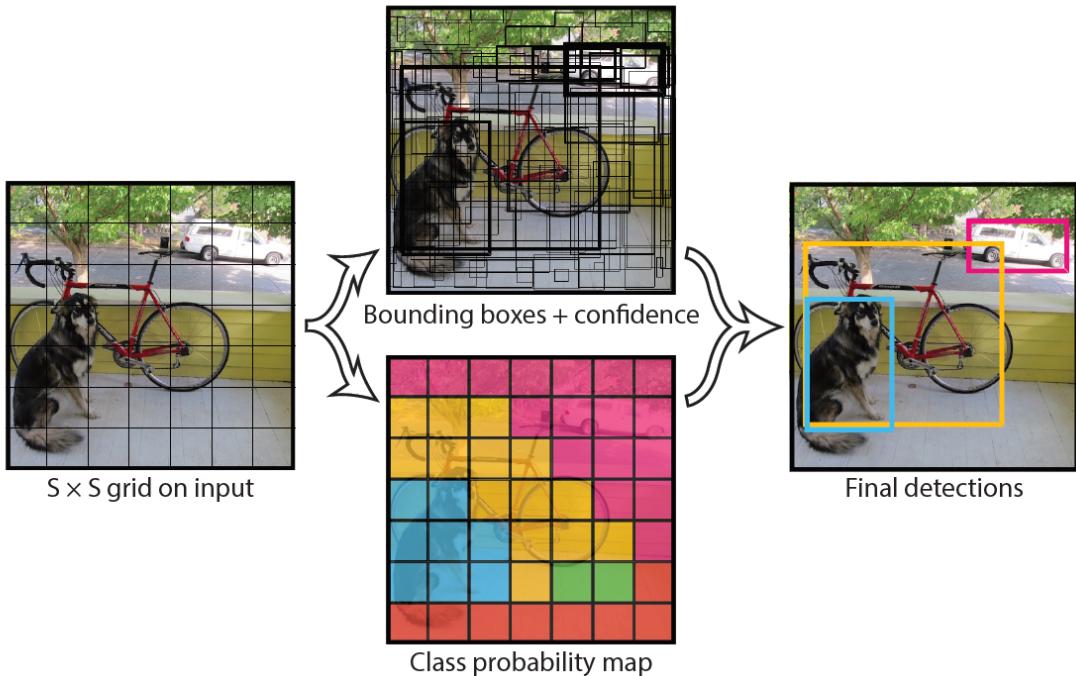


Abbildung 2.18.: Vereinfachte Darstellung des YOLO Algorithmus [26]

Die vorhergesagten Werte werden in einem $S \times S \times (B * 5 + C)$ Tensor kodiert, wobei S und B wie zuvor beschrieben durch das Gitter und die Bounding Boxen festgelegt sind und C die Anzahl der Klassen beschreibt. Abbildung 2.19 zeigt den Aufbau des CNN von *YOLO* für die Detektion. Es besteht aus 24 Convolutional Layern gefolgt von zwei *Fully-Connected* Layern. Für die Genauigkeit bei der Detektion wird die Auflösung des Eingangsbildes verdoppelt [26].

In dem Netzwerk in Abbildung 2.19 wird ein Bild mit einer Auflösung von 224×224 Pixeln verwendet und die vorhergesagten Werte im $7 \times 7 \times 30$ Tensor ausgegeben.

Die mittlerweile dritte und aktuelle Version von *YOLO* weist enorme Verbesserungen auf, gerade im Bezug auf die Erkennung von sehr kleinen Objekten wie zum Beispiel einzelne Vögel in einem Schwarm [27].

TODO Yolo v3: Veränderung neuronales Netz

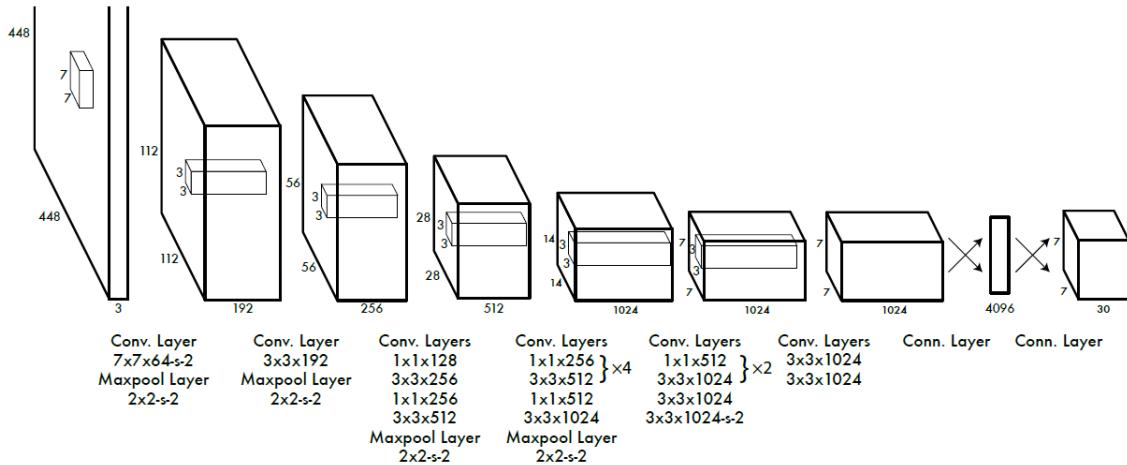


Abbildung 2.19.: YOLO Architektur [26]

2.7. Datensatzformate

Basierend darauf, welcher Objektdetektor trainiert werden soll, muss der zum Training verwendete Datensatz in einem bestimmten Format vorliegen. Zum Trainieren des *SSDs* wird das sogenannte *Pascal Visual Object Classes* (PascalVOC) Format benötigt.

Es definiert eine Unterteilung in *Annotations*, *ImageSets* und *JPEGImages*. Während in dem Ordner *JPEGImages* alle Bilder des Datensatzes vorhanden sind, befindet sich unter anderem die Information über die vorhandenen Objekte in dem Bild im Ordner *Annotations*. Für jedes Bild des Datensatzes werden die Informationen in einer gleichnamigen XML-Datei abgelegt (siehe Listing 2.1).

```

1 <annotation>
2   <folder>JPEGImages</folder>
3   <filename>000001.jpg</filename>
4   <path>..\JPEGImages\000001.jpg</path>
5   <source>
6     <database>Unknown</database>
7   </source>
8   <size>
9     <width>4608</width>
10    <height>2112</height>
11    <depth>3</depth>
```

```

12      </size>
13      <segmented>0</segmented>
14      <object>
15          <name>Saskia Wasser Groß</name>
16          <pose>Unspecified</pose>
17          <truncated>0</truncated>
18          <difficult>0</difficult>
19          <bndbox>
20              <xmin>1575</xmin>
21              <ymin>95</ymin>
22              <xmax>3163</xmax>
23              <ymax>635</ymax>
24          </bndbox>
25      </object>
26  </annotation>

```

Listing 2.1: PascalVOC Bildannotation

Neben allgemeinen Metainformationen über das Bild befindet sich hier ebenso eine Liste aller markierten Objekte. Pro Objekt wird die Klassifikationskategorie, die Ausrichtung (z.B. „Frontal“), die Information über vollständiges Erscheinen im Bild, die Information über schwere Erkennbarkeit und die Bounding Box angegeben. Im Ordner *ImageSets/-Main* wird eine Unterteilung in Trainings- und Testdatensatz durch zwei Textdateien realisiert, die die Dateinamen der Bilddateien als Auflistung enthalten [28].

Das *YOLO* Format für den *YOLO* Objektdetektor definiert in einer *.names*-Datei alle im Datensatz vorhandenen Kategorien durch simple Auflistung der Bezeichner. Die Bilder werden zusammen mit ihren Annotationen in einem separaten Ordner abgelegt. Die Annotationen folgen hier dem Format:

< Kategorie – ID > < Zentrum – X > < Zentrum – Y > < Breite > < Hoehe >

Die Unterteilung in Trainings- und Testdatensatz erfolgt durch Referenzierung der Bildpfade in zwei getrennten Textdateien. Schließlich wird in einer *.data*-Datei der Pfad zu den beiden Textdateien und zur *.names*-Datei sowie die Anzahl an Kategorien gespeichert [29].

2.8. Cloud Infrastruktur

Das Trainieren eines *Deep Learning* Modells ist gerade bei großen CNN Architekturen rechenaufwendig. Tensor Operationen wie Matrixmultiplikationen und Konvolutionen erfordern im Rahmen des maschinellen Lernens hohe Parallelisierung und Taktfrequenzen, um in absehbarer Zeit gute Ergebnisse zu liefern. Die Rechenkapazität normaler Desktop-PCs reicht meist nicht aus, um performantes *Deep Learning* betreiben zu können.

Abhilfe bieten Software-as-a-Service (SaaS) bzw. Platform-as-a-Service (PaaS) Angebote wie *Amazon SageMaker*, *Google Cloud Platform Cloud AI* oder *Azure ML Services* oder aber auch Start-ups wie *FloydHub*. Diese bieten Infrastruktur in unterschiedlichen Zonen je nach Standpunkt der Rechenzentren zum Trainieren an sowie eine Plattform zum Verwalten der *Deep Learning* Prozesse.

Trainingshardware

Gerade GPUs bieten sich aufgrund ihres hohen Parallelisierungsvermögens gegenüber herkömmlichen CPUs an. Insbesondere *NVIDIA* nimmt hierbei eine Vorreiterrolle in der Produktion von Server-GPUs ein. Die *Compute Unified Device Architecture* (CUDA) von *NVIDIA* ermöglicht als Programmiermodell und paralleler Computing Plattform das Auslagern von Rechenprozessen auf GPUs. Das *CUDA* Toolkit beinhaltet GPU beschleunigte Bibliotheken, einen Compiler, Entwicklungswerkzeuge sowie die eigentliche *CUDA* Laufzeit und wird von vielen *Deep Learning* Bibliotheken genutzt, wie z.B. *PyTorch* [30, 31]. Hauptvergleichskriterien zwischen GPUs sind hierbei der Grad der möglichen Parallelisierung und die reine Rechenleistung im Verhältnis zum Stromverbrauch.

Neben GPUs existieren seit 2015 die von Google entwickelten *Tensor Processing Units* (TPUs). Diese sind speziell entwickelte, anwendungsspezifische integrierte Schaltung (engl.: *Application-Specific Integrated Circuit*) (ASIC) für Arbeitslasten im maschinellen Lernen [32].

Eine weitere Steigerung versprechen Microsofts Field Programmable Gate Arrays (FPGAs), die allerdings nicht weiter im Rahmen dieser Arbeit betrachtet werden sollen [33].

Amazon Web Services SageMaker

Amazon Web Services (AWS) bietet mit *SageMaker* eine integrierte Plattform zum Trainieren und Bereitstellen von *Deep Learning* Modellen. Zentrales Alleinstellungsmerkmal ist das einheitliche Toolset, in dem alle Arbeitsprozesse rund um ein *Deep Learning* Modell integriert abgebildet werden können. Es ist somit nicht mehr nötig, unterschiedliche Tools und Arbeitsabläufe zusammenfügen, was zuvor zeitaufwändig und fehleranfällig war [34].

Außerdem bietet *AWS SageMaker* die ersten vollständig integrierte Entwicklungsumgebung für Machine Learning, „*Amazon SageMaker Studio*“. Zum Erstellen der *Deep Learning* Modelle werden sogenannte *Amazon Sagemaker Notebooks* genutzt, eine Ableitung klassischer *Jupyter Notebooks*. Unterstützte Frameworks sind TensorFlow, PyTorch, Apache MXNet, Chainer, Keras, Gluon, Horovod, Scikit-Learn und Deep Graph Library [34].

AWS SageMaker bietet verschiedene Instanztypen an, die sich je nach Anzahl an vCPUs und GPUs unterscheiden. Auch der vorhandene Arbeitsspeicher, Grafikkartenspeicher und die Netzwerkleistung kann durch die Vielzahl an angebotenen Instanzen nach individuellen Bedürfnissen gewählt werden [35].

Google Cloud Platform AI Platform

AI Platform ist das Konkurrenzprodukt zu *AWS SageMaker* von der *Google Cloud Platform* (GCP). Die Plattform bietet ebenso verschiedene Komponenten für das *Deep Learning* an. Hierzu gehören *AI Platform Notebooks*, ein Dienst mit einer integrierten JupyterLab-Umgebung, *Deep Learning* Virtual Machines (VM) mit vorinstallierten *Deep Learning* Frameworks, verteiltes Training mit automatischer Hyperparameter-Abstimmung durch den *AI Platform Training* Dienst oder *AI Platform Prediction* zum Bereitstellen trainierter Modelle [36].

Hervorzuheben sind allerdings Googles TPU Hardwarebeschleuniger, die für Projekte im *TensorFlow* Framework für jede Instanz mobilisiert werden können. Sie sind darauf ausgelegt ein optimales Preis-/Leistungsverhältnis beim Trainingen von *Deep Learning* Modellen zu erreichen [36].

Google Colab

Google Colaboratory, kurz *Google Colab*, ist eine kostenfreie, cloudbasierte *Jupyter Notebook* Umgebung von Google. Dokumente, die in Google Colab erstellt werden, werden automatisch mit *Google Drive* synchronisiert. Die Laufzeit ist frei konfigurierbar zwischen Python 2 und 3 bzw. zwischen einfachem CPU, GPU oder TPU Computing. Nachteil an dem kostenfreien *Google Colab* ist, dass zugewiesene Hardwareresourcen mit weiteren Nutzern geteilt werden müssen und so nicht die volle Rechenleistung für den individuellen Entwickler zur Verfügung stehen [37].

Auch können keine längerfristigen Trainingsjobs ausgeführt werden, ohne dass nach 90 Minuten der Client von dem zugewiesenen Server getrennt wird [38].

Microsoft Azure

Microsoft Azures Angebot für *Deep Learning* in der Cloud ist zunächst wenig transparent. Sie bieten ebenso wie Amazon und Google das Trainieren und Bereitstellen von *Deep Learning* Modellen an und zudem eine einige DevOps Landschaft für solche Arbeitsprozesse. Auch werden Frameworks wie *TensorFlow* oder *PyTorch* unterstützt sowie das Programmieren in *Jupyter Notebooks* [39].

FloydHub

FloydHub, ein kalifornisches Start-up, bietet eine Data-Science Plattform zum Trainieren und Bereitstellen von *Deep Learning* Applikationen. FloydHub erlaubt es Anwendern, sich auf reines *Deep Learning* zu konzentrieren, während es die Arbeit rund um den *Deep Learning* Lebenszyklus abnimmt. Hierzu gehört das Bereitstellen der entsprechenden Hardware, das Installieren von Treibern oder das Integrieren verschiedener *Deep Learning* Bibliotheken, wie *TensorFlow*, *PyTorch* oder *Keras* [40].

Mit Hilfe des von FloydHub bereitgestellten Command Line Interfaces (CLI) kann ein lokales Projekt zu einem FloydHub Projekt initialisiert werden. Anschließend können anhand einer Konfigurationsdatei Einstellungen über das Training spezifiziert werden (siehe Listing 2.2). Alternativ können diese auch über das CLI festgelegt werden.

```

1 machine: gpu2
2 env: pytorch-1.4
3 input:
4   - destination: input
5     source: <username>/datasets/smartwarehousessd/3
6     - <username>/datasets/smartwarehousessd/3:ssd
7 description: Job to train the SSD
8 max_runtime: 3600
9 command: python train.py

```

Listing 2.2: Konfigurationsdatei zum Trainingsjob auf FloydHub

Hierbei kann zwischen der K80 (gpu) oder der V100 (gpu2) GPU für das Training gewählt werden. Auch die *Deep Learning* Laufzeitumgebung muss spezifiziert werden. Bei Bedarf auch zusätzliche Bibliotheken in einer *floyd_requirements.txt*-Datei zur Installation mit angegeben werden. Anschließend muss der Datensatz referenziert werden, mit dem das Modell trainiert werden soll.

Dieser Datensatz wird separat hochgeladen, da sich dieser im Gegensatz zum Programmcode nur selten ändert. FloydHub implementiert auf seiner Plattform eine Art Pfadsystem, unter dem Datensätze und Projekte abgespeichert werden. Diese Pfade werden in der Konfigurations-Datei zur Referenzierung genutzt. Um auch im Programmcode auf den Datensatz zuzugreifen, muss ein Mountname definiert werden. In obigen Beispiel wird dem Datensatz unter Verzeichnis *<username>/datasets/smartwarehousessd/3* der Mountname *ssd* gegeben. Das Verzeichnis zum Einlesen der Daten ist anschließend im Code unter */floyd/input/ssd/* erreichbar.

Mit dem CLI Befehl *floyd run* wird der Programm Code auf die Plattform hochgeladen und der in der Konfigurationsdatei angegebene Befehl ausgeführt. Daraufhin wird ein Job erstellt, versioniert und ausgeführt. Während der Job ausgeführt wird, wird dem Nutzer ein Einblick in die Konsolenausgabe gewährt sowie in Metriken zur Hardwareauslastung. In der Jobhistorie kann im Nachhinein jeder Job mit dem damals aktuellen Programmcode und Datensatz eingesehen werden. Auch Datensätze werden versioniert. Schreibrechte sind auf das Verzeichnis */floyd/home* begrenzt, hier können Zwischenspeicherpunkte des Modells abgelegt werden.

Neben klassischen Trainingsjobs können *Deep Learning* Modelle auch ganz einfach in

Jupyter Notebooks erstellt werden. Hierzu muss in einem Projekt ein Workspace angelegt werden.

3. Konzeption

Um eine Basis zur Umsetzung des *Smart Warehouse* Szenarios zu schaffen, sind zunächst einige konzeptionelle Überlegungen notwendig, die in diesem Kapitel betrachtet werden sollen. Sie beschränken sich im Wesentlichen auf sechs übergeordnete Themenbereiche, die in Abbildung 3.1 dargestellt sind.

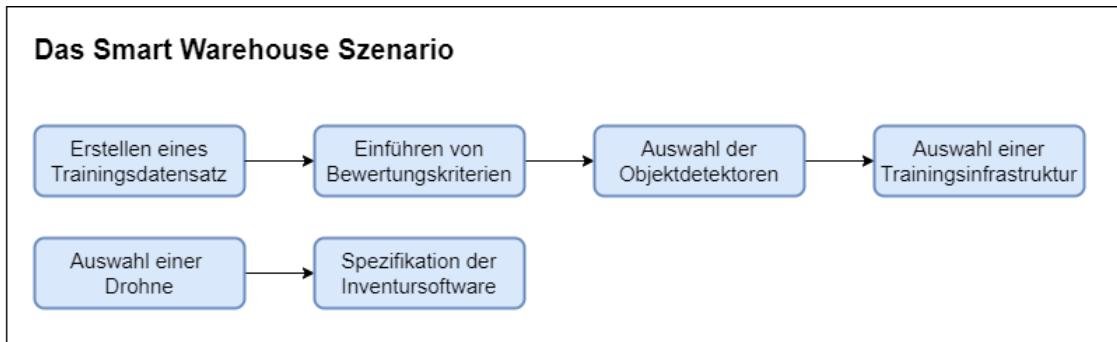


Abbildung 3.1.: Konzeptionelle Schritte

Daraus ergeben sich zwei Neuheitswerte. Zum einen soll evaluiert werden, ob die bestehenden Objektdetektoren für industrielle Anwendungsszenarien grundsätzlich geeignet sind, zum anderen, ob das spezifische Anwendungsszenario zur automatisierten Durchführung einer Inventur von Warenhäusern mit einer Drohne prototypisch umsetzbar ist.

Für das *Smart Warehouse* Szenario ist demnach zum einen die Entwicklung eines *Deep Learning* Modells notwendig, das auf das Anwendungsszenario einer Inventur von Warenhäusern spezialisiert ist. Hierzu muss ein geeigneter Trainingsdatensatz erstellt werden, auf dem die ausgewählten Objektdetektoren trainiert werden sollen. Um die Objektdetektoren miteinander vergleichen zu können, ist zudem die Einführung von Bewertungskriterien notwendig. Auf Basis dieser Kriterien wird später evaluiert, ob die ausgewählten Objektdetektoren sich allgemein für einen industriellen Einsatz eignen. Welche Objektdetektoren überhaupt im Sinne des *Smart Warehouse* Szenarios miteinander verglichen werden sollen, ist im nächsten Schritt „Auswahl der Objektdetektoren“ zu beschließen. Anschließend stellt sich nur noch die Frage, auf welcher Infrastruktur die Objektdetektoren trainiert werden sollen.

Um das zweite Ziel der Arbeit zu erarbeiten, kann auf den Ergebnissen des ersten Ziels aufgebaut werden. Eine Drohne soll die benötigten Daten zur Inferenz liefern. Welche Drohne dies sein soll, ist im ersten Schritt zu erarbeiten. Die Drohne soll über eine Webapplikation ansprechbar sein, auf der zudem die Ergebnisse der Inventur präsentiert werden sollen. Diese Inventursoftware ist im letzten Schritt in ihrem Umfang zu spezifizieren.

3.1. Erstellen eines Trainingsdatensatzes

Das *Smart Warehouse* lehnt sich an ein großes Warenhaus an, bei dem Produkte nicht in Kartons verpackt, sondern als ganzes auf Regalen angeordnet sind, ähnlich wie bei Warenhäusern wie *Baumarkt* oder *Selgros*. Bei dem Aufbau des Trainingsdatensatzes hat die Machbarkeitsstudie allerdings nicht zum Ziel, ein solches Warenhaus vollständig im Datensatz abzubilden, sondern wesentlich den Datensatz so umfangreich zu wählen, um eine generelle Umsetzbarkeit des *Smart Warehouse* Szenarios zu beweisen. Im Rahmen des Projektes wurde sich deshalb exemplarisch auf Getränkeflaschen eines Warenhauses konzentriert. Dabei wurden neun Kategorien festgelegt (siehe Abbildung 3.2).

Der Datensatz besteht aus 1078 manuell annotierten Bildern und dient später dazu die ausgewählten Objektdetektoren zu trainieren. Die Bilder besitzen eine Auflösung von 2112x4608 Pixeln mit einer Farbtiefe von 24 Bit. Alle neuen Kategorien sind nahezu gleich häufig im Datensatz vorhanden. Durch das Tool *LabelImg* wurden die Daten sowohl für das *PascalVOC*, als auch das *YOLO* Format annotiert. Im initialen Datensatz sind auf 75% der Bilder die Objekte der jeweiligen Kategorien einzeln und klar erkennbar abgebildet. Hierdurch wird erhofft, dass Modell zunächst auf die Muster der jeweiligen Objekte zu trainieren. In 12,5% der Bilder sind die Objekte der jeweiligen Kategorien ebenso einzeln, allerdings mit unterschiedlichen Hintergründen, Beleuchtungsverhältnissen, Blickwinkeln und Entferungen abgebildet. Je nach Umgebung wurden Bilder dieses Anteils als schwer erkennbar markiert. Um das Warenhaus zu simulieren, sind in den letzten 12,5% der Bilder die Objekte auf Regalen angeordnet, jeweils hintereinander oder in Getränkekästen (siehe Abbildung 3.3).



Abbildung 3.2.: Die neun Datensatz-Kategorien

3.2. Einführen von Bewertungskriterien

Um Objektdetektoren miteinander vergleichbar zu machen und um deren Potential zum industriellen Einsatz zu bewerten, müssen konkrete Bewertungskriterien eingeführt werden.

Präzision

Zur Messung der Präzision wird die Metrik *mAP* verwendet. Dies garantiert eine gute Vergleichbarkeit mit den veröffentlichten Leistungsmerkmalen der Objektdetektoren.

Reaktionsvermögen

Um eine Verarbeitung in Echtzeit zu ermöglichen, muss gewährleistet sein, dass die Inferenzgeschwindigkeit mit dem Modell mit der eingehenden Bildrate einhergeht. Als



Abbildung 3.3.: SmartWarehouse Regal

Maßstab dafür dient die *Frames Per Second* (FPS) Zahl. Echtzeitfähigkeit in der Machbarkeitsstudie ist so definiert, dass die Inferenz mit dem Modell mindestens so schnell ablaufen muss, dass Änderungen in der Umgebung rechtzeitig von Objektdetektor noch wahrgenommen werden können.

Trainingsverhalten

Unter dem Punkt Trainingsverhalten wird zusammengefasst, wie schnell sich die einzelnen Modelle mit den unterschiedlichen Objektdetektoren trainieren lassen. Hierbei wird besonderer Fokus darauf gelegt, wie viele Trainingsepochen notwendig sind, bis der Gradient der Fehlerfunktion des neuronalen Netzes keine merkenswerten Fortschritte auf Basis des verwendeten Datensatzes mehr erzielt. Es soll aber auch betrachtet werden, wie mit doppelt erkannten Objekten während des Trainingsprozesses umgegangen wird.

Inferenzverhalten

Im Zuge der Evaluation des Inferenzverhaltens werden drei Kriterien betrachtet.

- Das Verhalten bei besonderen Beleuchtungsverhältnissen wie unterbeleuchteten

oder überbeleuchteten Gegenden.

- Das Verhalten bei extremen Blicklagen auf Basis der Entfernung und des Winkels zum detektierenden Objekt.
- Das Verhalten bei nicht vollständig sichtbaren Objekten, z.B. bei Verdeckung.

3.3. Auswahl der Objektdetektoren

Für das *Smart Warehouse* Szenario soll eine Auswahl zwischen den vier Detektoren *Faster R-CNN*, *Mask R-CNN*, *SSD* und *YOLO* getroffen werden. Als Vergleichsbasis dienen die bereits veröffentlichten Benchmarkergebnisse. Zur Evaluation der Machbarkeitsstudie werden die zuvor eingeführten Bewertungskriterien auf die aus dieser Auswahl resultierenden Objektdetektoren angewendet.

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 600
Fast YOLO	52.7	155	1	98	448 × 448
YOLO (VGG16)	66.4	21	1	98	448 × 448
SSD300	74.3	46	1	8732	300 × 300
SSD512	76.8	19	1	24564	512 × 512
SSD300	74.3	59	8	8732	300 × 300
SSD512	76.8	22	8	24564	512 × 512

Abbildung 3.4.: Vergleich SSD auf PascalVOC [12]

Abbildung 3.4 sind die Referenzergebnisse aus der wissenschaftlichen Veröffentlichung des *SSDs* [12]. Die Ergebnisse zeigen, wie die Objektdetektoren *SSD*, *YOLO*, *Fast YOLO* und *Faster R-CNN* untereinander abschneiden. Jeder der Detektoren besitzt eigene Charakteristika bezüglich der benötigten Auflösung der zu verarbeitenden Bilder, der Anzahl der generierten Bounding Boxen und der Batch Größe während des Trainings. Nach diesen Ergebnissen ist eindeutig festzustellen, dass der *SSD* bezüglich *mAP* mit 74.3% bzw. 76.8% am besten abschneidet. *Faster-RCNN* kann zwar mit 73.2% bezüglich der *mAP* mithalten, ist allerdings mit nur 7 FPS nicht zur schnellen Inferenz ausgelegt. *YOLO* schneidet in beiden Kategorien schlechter als der *SSD* ab, er erzielt wesentlich eine *mAP* von 66.4% und eine Framerate von 21 FPS. Dem *SSD* gelingt es also, ein gutes Verhältnis zwischen Präzision und Reaktionsvermögen zu bewahren. Durch den

Verzicht auf den Schritt der Generierung von Bounding Box Vorschlägen und des *Poolings* kann *SSD* deutlich schneller ablaufen als die Vergleichsdetektoren, während durch das Vordefinieren von Bounding Boxen ebenso eine hohe Präzision erzielt werden kann [12].

Allerdings ergibt sich vor allem für kleine Objekte ein erschwertes Detektionsvermögen, da diese in den höherliegenden Convolutional Layern untergehen. Als Lösung hierfür kann eine erhöhte Inputgröße gewählt werden (vgl. *SSD512*) oder *Data Augmentation* für den Lernprozess angewandt werden [12].

Diese Probleme treten bei Netzen der *R-CNN* Familie nicht auf. Wird der *Faster-RCNN* Objektdetektor mit *Mask R-CNN* zur *instanzbasierten Segmentierung* auf Basis des *Common Objects in Context* (COCO) Datensatzes verglichen, so ergibt sich für *Mask R-CNN* mit 38.2% mAP nur eine geringe Verbesserung gegenüber *Faster R-CNN* mit 37.3%. Für diesen Benchmark wurde wohlberichtet das *RoI-Pooling Layer* des *Faster R-CNN* mit einem *RoI-Align Layer* zur besseren Vergleichbarkeit mit dem *Mask R-CNN* Objektdetektor ausgetauscht. Dennoch bleibt auch beim *Mask R-CNN* das Problem eines langsameren Inferenzverhaltens gegenüber dem *SSD* oder *YOLO* offen. Für einen generell industriellen Einsatz könnte dieses langsame Inferenzverhalten möglicherweise problematisch sein, doch für das konkrete Szenario von stehenden Getränkeflaschen in der Machbarkeitsstudie ist eine starke Gewichtung der FPS Metrik zunächst mit Vorsicht zu betrachten [41].

Ein weiteres Auswahlkriterium stellt dar, wie gut die Detektoren aufgesetzt und auf eigens erstellte Datensätze umkonfiguriert werden können. Nach Betrachtung mehrerer Repositories ließen sich der *YOLO* und *SSD* Objektdetektor einfach aufsetzen und auf eigene Datensätze anpassen, stellten bei *Faster R-CNN* und *Mask R-CNN* vermehrt auf Probleme gestoßen wurde. So ist Facebooks Implementierung von *Mask R-CNN* „*Detectron*“ beispielsweise nur auf Linux oder macOS lauffähig. Abgeleitete Repositories sind bereits als *deprecated* deklariert und werden nicht mehr gewartet. Ein manuelles Aufsetzen dieser Implementierungen ist nur unter großem Aufwand möglich und wurde aufgrund der limitierten Zeit abgebrochen. Auch die Referenzimplementierung von *Faster R-CNN* ist bereits als *deprecated* deklariert und verweist auf die *Mask R-CNN* Implementierung *Detectron*. Nebenläufige Implementierungen sind ebenso als *Legacy* Implementierungen vermerkt und nur zeitaufwändig manuell aufsetzbar, sofern sie die Windows-Plattform unterstützen.

Aufgrund dieser Umstände und des schlechteren Abschneidens in der zeitkritischen Modellinferenz wurden *YOLO* und *SSD* als die beiden Detektoren ausgewählt, die im *Smart Warehouse* Szenario genutzt werden sollen. Für *YOLO* wird die Referenzimplementierung im *Darknet* Framework gewählt. Dieses lässt sich ebenso einfach auf eigene Datensätze anpassen im Gegensatz zur Referenzimplementierung des *SSD* im *Caffe* Framework. Deswegen wurde sich bei dem *SSD* für eine Custom Implementierung in *PyTorch* entschieden.

3.4. Auswahl der Trainingsinfrastruktur

Bei der Auswahl der Trainingsinfrastruktur wurden zunächst die Cloud PaaS-Angebote in Betracht gezogen. Diese ermöglichen meist eine weit bessere Performance als lokales Training. Wichtig bei der Auswahl war hierbei

- möglichst niedrige Betriebskosten,
- ein diverses Angebot an Hardware-Beschleunigern und
- ein einfaches Aufsetzen der Trainingsinfrastruktur.

Insbesondere sollten die Testversionen der jeweiligen Angebote zu Nutze gemacht werden, um niedrige Betriebskosten zu erreichen. In Tabelle 3.1 sind die Ergebnisse der Untersuchung dargestellt.

	AWS	GCP	Azure	FloydHub	Colab
Nutzungsrahmen	50 Std.	300\$	200\$	2 Std.	Kostenlos
Hardware-Beschleuniger	Nein	Ja	Ja	Ja	Ja
Setup-Komplexität	Hoch	Mittel	Mittel	Einfach	Einfach

Tabelle 3.1.: Vergleich der SaaS-Angebote der Cloud Anbieter

Amazon SageMaker bietet hierbei für 50 Stunden eine *ml.m4.xlarge* Instanz für Modelltrainingszwecke an [42]. Da diese allerdings nur 4 vCPUs und 16 GiB Arbeitsspeicher umfasst, also keinerlei Cloud GPU als Hardwarebeschleuniger angeboten wird, wurde das Angebot wieder verworfen [35].

Auf Empfehlung wurde anschließend die *GCP* betrachtet. Diese bietet mit 300\$ Startguthaben für 12 Monate ein lukratives Angebot zum Ausprobieren von beliebigen *GCP* Produkten [43]. Die Benutzung der *Deep Learning VM* bietet zudem eine native Unterstützung des *PyTorch* Frameworks, was von der *SSD* Implementierung genutzt wird, und zugleich eine Auswahl aus vier gängigen Cloud GPUs, der *NVIDIA Tesla K80*, *NVIDIA Tesla P100*, *NVIDIA Tesla T4* und der *NVIDIA Tesla V100*. Um die Konfiguration der *Deep Learning VM* allerdings mit Auswahl einer Cloud GPU abschließen zu können, muss zunächst das mit dem Account verknüpfte Kontingent erhöht werden. Hierzu konnte an das *GCP* Support Team ein offizieller Antrag gestellt werden. Aufgrund der geringen Kaufhistorie wurde der Antrag allerdings abgelehnt.

Microsoft Azure bietet für 200\$ bei einer Laufzeit von 30 Tagen Zugang zu allen *Microsoft Azure* Diensten [44]. Darunter gehört eine *NC6* Instanz mit sechs vCPUs und einer *NVIDIA Tesla K80* [45]. Da *Microsoft Azures* Angebot allerdings nur sehr oberflächlich beschrieben wurde, wurde sich letzten Endes auch gegen *Microsoft Azure* entschieden.

Als letzter Anbieter wurde *FloydHub* getestet. Hervorzuheben ist die besonders einfache Vorgehensweise bei der Account Erstellung und dem Aufsetzen der Trainingsinfrastruktur, was bereits im Grundlagen Kapitel beschrieben wurde. *FloydHub* bietet 20 Stunden CPU Trainingszeit bzw. 2 Stunden GPU Trainingszeit auf einer *NVIDIA Tesla K80* [46]. Neben einer *NVIDIA Tesla K80* konnte ebenso Trainingszeit auf einer *NVIDIA Tesla V100* erworben werden. Zudem wurde das verwendete *PyTorch* Framework unterstützt. Aufgrund der einfachen Handhabung wurde sich trotz der erhöhten Kosten für *FloydHub* entschieden.

Während des Trainings mit der *NVIDIA Tesla K80* fiel allerdings auf, dass die Wahl dieser GPU keine großen Performance Verbesserungen brachte. Dies veranlasste eine Gegenüberstellung gängiger Cloud GPUs mit lokalen GPUs, allen voran den bereits vorhandenen Desktop-Grafikkarten *GeForce GTX 1080* und *Titan RTX* (siehe Tabelle 3.2) [47].

	K80	P100	T4	V100	GTX 1080	TITAN RTX
CUDA Cores	2496	3584	2560	5120	2560	4608
Tensor Cores	/	/	320	640	/	576
TeraFLOPS (Single Precision)	4,113	9,526	8,141	14,13	8,873	16,31
Memory Bandwidth (GB/sec)	240,6	732,2	320	897	320,3	672
Suggested Power Supply Unit	700	600	350	600	450	600

Tabelle 3.2.: Vergleich von GPUs nach Rechenleistung

Hierbei fällt auf, dass im Grad der Parallelisierung eine *NVIDIA Tesla K80* zwar mit den vorhandenen lokalen Grafikkarten mithalten kann, in der Anzahl an Rechenoperationen pro Sekunden allerdings weit schlechter abschneidet. Damit sich das Training in der Cloud nach Performance lohnt, muss demnach mindestens eine *NVIDIA Tesla V100* verwendet werden. Da diese allerdings mit 42\$ für zehn Stunden mehr als dreimal so teuer als eine *NVIDIA Tesla K80* für 12\$ ist und zusätzlich zu den GPU Kosten noch monatliche Account-Gebühren berechnet werden¹, wurde sich nach nun nach Kosten-Nutzen Abwägung letzten Endes auf lokales Training festgelegt. Dies ist ebenso hinsichtlich des Trainings des *YOLO* Objektdetektors besser, da das sehr spezifische *Darknet* Framework, das in der Implementierung genutzt wird, bisher von noch keinem Cloud Anbieter unterstützt wurde. Das Trainieren des *YOLO* Objektdetektors in der Cloud hätte demnach eine Umentscheidung auf eine Alternativ-Implementierung in beispielsweise *TensorFlow* oder *PyTorch* nötig gemacht. Werden noch andere Anpassungen in der Programmlogik mit einbezogen, z.B. dass ein Zugriff auf das Dateisystem beim Erstellen von Dateien in der *SSD* Implementierung in der Cloud Umgebung nicht möglich ist, so kommen zusätzlich zeitliche Bedenken mit auf. Ein lokales Trainieren bietet unter den genannten Voraussetzungen somit eine weitaus bessere Umgebung.

Auch wurden Überlegungen zum Training in *Google Colab* unternommen, da diese einfach ein Training mit TPUs ermöglichen. Diese Art von Spezialhardware erreicht pro TPU-Kern eine Rechenleistung von bis zu 92 TOPS [48]. Werden 2048 solcher TPU-Kerne zu einem TPU-Pod zusammen geschlossen, so ergibt sich eine Rechenleistung von über 100 PetaFLOPS [49]. Zudem ist die größere Rechenleistung gleichzeitig effizienter als herkömmliche GPUs (siehe Abbildung 3.5).

¹ Je nach Account kann eine unterschiedliche Anzahl an Projekten erstellt und Speicherplatz verwendet werden. Die *Beginner* Ausstattung von einem Projekt und 10 GB Speicher ist allerdings kostenfrei.

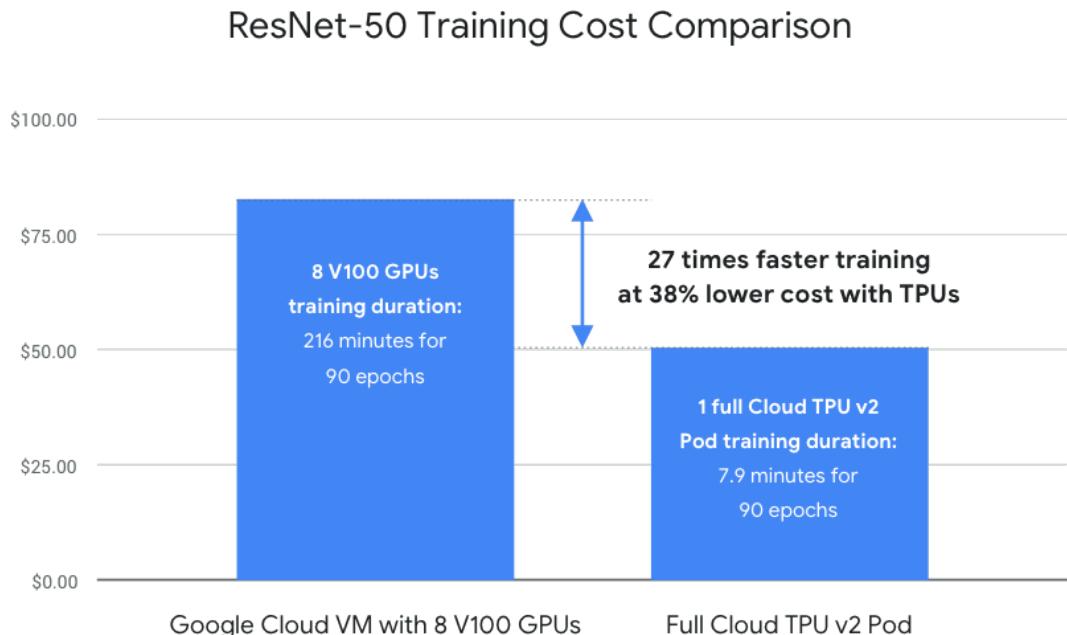


Abbildung 3.5.: Vergleich V100 - TPU Pod [32]

Da allerdings in *Google Colab* nach 90 Minuten ein Trainingsjob beendet wird und die Rechenressourcen neuen Nutzern zugewiesen werden, war diese Art des Trainings ebenfalls nicht möglich.

3.5. Auswahl einer Drohne

Bei der Auswahl der Drohne sind vor allem zwei Bereiche zu betrachten:

- Die gesetzlichen Rahmenbedingungen zur Drohne und
- Die technischen Anforderungen an die Drohne.

Was die gesetzlichen Anforderungen betrifft, so wurden im Juni 2019 von der *European Aviation Safety Agency* (EASA) einheitliche Regeln veröffentlicht, die den Drohnenbetrieb in der Europäischen Union einheitlich regeln sollen. Da den Mitgliedsstaaten ein Zeitraum von einem Jahr zur Umsetzung der Regularien zugesprochen wurde, gelten in Deutschland weiterhin die Vorschriften der deutschen Drohnen-Verordnung von 2017 [50].

Diese schreiben unter anderem vor [51]:

- Eine Kennzeichnungspflicht ab einem Startgewicht von über 250g,
- Eine maximale Flughöhe von 100 Metern über dem Grund,
- Eine Haftpflichtversicherung und
- Flugverbotszonen (Wohngrundstücke, Naturschutzgebiete, etc.).

Aufgrund dieser Auflagen wurde der Beschluss gefasst, dass die auszuwählende Drohne nur innerhalb von geschlossenen Wohnräumen im privaten Betrieb genutzt werden soll und zudem ein Startgewicht von unter 250g besitzen soll. Um eine studentische Machbarkeitsstudie durchführen zu können, ist eine solche eingeschränkte Anwendungsumgebung allemal ausreichend.

Die Programmierbarkeit der Drohne gehört zu der wichtigsten technischen Anforderung an die Drohne. Zudem soll sie eine integrierte Kamera aufweisen, die in der Lage ist, einen Videostream zur Laufzeit zur Verarbeitung bereit zu stellen. Auch die Akkulaufzeit und Robustheit der Drohne wurden als Entscheidungskriterien aufgenommen. Bezuglich der Kostenübernahme konnte zuvor eine Einigung mit Vertretern des Informatik Labors der DHBW Karlsruhe erzielt werden. Die Kosten werden vollständig übernommen, solange sie sich im Rahmen eines Studienarbeit angemessenen Budgets befinden.

Die sehr speziellen Anforderungen ließen nur ein Drohnen Modell auf dem Markt zu, die *Ryze Tello EDU* Drohne. Die bietet eine Kamera mit 720p Übertragungsqualität an, eine Akkulaufzeit von 13 Minuten und ein *Python Software Development Kit* (SDK) zur Programmierung. Auch wirbt sie mit der Fähigkeit von präzisem Schweben, was gerade für die Objektdetektion von Vorteil sein könnte [52].

3.6. Spezifikation der Inventursoftware

Die Inventursoftware besteht aus einer einfachen Client-Server Anwendung. Die Client-Anwendung zeigt das Live-Drohnenbild mit den eingezeichneten, erkannten Objekten. Zudem wird die Anzahl an erkannten Objekten der spezifischen Klassen rechts daneben dargestellt. Wann die Drohne die Inventur durchführen soll, wird auf Initiative des Benutzers gestartet.

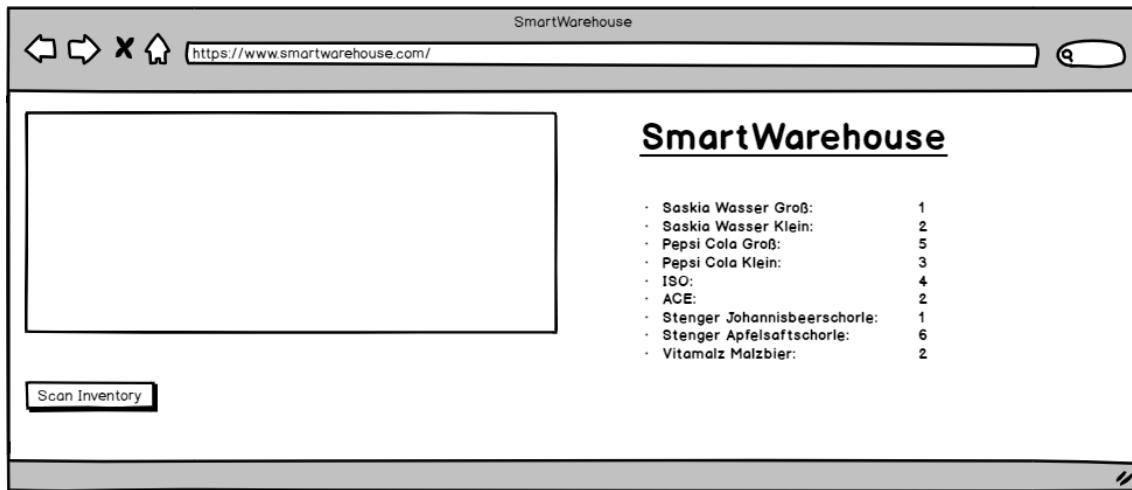


Abbildung 3.6.: SmartWarehouse User Interface

Der Server, geschrieben im *Django* Framework für Python, verbindet sich zu Drohne und weißt deren Flugsequenzen an. Auf ihm ist das *Deep Learning* Modell deployt. Der Server inferiert die von der Drohne empfangenen Video-Stream-Frames mit dem Modell und gibt diese anschließend an die Client Applikation weiter. Eine *REST* Schnittstelle gibt Aussage über die Bestandsdaten des Warenhauses.

4. Realisierung

Im folgenden Abschnitt soll auf die Implementierung des *Smart Warehouse* Szenarios eingegangen werden. Insbesondere werden Probleme während der Umsetzung der beiden Objektdetektoren *SSD* und *YOLO* betrachtet, Herausforderungen im Rahmen der Drohnen Anbindung besprochen und letztendlich die Realisierung der Dashboard-Webapplikation zur Durchführung der Inventur aufgezeigt.

4.1. Umsetzung der Objektdetektoren

SSD

Für den *SSD* wurde wie bereits erläutert nicht die ursprüngliche Referenzimplementierung im *Caffe* Framework verwendet, sondern eine Custom-Implementierung in *PyTorch*. Da der erstellte Datenbestand nur 1078 gelabelte Daten enthält, wurde zusätzlich zur Custom-Implementierung ein fünffaches Kreuzvalidierungsverfahren realisiert. Es dient dazu ein höheres Abstraktionsvermögen des Modells auf dem geringen Datenbestand zu erreichen. Auch unterstützte die Referenzimplementierung keine Validierung durch zuvor ungesehene Daten. Die Modellklassen zur Abstraktion des Datenbestandes wurden dahingehend angepasst.

Um ein lokales Training auf der *NVIDIA GTX 1080* GPU zu ermöglichen, wurde zudem *CUDA* Version 10.1 verwendet. Trainiert wurde mit folgenden Hyperparametern:

- Batch Größe: 16
- Lernrate: $1.0 \cdot e^{-3}$
- Momentum: 0.9
- Epochen: 500
- Gradientenverfahren: Stochastic Gradient Descent
- Kostenfunktion: Smooth L1

Das Basisnetzwerk des *SSDs* besteht aus einem auf *ImageNet* vortrainierten *VGG16*. Die restlichen *Convolutional Layer* sind *Xavier* initialisiert.

Die Hyperparameter sind nahezu gleich zu denen in der ursprünglichen wissenschaftlichen Veröffentlichung. Wesentlich die Batch Größe wurde für größere Stabilität von 32 auf 16 heruntergesetzt. Auch in der Evaluierung wurde die Batch Größe von 64 auf 48 herunter gesetzt, da die Eingangsdaten eine weitaus höhere Auflösung als die ursprünglich im *PascalVOC* verwendeten Daten haben. Andernfalls wird Gefahr gelaufen einen Speicherüberlauf zu erzielen.

Ursprünglich wurden 500 Epochen für das Training vorgesehen für jeden der Kreuzvalidierungsschritte. Da allerdings beim Training schon nach knapp über hundert Epochen sich der Gradient der Kostenfunktion nur träge veränderte, wurde im Sinne des *Early Stoppings* nach 108 Epochen das Training vorzeitig beendet, um *Overfitting* zu vermeiden. Das niedrigste Ergebnis der Kostenfunktion betrug 1.726. Es ergab eine *mAP* von 78.2%, leicht über den Referenzergebnissen von *PascalVOC*.

Wird nun das trainierte Modell auf echte Daten angewendet, so fällt auf, dass manche Objekte doppelt detektiert werden. Um dieses Problem zu lösen, gibt es zwei Möglichkeiten.

Als erstes kann bei der Detektion der minimale *confidence score* angegeben werden, ab wann eine Detektion offiziell als solche wahrgenommen wird. Hier liegt die Herausforderung darin, einen optimalen Wert zu finden, sodass verdeckte Objekte noch als solche erkannt werden, aber doppelt erkannte Objekte nicht mehr auftreten. Der *confidence score* wurde nach mehrmaligem Iterieren auf 0.75 festgelegt.

Die zweite Möglichkeit besteht darin, die maximale Überlappung zweier Bounding Boxen festzulegen. Somit werden doppelte Bounding Boxen, die sich flächenmäßig über einem gewissen Grenzwert überlappen, auf eine Bounding Box reduziert. Er stellte ich sich Parameter von 0.5 als geeignet heraus.

Außerdem wurden Inkonsistenzen im Detektionsverhalten festgestellt:

Merke:

- Bilder mit Labels in der Arbeit
- Soll beweisen, dass etwas entstanden ist (Krassen Eindruck vermitteln)



Abbildung 4.1.: Detektion mehrerer Wasserflaschen

So ist es schwierig, teilweise verdeckte Objekte zu detektieren (siehe Abbildung 4.1). Dies mag aber auch daran liegen, dass im Trainingsdatensatz nicht genügend Daten vorhanden waren, die solche Fälle abdecken.

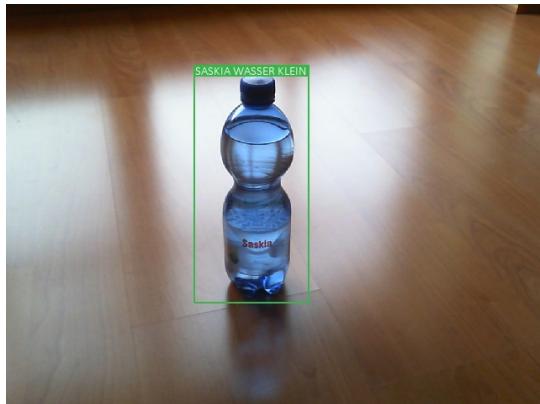


Abbildung 4.2.: Detektion einer nahen Wasserflasche

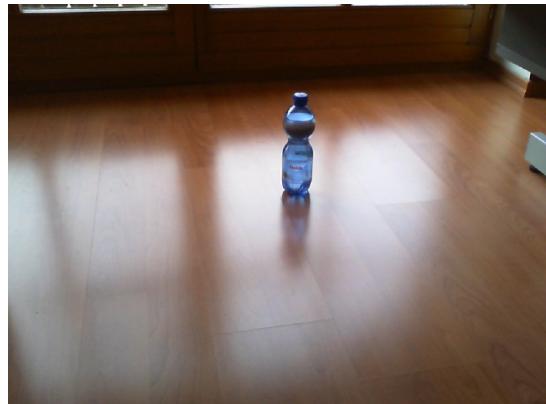


Abbildung 4.3.: Detektion einer entfernten Wasserflasche

- Keine Ergebnisse evaluieren, keine Detektoren bewerten

YOLO

Robin

4.2. Drohnen Anbindung

Drohnen Interface

Robin

Modellinferenz

Der Video Stream der Drohne kann mittels der *openCV* Klasse *VideoCapture* über das UDP Protokoll angesprochen werden. Bei der Inferenz fällt allerdings entgegen der Erwartungen auf, dass die Inferenz überdurchschnittlich langsam verläuft. Das Problem lässt sich auf die synchrone Arbeitsweise des bisherigen Detektionsalgorithmus zurück führen, bei dem erst ein neuer Frame des Videostreams angefragt wird, sobald das aktuelle Bild durch die Vorverarbeitung gelaufen und durch das Modell inferiert wurde.

Um dem entgegen zu wirken, wurde ein Bufferkonzept in einem parallelen Thread realisiert, der einzelne Frames zeitgleich zur Inferenz anfragt und zwischenspeichert. Ist der Buffer voll, so wird die Frameanfrage ausgesetzt. Dadurch konnte die FPS Anzahl von maximal 18 auf die vollen 30 gesteigert werden.

4.3. Dashboard Entwicklung

Der Server wurde mit dem *Flask* Framework in Python implementiert. Er führt den Inferenzalgorithmus des *SSDs* bzw. des *YOLO* Objektdetektors bei jeder Anfrage an eine vordefinierte Route aus und sendet das inferierte Bild mit den Bounding Boxen zurück an den Client. Der Client wurde mit dem *Bootstrap* Framework designed.

4.4. Zählalgorithmus

5. Ergebnisse

- Ergebnisse auswerten

6. Bewertung

- Ergebnisse bewerten und diskutieren
- Eignet sich YOLO und SSD für die Industrie?
- Ist das SmartWarehouse Szenario umsetzbar?

7. Zusammenfassung und Ausblick

- Klare Darstellung, was die Arbeit geliefert hat
- Was liefert die Arbeit, was bisher noch nicht bekannt war (Mehrwert, auch wenn etwas nicht geht)?
- Ca. 2-4 Punkte: Zukünftige Ziele

Literaturverzeichnis

- [1] Aurélien Géron. *Machine Learning mit Scikit-Learn & TensorFlow: Konzepte, Tools und Techniken für intelligente Systeme: Übersetzung von Kristian Rother.* 1. Aufl. Heidelberg: dpunkt.verlag GmbH, 2018.
- [2] MathWorks. *Deep Learning: Drei Dinge, die Sie wissen sollten.* MathWorks, 2019. URL: <https://de.mathworks.com/discovery/deep-learning.html> (Einsichtnahme: 12.10.2019).
- [3] dokt. innovation. *inventAIRyX.* dokt. innovation Homepage, 2019. URL: <https://www.doks-innovation.com/inventory-x/> (Einsichtnahme: 26.10.2019).
- [4] Ravindra Parmar. *Detection and Segmentation through ConvNets.* Towards Data Science, 2018. URL: <https://towardsdatascience.com/detection-and-segmentation-through-convnets-47aa42de27ea> (Einsichtnahme: 07.03.2020).
- [5] Priya Dwivedi. *Semantic Segmentation — Popular Architectures.* Towards Data Science, 2019. URL: <https://towardsdatascience.com/semantic-segmentation-popular-architectures-dff0a75f39d0> (Einsichtnahme: 07.03.2020).
- [6] Philippe Lucidarme. *Simplest perceptron update rules demonstration.* Homepage Blog, 2017. URL: <https://www.lucidarme.me/simplest-perceptron-update-rules-demonstration/> (Einsichtnahme: 26.10.2019).
- [7] Wikipedia. *Deep Learning.* Wikipedia, 2019. URL: https://de.wikipedia.org/wiki/Deep_Learning (Einsichtnahme: 27.01.2019).
- [8] David E. Rumelhart / Geoffrey E. Hinton / Ronald J. Williams. *Learning Internal Representations by Error Propagation.* Hrsg. von University of California, San Diego. 09/1985. URL: <https://apps.dtic.mil/dtic/tr/fulltext/u2/a164453.pdf> (Einsichtnahme: 26.10.2019).
- [9] Xavier Glorot, Y. B. „Understanding the difficulty of training deep feedforward neural networks“. Diss. Montréal: Universite de Montréal, 2010. URL: <http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf> (Einsichtnahme: 26.10.2019).

- [10] Imad Dabbura. *Gradient Descent Algorithm and Its Variants*. Towards Data Science, 2017. URL: <https://towardsdatascience.com/gradient-descent-algorithm-and-its-variants-10f652806a3> (Einsichtnahme: 26. 10. 2019).
- [11] Danqing Liu. *A Practical Guide to ReLU*. Medium, 2017. URL: <https://medium.com/@danqing/a-practical-guide-to-relu-b83ca804f1f7> (Einsichtnahme: 26. 10. 2019).
- [12] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Chang-Yang Fu, Alexander C. Berg. „SSD: Single Shot MultiBox Detector“. Diss. 2016-12-29. URL: <https://arxiv.org/pdf/1512.02325.pdf> (Einsichtnahme: 02. 11. 2019).
- [13] Mark Everingham, J. W. *The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Development Kit*. 2007. URL: <http://host.robots.ox.ac.uk/pascal/VOC/voc2007/html/doc/voc.html#SECTION00030000000000000000> (Einsichtnahme: 08. 02. 2020).
- [14] Mark Everingham, J. W. *The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Development Kit*. 2012. URL: http://host.robots.ox.ac.uk/pascal/VOC/voc2012/html/doc/devkit_doc.html#SECTION00030000000000000000 (Einsichtnahme: 08. 02. 2020).
- [15] Daphne Cornelisse. *An intuitive guide to Convolutional Neural Networks*. freeCodeCamp Homepage, 2018. URL: <https://medium.freecodecamp.org/an-intuitive-guide-to-convolutional-neural-networks-260c2de0a050> (Einsichtnahme: 26. 10. 2019).
- [16] Abhineet Saxena. *Convolutional Neural Networks (CNNs): An Illustrated Explanation*. XRDS Crossroads - The ACM Magazine for Students, 2016. URL: <https://blog.xrds.acm.org/2016/06/convolutional-neural-networks-cnns-illustrated-explanation/> (Einsichtnahme: 26. 10. 2019).
- [17] Leonadro Araujo Santos. *Pooling Layer: Introduction*. GitBook, 2018. URL: https://leonardoaraujosantos.gitbooks.io/artificial-intelligence/content/pooling_layer.html (Einsichtnahme: 26. 10. 2019).
- [18] Jonathan Hui. *mAP (mean Average Precision) for Object Detection*. Medium, 2018. URL: https://medium.com/@jonathan_hui/map-mean-average-precision-for-object-detection-45c121a31173 (Einsichtnahme: 15. 02. 2020).

- [19] Tarang Shah. *Measuring Object Detection models - mAP - What is Mean Average Precision?* Tarang Shabs Blog, 2018. URL: <https://tarangshah.com/blog/2018-01-27/what-is-map-understanding-the-statistic-of-choice-for-comparing-object-detection-models/> (Einsichtnahme: 15.02.2020).
- [20] Adrian Rosebrock. *Intersection over Union (IoU) for object detection.* pyimagesearch, 2016. URL: <https://www.pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection/> (Einsichtnahme: 02.11.2019).
- [21] Rohith Gandhi. *R-CNN, Fast R-CNN, Faster R-CNN, YOLO — Object Detection Algorithms.* Towards Data Science, 2018. URL: <https://towardsdatascience.com/r-cnn-fast-r-cnn-faster-r-cnn-yolo-object-detection-algorithms-36d53571365e> (Einsichtnahme: 08.03.2020).
- [22] Umer Farooq. *From R-CNN to Mask R-CNN.* Medium, 2018. URL: https://medium.com/@umerfarooq_26378/from-r-cnn-to-mask-r-cnn-d6367b196cf (Einsichtnahme: 08.03.2020).
- [23] Dhruv Parthasarathy. *A Brief History of CNNs in Image Segmentation: From R-CNN to Mask R-CNN.* Medium, 2017. URL: <https://blog.athelas.com/a-brief-history-of-cnns-in-image-segmentation-from-r-cnn-to-mask-r-cnn-34ea83205de4> (Einsichtnahme: 08.03.2020).
- [24] Andrew Ng. *Anchor Boxes.* Coursera, 2019. URL: <https://www.coursera.org/lecture/convolutional-neural-networks/anchor-boxes-yNwO0> (Einsichtnahme: 02.11.2019).
- [25] MathWorks. *vgg16.* MathWorks Homepage, 2019. URL: <https://de.mathworks.com/help/deeplearning/ref/vgg16.html;jsessionid=bf0fea41a0c7700184672711881f> (Einsichtnahme: 02.11.2019).
- [26] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, Ali Farhadi. „You Only Look Once: Unified, Real-Time Object Detection“. Diss. 2019. URL: <https://arxiv.org/pdf/1506.02640.pdf> (Einsichtnahme: 10.11.2019).
- [27] Joseph Redmon, A. F. „YOLOv3: An Incremental Improvement“. In: (2018). URL: <https://pjreddie.com/media/files/papers/YOLOv3.pdf> (Einsichtnahme: 04.02.2020).

- [28] Renu Khandelwal. *COCO and Pascal VOC data format for Object detection*. Towards Data Science, 2019. URL: <https://towardsdatascience.com/coco-data-format-for-object-detection-a4c5eaf518c5> (Einsichtnahme: 02.02.2020).
- [29] Arun Ponnusamy. *Preparing Custom Dataset for Training YOLO Object Detector*. Arun Ponnusamy Homepage, 2019. URL: <https://www.arunponnusamy.com/preparing-custom-dataset-for-training-yolo-object-detector.html> (Einsichtnahme: 02.02.2020).
- [30] NVIDIA. *TRAIN MODELS FASTER*. NVIDIA Homepage, 2020. URL: <https://developer.nvidia.com/cuda-zone> (Einsichtnahme: 09.02.2020).
- [31] PyTorch. *GET STARTED*. PyTorch Homepage, 2020. URL: <https://pytorch.org/get-started/locally/> (Einsichtnahme: 09.02.2020).
- [32] Google Cloud. *Cloud TPU*. Google Cloud Homepage, 2020. URL: <https://cloud.google.com/tpu/?hl=de> (Einsichtnahme: 09.02.2020).
- [33] Karl Freund. *Microsoft: FPGA Wins Versus Google TPUs For AI*. Forbes, 2017. URL: <https://www.forbes.com/sites/moorinsights/2017/08/28/microsoft-fpga-wins-versus-google-tpus-for-ai/#781e2bf53904> (Einsichtnahme: 09.02.2020).
- [34] Amazon Web Services. *Amazon SageMaker*. AWS Homepage, 2020. URL: <https://aws.amazon.com/de/sagemaker/> (Einsichtnahme: 12.03.2020).
- [35] Amazon Web Services. *Amazon SageMaker-ML-Instance-Typen*. AWS Homepage, 2020. URL: <https://aws.amazon.com/de/sagemaker/pricing/instance-types/> (Einsichtnahme: 14.03.2020).
- [36] Google Cloud Platform. *Produkte für künstliche Intelligenz und maschinelles Lernen: AI Platform*. GCP Homepage, 2020. URL: <https://cloud.google.com/products/ai> (Einsichtnahme: 14.03.2020).
- [37] Google Colaboratory. *Welcome to Colaboratory!* Google Colaboratory Homepage, 2020. URL: <https://colab.research.google.com/notebooks/welcome.ipynb> (Einsichtnahme: 14.03.2020).
- [38] Google Cloud. *Colaboratory-Notebooks*. Google Cloud Dokumentation, 2020. URL: <https://cloud.google.com/automl-tables/docs/notebooks> (Einsichtnahme: 14.03.2020).

- [39] Microsoft Azure. *Azure Machine: Machine-Learning-Dienst für Unternehmen zur schnelleren Erstellung und Bereitstellung von Modellen*. Microsoft Azure Homepage, 2020. URL: <https://azure.microsoft.com/de-de/services/machine-learning/> (Einsichtnahme: 14. 03. 2020).
- [40] FloydHub. *Home*. FloydHub Documentation, 2020. URL: <https://docs.floydhub.com/> (Einsichtnahme: 15. 02. 2020).
- [41] Intan Purnamasar. *Vergleich von Algorithmen zur Objekterkennung für die Anwendung in Abbildungen komplexer Energiesystem*. RWTH Aachen University, 2018. URL: https://www.matse.itc.rwth-aachen.de/dienste/public/show_document.php?id=18753 (Einsichtnahme: 14. 03. 2020).
- [42] Amazon Web Services. *Kostenloses Kontingent für AWS*. AWS Homepage, 2020. URL: <https://aws.amazon.com/de/free/?all-free-tier.sort-by=item.additionalFields.SortRank&all-free-tier.sort-order=asc> (Einsichtnahme: 14. 03. 2020).
- [43] Google Cloud Platform. *Kostenlose Stufe der Google Cloud Platform*. GCP Homepage, 2020. URL: <https://cloud.google.com/free> (Einsichtnahme: 14. 03. 2020).
- [44] Microsoft Azure. *Create your Azure free account today: Get started with 12 months of free services*. Microsoft Azure Homepage, 2020. URL: <https://azure.microsoft.com/en-us/free/> (Einsichtnahme: 14. 03. 2020).
- [45] Microsoft Azure. *Virtuelle Windows-Computer – Preise*. Microsoft Azure Homepage, 2020. URL: <https://azure.microsoft.com/de-de/pricing/details/virtual-machines/windows/> (Einsichtnahme: 12. 03. 2020).
- [46] FloydHub. *Plans*. FloydHub Documentation, 2020. URL: <https://docs.floydhub.com/faqs/plans/#what-is-in-the-trial-plan> (Einsichtnahme: 25. 01. 2020).
- [47] TechPowerUp. *GPU Specs Database*. TechPowerUp Homepage, 2020. URL: <https://www.techpowerup.com/gpu-specs/> (Einsichtnahme: 09. 02. 2020).
- [48] Harald Bögeholz. *Künstliche Intelligenz: Architektur und Performance von Googles KI-Chip TPU*. Heise Online, 2017. URL: <https://www.heise.de/newsticker/meldung/Kuenstliche-Intelligenz-Architektur-und-Performance-von-Googles-KI-Chip-TPU-3676312.html> (Einsichtnahme: 09. 02. 2020).
- [49] Google Cloud. *Systemarchitektur*. Google Cloud Dokumentation, 2020. URL: <https://cloud.google.com/tpu/docs/system-architecture> (Einsichtnahme: 09. 02. 2020).

- [50] EASA. *EU wide rules on drones published*. EASA Homepage, 2019. URL: <https://www.easa.europa.eu/newsroom-and-events/press-releases/eu-wide-rules-drones-published> (Einsichtnahme: 28.03.2020).
- [51] Drohnen.de. *Vorschriften, Genehmigungen für die Nutzung von Drohnen und Multicoptern*. Drohnen.de Homepage, 2020. URL: <https://www.drohnen.de/vorschriften-genehmigungen-fuer-die-nutzung-von-drohnen-und-multicoptern/> (Einsichtnahme: 28.03.2020).
- [52] RyzeRobotics. *Tello EDU*. RyzeRobotics Homepage, 2020. URL: <https://www.ryzerobotics.com/de/tello-edu?from=store-product-page> (Einsichtnahme: 28.03.2020).

A. Anhang

Abbildungen