



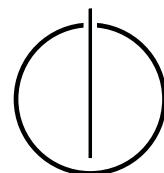
DEPARTMENT OF INFORMATICS

TECHNICAL UNIVERSITY OF MUNICH

Bachelor's Thesis in Informatics

# **Implementing a mobile app for object detection**

David Drews







DEPARTMENT OF INFORMATICS  
TECHNICAL UNIVERSITY OF MUNICH

Bachelor's Thesis in Informatics

**Implementing a mobile app for object detection**

**Entwicklung einer mobilen App zur  
Objekterkennung**

Author: David Drews  
Supervisor: Univ.-Prof. Dr. Hans-Joachim Bungartz  
Advisor: Severin Reiz, M.Sc.  
Submission Date: 15th of August 2021





I confirm that this bachelor's thesis is my own work and I have documented all sources and material used.

Munich, 15th of August 2021

David Drews



---

## **Abstract**

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.



# Contents

<b>Abstract</b>	<b>vii</b>
<b>1. Motivation</b>	<b>1</b>
1.1. Growing Support for Running Machine Learning Operations on Mobile Platforms	1
1.2. Offline Usability	1
1.3. Improved Privacy	2
<b>2. Background Theory</b>	<b>3</b>
2.1. Important Concepts	3
2.1.1. Machine Learning	3
2.1.2. Artificial Neural Networks	5
2.1.3. Deep Learning	5
<b>3. Object Detection</b>	<b>6</b>
3.1. How Object Detection Differs From Related Tasks	6
3.1.1. Semantic Segmentation	6
3.1.2. Image Classification	6
3.1.3. Object Detection	7
3.1.4. Instance Segmentation	7
3.2. Object Detection Frameworks	7
3.2.1. Backbones	7
3.2.2. Two-Stage Detectors	8
3.2.3. R-CNN	9
3.2.4. Single-Stage Detectors	9
3.3. SSD MobileNet v1	9
3.3.1. MobileNet Architecture	10
3.3.2. Introduction to SSD	10
3.3.3. The SSD Head	10
3.3.4. Architecture	11
3.3.5. Latest Improvements	11
<b>4. App Development</b>	<b>12</b>
4.1. Previous State of the Application	12
4.2. Development Goals	12
4.2.1. Migration From Java to Kotlin	12
4.2.2. New Functionality: Object Detection	13
4.3. Implementing Object Detection Based on The TensorFlow Lite Framework	13
4.3.1. Some Deep	14
4.3.2. Dive Into	14

4.3.3. Object Detection . . . . .	14
4.3.4. Implementation Fun . . . . .	14
4.3.5. Next steps in the development of TUM-Lens . . . . .	14
<b>5. Results</b>	<b>16</b>
5.1. Performance Evaluation . . . . .	16
5.2. Accuracy . . . . .	16
5.3. Possible Applications . . . . .	16
<b>A. Screenshots of the Application</b>	<b>17</b>
<b>B. Tips With Greetings From the Chair</b>	<b>18</b>
B.1. Tips . . . . .	18
B.1.1. How to Describe . . . . .	18
B.1.2. How to Quote . . . . .	18
B.1.3. How to Math . . . . .	18
B.2. Environments . . . . .	19
B.2.1. How to Figure . . . . .	19
B.2.2. How to Algorithm . . . . .	19
B.2.3. How to Code . . . . .	21
B.2.4. How to Table . . . . .	21
<b>Bibliography</b>	<b>24</b>
<b>Acronyms</b>	<b>28</b>
<b>Glossary</b>	<b>29</b>

# 1. Motivation

The aim of this work was to further develop the Android app TUM-Lens [18]. The core functions of the app include the analysis of images that are captured via the camera of the Android device and transmitted to the app as a live feed. For an optimal user experience, the analysis of the images must take place in near real time. This is the only way to ensure that the analysis results displayed always match the current content of the camera feed, which can change very quickly due to panning of the camera by its user. While in many applications the analysis of image data can take place decentrally in powerful data centres, in the case of TUM-Lens the image analysis runs on the mobile device itself. With the completion of this work, image analysis now also includes object detection in addition to the classification of images.

## 1.1. Growing Support for Running Machine Learning Operations on Mobile Platforms

Support for the development of Machine Learning (ML) and also in particular deep learning applications for smartphones is growing steadily and from different directions at the same time. Developer-friendly frameworks such as TensorFlow, developed by Google Brain, or PyTorch, developed by Facebook's AI Research Lab, are among the best-known deep learning frameworks [12]. The release of TensorFlow Lite<sup>1</sup> 2017 [42] and PyTorch Mobile 2019 [31] show that mobile platforms increasingly come into focus of companies providing Machine Learning software. In recent year, device manufacturers and operating system developers also started to provide dedicated hardware and software components for mobile machine learning. Examples include Apple's Neural Engine [43], unveiled in 2017, or Android's Neural Networks API (NNAPI) [6]. Apple's Neural Engine is a hardware component optimised for Machine Learning requirements. Android's NNAPI, on the other hand, is an Android C application programming interface (API) for efficient computation of ML operations and provides a basic set of functions for higher-level ML frameworks. As a result of these developments, it is becoming easier for developers to build ML applications that run efficiently on mobile devices. This support was a major catalyst for the initial and further development of TUM-Lens in the context of two bachelor theses.

## 1.2. Offline Usability

TUM-Lens is more independent compared to many other Machine Learning based apps as it does not require an internet connection to use it. Often, apps and services by definition need a connection to the internet to perform their task. The Amazon voice assistant Alexa

---

<sup>1</sup><https://www.tensorflow.org/lite>

can answer simple voice commands to control smart home devices or check the time without an internet connection and thus already uses on-device Machine Learning. But even if Alexa could analyse and understand all voice commands locally, the request would still have to be forwarded to the Amazon servers in most cases. Due to the large number of possible queries, not all answers can be kept on the device, but must be retrieved from the Internet. Such queries include daily topics such as the weather report, traffic or the result of a sporting event. However, an internet connection is not required to use the full range of functions of TUM-Lens. All the information needed for image classification and object detection is stored locally on the device in the form of various already trained Artificial Neural Network (ANN). With the integration of the corresponding mobile frameworks, the image analysis can therefore be carried out locally on the device, making the app independent of an internet connection.

### 1.3. Improved Privacy

The use of on-device ML provides another mechanism for protecting personal data in the context of machine learning in addition to existing methods such as differential privacy. Due to the growing support for mobile ML applications mentioned above, but also due to the independently increasing power of mobile devices [13], not only the use of pre-trained ANNs becomes possible, but also the training of new ANNs on the mobile device itself becomes more and more relevant [26]. If the training process takes place locally on the device itself, no data needs to be transferred to external instances such as a company's servers. This makes it possible to develop applications that adapt more and more individually to the user as they are used, while guaranteeing maximum data protection. An example of the development of such an application is DeepType [45]. DeepType attempts to predict the next word used when the user enters the keyboard. While every user initially starts with the same pre-trained version of the ANN used by DeepType, the application continues to train this ANN with each input and thus adapts more and more to the characteristic input behaviour of the user without the text inputs ever leaving the device.

## 2. Background Theory

### 2.1. Important Concepts

Everybody is talking about Machine Learning (ML). It is already impossible to imagine our everyday life without the use of the term. Due to the multitude of contexts in which Machine Learning (ML) is spoken of, some justifiably and some unjustifiably, it is important to create a common understanding for the theoretical content of this work.

#### 2.1.1. Machine Learning

A popular definition of ML is attributed to Arthur Samuel describing it as the "field of study that gives computers the ability to learn without being explicitly programmed"<sup>1</sup>. Machine Learning algorithms circumvent this need for explicit programming by improving an internal model through data. This process is called training and the data used to train the model is often regarded to as the model's experience [28]. As depicted in figure 2.1, ML can be divided into the subfields supervised learning, unsupervised learning, semi-supervised learning and reinforcement learning.

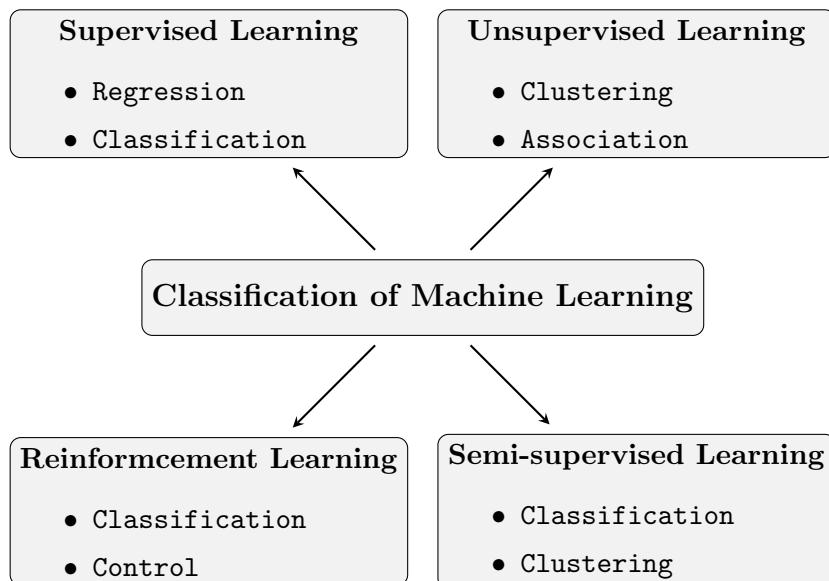


Figure 2.1.: The field of Machine Learning divided into subfields by the characteristics of the underlying learning process. Also indicates the learning problems that are typically tried to be solved by applying the respective learning process.

<sup>1</sup>Although cited in popular machine learning material like Andrew Ng's ML course at Stanford [29] the quote appears neither in Samuel's 1959 [34] nor his 1967 paper [35].

### **Supervised Learning**

In supervised learning, the learning machine is provided with input data as well as the output that is expected for the given input [22]. In the classical case of spam filtering, the input can be a collection of emails and the expected output is a label attached to each email that either classifies it as spam or as non-spam. The learning machine is then fed all e-mails as input data and learns to recognise which information in the input is important to produce the correct classification. As the system knows the correct answer for each training input, it can process an email, predict whether or not it is spam, and then use the known answer to change its weights in a way that will make it more likely to lead to a correct prediction and less likely to lead to a false prediction the next time it is presented with similar input.

### **Unsupervised Learning**

Detecting hidden patterns and structuring data is where unsupervised learning comes into play. Learners of this type don't need to be provided with an expected output while being trained [36]. A scenario for the application of unsupervised learning is the problem of dividing a customer base into subgroups in order to treat every subgroup according to their specific needs. An employee might help the machine learning system by providing the number of subgroups she wants the system to generate. The learner then builds up its representation of the internal structure of the entire data set with every input it processes. After having processed enough customers, it will most likely have identified the key metrics that distinguish customers into the different groups.

### **Semi-Supervised Learning**

One use for semi-supervised learning is cluster analysis, which was already used as an example in the previous paragraph. In the case of semi-supervised learning, the system no longer has to work out the different groups (also known as *clusters*) from the unlabelled data alone. Instead, it can use a small set of already labelled customers as a reference and build its internal representation of the entire dataset (labelled and unlabelled) around the clusters indicated by the pre-labelled data. This is especially useful because in many domains collecting or creating labelled data is difficult, expensive, or both [47].

### **Reinforcement Learning**

Reinforcement learning is "learning what to do - how to map situations to actions - so as to maximize a numerical reward signal" [40]. Systems are trained via reinforcement learning to learn how to behave in dynamic environments. The tasks in these environments can stretch from playing a video game [44] to driving an autonomous car [33]. These exemplary tasks show two characteristics that distinguish reinforcement learning from the other subfields of ML: The reward signal is often delayed and attribution to single actions is difficult. Only once a game is won or the car has arrived safely at its destination the system knows if all the decisions it made along the way lead to a positive outcome. *Trial-and-error* is therefore a term that summarises this learning paradigm quite precisely.

### 2.1.2. Artificial Neural Networks

An Artificial Neural Network - often just referred to as neural network - is a data processing concept that is inspired by biological neurons and their interconnectivity. As figures 2.2 and 2.3 show the artificial neurons (also called *nodes*) in an ANN are grouped in *layers*. There are three important types of layers: The *input layer*<sup>2</sup>, the *output layer* and an arbitrary number of *hidden layers* in between the input and output layer. Similar to neurons in human brains, nodes of different layers can be connected. In ANNs, the nodes exchange signals in the form of numbers. Each node outputs a number that is computed by applying a non-linear function to its inputs. The output signal can then be a new input for other nodes or it can be part of the result returned by the output layer. The connections between nodes are also known as *edges* and typically carry a weight. In the case of ANNs, the training process that is typical for all machine learning systems is the adjustment of these connection weights. The weights and other variables of the ANN are grouped under the term *parameters*. In summary, an ANN transforms an input vector into an output vector through a series of non-linear functions, where both the calculation of the output and the training process are characterised by the specific structure of the ANN and its parameters.

### 2.1.3. Deep Learning

Deep learning is a subarea of machine learning. Deep learning is characterised by the use of ANNs with many hidden layers. The more hidden layers a network has, the deeper it is. The deeper a network is and the more nodes the network has per layer, the more complex the computations that the ANN can successfully perform [10]. As the number of layers and nodes grows, so does the number of parameters. Their large number is the reason deep learning requires extensive amounts of data to provide adequate results compared to other sub-disciplines of machine learning. Networks of this genus have been given the ability to perform extraordinarily complex computations at the expense of a resource-intensive training process.

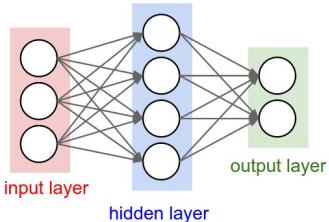


Figure 2.2.: 2-layered ANN. It is called fully connected as every node from the previous layer is connected to every node in the next layer.  
Source: [23]

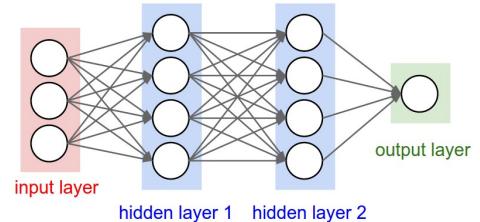


Figure 2.3.: 3-layered ANN. In ANNs, nodes in one layer are connected to nodes in other layers but not to other nodes in the same layer.  
Source: [23]

<sup>2</sup>Note that the input layer is not counted towards the total number of layers in an ANN.

# 3. Object Detection

## 3.1. How Object Detection Differs From Related Tasks

The field of Computer Vision (CV) encompasses numerous distinct problems and an even larger number of potential solutions. In the following, object detection as a typical task in the CV context is distinguished from other CV tasks that are closest to it in terms of learning objectives.

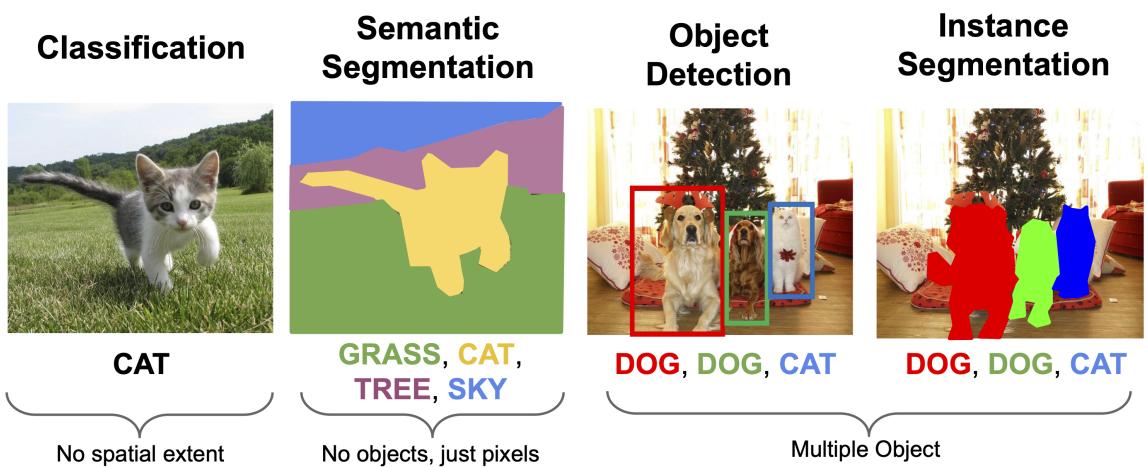


Figure 3.1.: Object detection differs conceptually from other related CV tasks regarding spatial information, the concept of objects and the number of detections in a scene.  
Source: [24]

### 3.1.1. Semantic Segmentation

In semantic segmentation, each pixel of an image is assigned a class. However, there are no objects. This means that if there are several objects of the same class in the image, all the associated pixels receive the same class label. Therefore, the different objects cannot be differentiated based on the result of the semantic segmentation.

### 3.1.2. Image Classification

In image classification, the result of the detection is a single class. In rare application variants, a bounding box for one object of the detected class is also returned - as a rule, however, no spatial extent is associated with the task of classification.

### 3.1.3. Object Detection

Object detection deals with the identification of any number of objects within an image. For each object, a class label and its position are returned as the coordinates of a rectangle enclosing the object. It is important here that, as depicted in Figure 3.1, several objects of the same class can also be recognised. In contrast to the task of semantic segmentation, the different objects of the same class can be distinguished.

### 3.1.4. Instance Segmentation

The instance segmentation essentially fulfils a better variant of the object detection. Again, several objects of different classes are detected and the positions of the classes are also part of the output. However, the positions are not marked by bounding boxes as in object detection, but each pixel belonging to an object receives a label of this object. The objects are therefore even more sharply separated from the image areas that do not contain an object and, in contrast to semantic segmentation, individual objects of the same class remain distinguishable.

## 3.2. Object Detection Frameworks

State of the art object detection frameworks run predominantly on deep learning architectures [2]. The bachelor's thesis preceding this work, *Implementing a TensorFlow-Slim based Android app for image classification* [19], already explains how Convolutional Neural Networks (CNNs) work and why they are the most important building blocks for deep learning frameworks solving perceptual tasks. This explanation will therefore not be pursued here once more. Instead, we survey a range of modern object detection frameworks. Each framework is defined by the combination of a backbone and a detector. The selection was not put together at random but is based on popularity and performance in standard object detection benchmarks [4].

### 3.2.1. Backbones

The term *backbone* in the context of ANNs might be used differently depending on the task pursued. In the case of visual tasks and this thesis, the backbone is the part of an object detection framework that extracts features from the input data. This encapsulation of the feature extraction task in its own set of CNNs allows the designer of the object detection pipeline to swap and test different backbones for the task at hand [4].

### 3. Object Detection

---

backbone	first publication	detectors
AlexNet	2012 [21]	HyperNet [20]
VGG-16	2014 [38]	PFPNet-R512 [K]
GoogLeNet	2014 [41]	YOLOv1 [32]
ResNets	2015 [14]	BlitzNet512 [8], CoupleNet [48], RetinaNet [25], Faster R-CNN [A], R-FCN ?
Inception-ResNet-V2	tbd	Faster R-CNN G-RMI [23], Faster R-CNN with TDM [24]
DarkNet-19	tbd	YOLOv2 [J]
MobileNet	2017 [17]	SSDv2

Table 3.1.: Overview of modern CNN backbone networks and detectors that build upon them.

TO-DO: further sources:

<https://scholar.google.com/scholar?q=DaiJ>: <http://arxiv.org/abs/1612.08242>

#### **AlexNet**

#### **ResNet**

While the baseline ANN described in section 2.1.2 only connects adjacent layers, Residual Neural Networks (ResNets) are a special subclass of ANN that introduces *shortcut connections* [14]. Shortcut connection are edges in the network that skip one or more layers.

#### **MobileNet**

As the name suggests, MobileNets were developed especially for mobile or embedded computer vision applications. The key innovation of the MobileNet architecture is the combination of two elements:

1. utilisation of depth-wise separable filters [37]
2. a network structure that enables the developer to scale the model size down by adjusting only two hyper-parameters

As a MobileNet is a building block of the detection framework implemented in TUM-Lens, a more detailed description can be found in section 3.3.1.

#### **3.2.2. Two-Stage Detectors**

This class of detectors separates the object detection task into two distinct stages [39]. The network of the first stage (named Region Proposal Network (RPN) in the context of the R-CNN detector family) uses the image data to generate region proposals. The second stage is a separate network that takes these region proposals (or ROI - short for regions of interest), potentially decreases the final number of ROI using mechanics that depend

on the specific detector, and then performs the classification on each of the final regions. The classified regions are then returned as the detected objects. Two-stage detectors tend to have a higher localization and object recognition accuracy than single-stage detectors but can only achieve that at the cost of considerably slower inference speed [15]. Table 3.1 shows a selection of popular detectors including R-CNN and some of its many successors, Fast-CNN and Faster-CNN.

### 3.2.3. R-CNN

R-CNN was Fast, Faster R-CNN, and other advancements like Mask-R-CNN 2014 [9]

### 3.2.4. Single-Stage Detectors

Detection frameworks are considered to have only a single stage when they consist of one deep neural network only and compute the objects (bounding box coordinates and category) in a single pass through that network. By eliminating the explicit generation of region proposals, single-stage detectors outperform their two-stage competitors with respect to speed while sacrificing a bit of detection accuracy, if at all [27]. This speed improvement makes single-stage detectors the preferred choice for applications running on mobile or embedded devices or in applications that require real-time image detection. Since the object detector implemented in TUM-Lens is described thoroughly in section 3.3, the following overview of selected single-state detectors will help us to understand how the app's object detector differs from other possible options.

#### RetinaNet

2018

#### YOLO

YOLO (You Only Look Once) ; [3]

## 3.3. SSD MobileNet v1

The SSD MobileNet v1 is the TensorFlow-Lite model used for the new object detection functionality integrated into TUM Lens. The model is pre-trained on the COCO dataset<sup>1</sup> and available on TensorFlow Hub<sup>2</sup>. Following the data flow through the object detection framework, we will first have a closer look at the architecture of the MobileNet used as the backbone network. Secondly, we introduce the core concepts of the Single Shot MultiBox Detector (SSD) followed by an detailed description of the SSD architecture. Finally, we conclude with an outlook on the latest improvements to MobileNets, SSDs and their combined applications.

---

<sup>1</sup><https://cocodataset.org/>

<sup>2</sup>[https://tfhub.dev/tensorflow/lite-model/ssd\\_mobilenet\\_v1/1/metadata/2](https://tfhub.dev/tensorflow/lite-model/ssd_mobilenet_v1/1/metadata/2)

### 3.3.1. MobileNet Architecture

TBD More Intro

Firstly, using depth-wise separable filters leads to a significantly lower number of parameters the model needs to learn during the training phase. Secondly, the easy shrinking of the model by adjusting the two hyper-parameters (named *width multiplier* and *resolution multiplier*) allows the model-builder to scale the model down to exactly the size that is appropriate for solving the problem the model is applied to [17].

### 3.3.2. Introduction to SSD

The SSD consists of a backbone and the *SSD head*. In theory, the backbone can be any of the networks mentioned in section 3.2.1 although the choice is not limited to this selection. When SSD was first introduced in 2015, the authors used the VGG-16 network as the backbone [27]. In the case of the specific SSD variant implemented in TUM-Lens, MobileNet is used as a backbone. The reasons for the particular applicability of MobileNet in our Android app use case are described in section 3.2.1. Independent of the specific type of backbone, it is first trained on publicly available image data, e.g. from a database like ImageNet<sup>3</sup>. Once the network is pre-trained the final layers that originally handled the classification are then replaced by the SSD head.

### 3.3.3. The SSD Head

The SSD head is the part of the detector that computes category scores and box locations. It is a characteristic of SSD to work with a fixed set of default bounding boxes. To account for the different sizes in which objects might appear in an image, SSD applies a sequence of comparatively small convolutional filters to the feature map returned by the backbone network. Each of the feature maps obtained by this process represents a different receptive field and allows the network to detect objects at different scales. Consequently, SSD computes predictions for all of these feature maps and for each of its predefined aspect ratios. It then fuses the predictions of the different aspect ratios and scales in order to improve detection accuracy. The range of aspect ratios and scales, although predetermined, enables the SSD to detect objects in various shapes and sizes. Its multi-scale feature maps are a core differentiator from other single-stage detectors like YOLO [27, 32].

---

<sup>3</sup><https://image-net.org/>

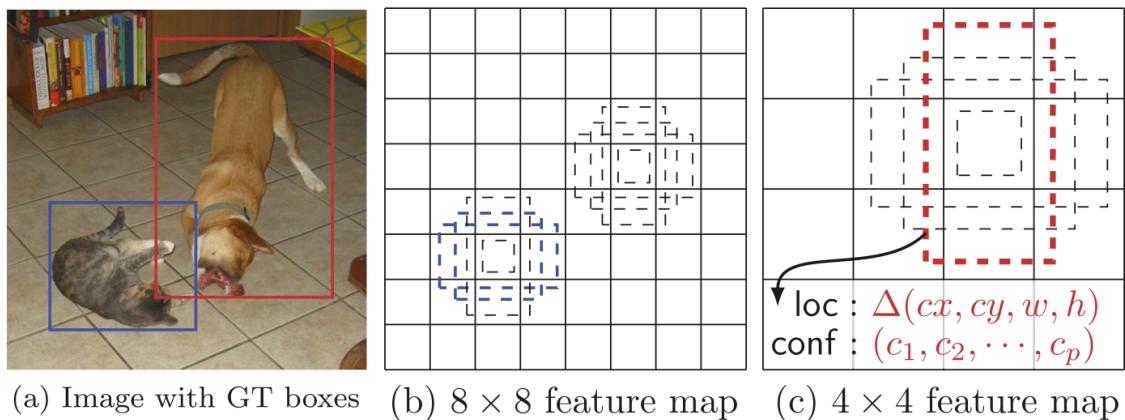


Figure 3.2.: The combination of bounding boxes in varying aspect ratios and multi-scale feature maps allows the SSD to detect objects in different sizes and orientations.  
Source: [27]

### 3.3.4. Architecture

### 3.3.5. Latest Improvements

# 4. App Development

## 4.1. Previous State of the Application

This practical part of this work was built in top of already existing Android app TUM-Lens in its version 1.0. The core feature of this previous app version is image classification. The app is already capable of classifying the images received from the live camera stream in real-time. It can also classify images that are loaded from the disk of the Android device instead of images received from the camera stream. Moreover, the user can choose between different TensorFlow-Slim<sup>1</sup> models to classify the images. This enables the user to compare the speed of similar detectors and also change the type of objects that can be detected as some of the available networks were trained on different data [19]. The most impactful design decision of the former developer of TUM-Lens was to run the entire classification locally on the Android device.

## 4.2. Development Goals

### 4.2.1. Migration From Java to Kotlin

During their Google I/O conference 2017 Google announced that they are making the programming language Kotlin a first-class citizen in Android [7]. Two years later, Google refined this statement announcing that Android development will "be increasingly Kotlin-first" and, at the present day, the Google Developers website recommends developers to choose Kotlin when they start building a new Android app [5].

#### Advantages of Kotlin

Continuing the development of TUM-Lens in Kotlin has several advantages. As it is an app currently solely developed by students, the enhanced expressiveness and conciseness over Java makes it easier for fellow students to understand the entire application within the scope of a bachelor's or master's thesis. Moreover, students can learn more in the given amount of time they spend working on the app because they can focus more strongly on the implementation of new features instead of having to write boilerplate Java code. As students are more error-prone than professional developers they surpassingly benefit from the safety features integrated into the Kotlin language. The most common cause for crashes of Java applications is the null pointer exception [16, 46]. With Kotlin being a strongly typed language a lot of these null pointer exceptions will be avoided because Kotlin's safety features help the developer to identify potential sources of such errors very easily. Google itself claims that Android apps are 20% less likely to crash when they contain Kotlin code [5].

---

<sup>1</sup><https://github.com/google-research/tf-slim>

### Pitfalls of Android Studio's Automatic Code Conversion

#### Code Reduction in Numbers

Using the npm distribution of a tool named Count Lines of Code (cloc)<sup>2</sup> we can easily count lines of Java and Kotlin code while ignoring blank lines and comment lines. This enables us to monitor how the shift from Java to Kotlin influenced the number of lines of code. With Perl, Node and npm installed on our machine, we run the following command which prints its output to stdout:

```
1 lens-david % npx cloc --by-file --git-diff-rel --match-d='main' --
    include-lang=Java,Kotlin -csv f6a18e0f 031a26a7
```

Listing 4.1: Command line input to print analysis of lines of code before and after the Java to Kotlin conversion. The output is also saved as cloc.csv and can be found in lens-david/thesis/raw\_data.

The `--by-file` option enables us to directly compare each Java file with Kotlin equivalent as it changes the output so that the results are returned for every source file encountered. `--git-diff-rel` interprets the command line arguments as git targets and compares only files which have changed in either commit which is exactly what we need to measure the change in the number of lines of code. Only files in the main directory are of our concern so we let cloc only search this directory using the option `--match-d='main'`. This excludes e.g. the test directory from being searched. With `--include-lang=Java,Kotlin` we limit the output to files written in the languages we seek to collate. Other files like Android's XML layout files are not of interest for this analysis. Finally, we provide two git commits. `f6a18e0f` is the last commit pushed by the previous developer of TUM-Lens and `bf0dbf7f` is the more recent commit that does not contain Java files anymore but replaced them with their respective Kotlin substitute.

#### 4.2.2. New Functionality: Object Detection

As mentioned in section 4.1, the existing image classification does not outsource the classification task to some remote server but performs the classification offline on the device itself. We adhere to the decision and implement the object detection analogously. The detection task is carried out locally using the same TensorFlow-Lite API as the image classification. We also display our detection on top of the live-stream received from the camera so that the user of the app can experience the object detection in real-time. All the same, there is some discrepancy to the existing set of functionalities built around the image classification. Only one model for object detection is currently available within the app and the detection logic cannot be applied to images that have been loaded from the storage of the device.

### 4.3. Implementing Object Detection Based on The TensorFlow Lite Framework

Graphic showing an overview of all app flows (old and new)

---

<sup>2</sup>Original cloc project by Al Danial: <https://github.com/AlDanial/cloc>; cloc npm distribution by Kent C. Dodds: <https://www.npmjs.com/package/cloc>

## *4. App Development*

---

**4.3.1. Some Deep**

**4.3.2. Dive Into**

**4.3.3. Object Detection**

**4.3.4. Implementation Fun**

**4.3.5. Next steps in the development of TUM-Lens**

Files migrated from Java to Kotlin (extension omitted)	Delta in lines of code
CameraRoll	-34
Classifier	4
ListSingleton	-9
PermissionDenied	-8
StartScreen	-12
ViewFinder	-7
fragments/CameraRollPredictionsFragment	9
fragments/CameraSettingsFragment	-15
fragments/ModelSelectorFragment	5
fragments/PredictionsFragment	4
fragments/ProcessingUnitSelectorFragment	-7
fragments/SmoothedPredictionsFragment	6
fragments/ThreadNumberFragment	-6
helpers/App	1
helpers/CameraEvents	0
helpers/FreezeAnalyzer	-6
helpers/FreezeCallback	0
helpers/ImageUtils	43
helpers/Logger	6
helpers/ModelConfig	-2
helpers/ProcessingUnit	-2
helpers/Recognition	-26
helpers/ResultItem	-21
helpers/ResultItemComparator	-2
<b>cumulative delta over all relevant files</b>	<b>-79</b>

Table 4.1.: This table shows the results from the command line prompt in listing 4.1. Packages have been indicated as a prefix to the file name to resemble the original project structure of TUM-Lens v1.0. Overall, the total size of the codebase shrank because of the migration from Java to Kotlin. This is particularly impressive as understandably further logic needed to be added to existing classes in order to account for the new object detection functionality.

# **5. Results**

## **5.1. Performance Evaluation**

## **5.2. Accuracy**

## **5.3. Possible Applications**

Die Anwendungsbereiche von Computer Vision sind zahlreich. Zu den regelmäßigen Aufgaben im Bereich

Organisation von Fotos auf dem Smartphone, ohne dass diese an die Server von Apple, & Google Co geschickt werden müssen (Gruppieren von Fotos, die zu einem Urlaub gehören, Gesichtserkennung, Objekterkennung)

Autonome Autos werden unabhängig von einer Verbindung zum Internet (5G, shared medium, bleibt das Auto im Tunnel dann stehen?)

## A. Screenshots of the Application

Compare old and new

## B. Tips With Greetings From the Chair

Here are tips along the way:

### B.1. Tips

#### B.1.1. How to Describe

When listing several points you have three basic options:

- |               |                |  |
|---------------|----------------|--|
| • itemize     | 1. itemize     | <b>itemize</b> short, unordered  |
| • enumerate   | 2. enumerate   | <b>enumerate</b> short ordered   |
| • description | 3. description | <b>description</b> listing of descriptions. Also nice for longer ones. |

#### B.1.2. How to Quote

”This is a quote!”

- Citations to a source can be made like this `\cite{grat117task} = [11]`  
Always join text and the citation with a non-breaking space: `text~\cite{foo}`.
- Referencing Sections, Figures, Tables, Formulas: `\autoref{sec:tips} = Appendix B.`
- Footnotes for url or further notes: `\footnote{\url{https://www.top500.org}} = 1`

#### B.1.3. How to Math

Use the align environment for equations especially if you want to align them somehow.

$$1 + 1 \neq 3 \tag{B.1}$$

$$\left(\frac{10}{1}\right) - 9 = 1 \tag{B.2}$$

---

<sup>1</sup><https://www.top500.org>

## B.2. Environments

### B.2.1. How to Figure

Anything can also be put in multiple columns.

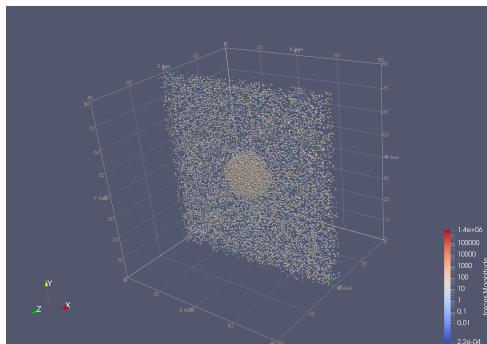


Figure B.1.: Some Caption. Always also include a source if it wasn't created by you!  
Source: [11]

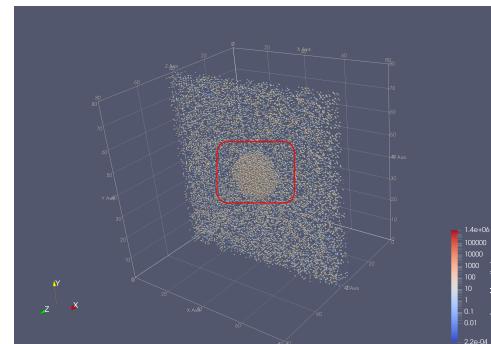
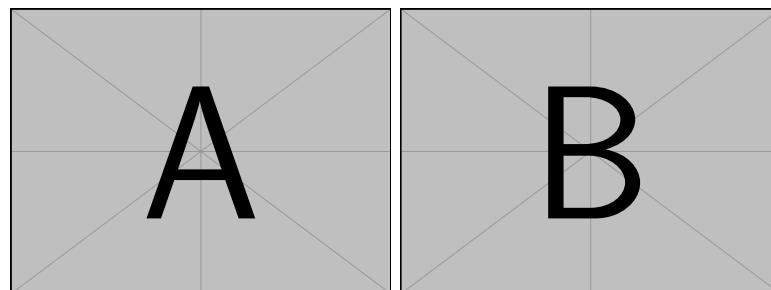


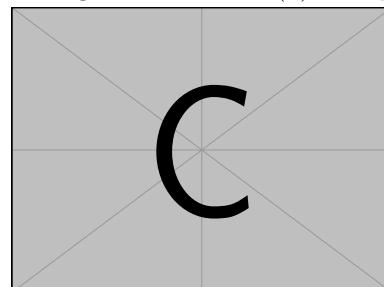
Figure B.2.: Figures can be drawn on or completely generated with tikz.

**Subfigures** If grouping of several pictures seems reasonable, think about using subfigures. This often comes in handy with plots.



(a) example-image-a

(b) example-image-b



(c) example-image-c

Figure B.3.: One caption to describe them all.

### B.2.2. How to Algorithm

---

**Algorithm 1:** Bogosort

---

**Input:** data array  
**Output:** data sorted

// Checks if array is sorted

1 **Function** is\_sorted(*data*):  
2   **for** *i*  $\leftarrow$  0 **to** *data.size()* - 1 **do**  
3     **if** *data*[*i*] > *data*[*i*+1] **then**  
4       **return** false  
5   **return** true

// actual algorithm

6 **Function** bogosort(*data*):  
7   **while** not is\_sorted(*data*) **do**  
8     **random.shuffle**(*data*)

---

Figure B.4.: some description what is happening

### B.2.3. How to Code

```

1 void runner(int type, void *data){
2     switch(type)
3     case taskType1:
4         // do stuff using data
5     case taskType2:
6         // do other stuff using data
7 }
```

Listing B.1: General form of a typical runner() function.

### B.2.4. How to Table

bla left	bla centered over two lines	bla right
bla left	bla centered	cell spanning two rows
cell spanning two columns		

Table B.1.: Fancy table that can contain line breaks and extended cells.

# List of Figures

2.1.	Classification of Machine Learning . . . . .	3
2.2.	2-Layered ANN . . . . .	5
2.3.	3-Layered ANN . . . . .	5
3.1.	Typical Computer Vision Tasks . . . . .	6
3.2.	Detecting Objects at Different Scales Through SSD's Feature Maps . . . . .	11
B.1.	Example Figure . . . . .	19
B.2.	Figure with tikz . . . . .	19
B.3.	One caption to describe them all. . . . .	19
B.4.	some description what is happening . . . . .	20

# List of Tables

3.1. Modern CNN Backbones and Detectors . . . . .	8
4.1. Java Files in TUM-Lens v1.0 And Change in Lines of Code After Their Conversion to Kotlin . . . . .	15
B.1. Some Table . . . . .	21

# Bibliography

- [1] *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Las Vegas, NV, USA). IEEE Computer Society, Dec. 2016. ISBN: 978-1-4673-8851-1.
- [2] Shivang Agarwal, Jean Ogier Du Terrail, and Frédéric Jurie. “Recent Advances in Object Detection in the Age of Deep Convolutional Neural Networks”. In: *arXiv preprint arXiv:1809.03193* (Sept. 2018).
- [3] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. “YOLOv4: Optimal Speed and Accuracy of Object Detection”. In: (Apr. 2020). URL: <https://arxiv.org/abs/2004.10934v1>.
- [4] “Convolutional Neural Networks Backbones for Object Detection”. In: *Image and Signal Processing. ICISP 2020. Lecture Notes in Computer Science* 12119 (June 2020), pp. 282–289. DOI: 10.1007/978-3-030-51935-3\_30.
- [5] Android Developers. *Android’s Kotlin-first approach*. 2021. URL: <https://developer.android.com/kotlin/first> (visited on 07/31/2021).
- [6] Android Developers. *Neural Networks API*. 2021. URL: <https://developer.android.com/ndk/guides/neuralnetworks> (visited on 07/31/2021).
- [7] Google Developers. *Google I/O Keynote (Google I/O ’17)*. May 2017. URL: <https://youtu.be/Y2VF8tmLFHw?t=5245> (visited on 07/31/2021).
- [8] Nikita Dvornik et al. “BlitzNet: A Real-Time Deep Network for Scene Understanding”. In: *Proceedings of the IEEE International Conference on Computer Vision* (Venice, Italy). Vol. 2017-October. Institute of Electrical and Electronics Engineers Inc., Oct. 2017, pp. 4174–4182. URL: <https://arxiv.org/abs/1708.02813v1>.
- [9] Ross Girshick et al. “Rich feature hierarchies for accurate object detection and semantic segmentation”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (Nov. 2014), pp. 580–587.
- [10] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. URL: <http://www.deeplearningbook.org>.
- [11] Fabio Alexander Gratl. “Task Based Parallelization of the Fast Multipole Method implementation of ls1-mardyn via QuickSched”. Master’s thesis. Garching: Institut für Informatik 5, Technische Universität München, Nov. 2017. URL: [https://www5.in.tum.de/pub/Gratl\\_MA\\_TaskBasedFMM.pdf](https://www5.in.tum.de/pub/Gratl_MA_TaskBasedFMM.pdf).
- [12] Jeff Hale. *Deep Learning Framework Power Scores 2018*. Sept. 2018. URL: <https://towardsdatascience.com/deep-learning-framework-power-scores-2018-23607ddf297a> (visited on 07/31/2021).

- [13] Matthew Halpern, Yuhao Zhu, and Vijay Janapa Reddi. “Mobile CPU’s rise to power: Quantifying the impact of generational mobile CPU design trends on performance, energy, and user satisfaction”. In: *Proceedings - International Symposium on High-Performance Computer Architecture* (Apr. 2016), pp. 64–76. DOI: 10.1109/HPCA.2016.7446054.
- [14] Kaiming He et al. “Deep Residual Learning for Image Recognition”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Las Vegas, NV, USA). IEEE Computer Society, Dec. 2016, pp. 770–778. ISBN: 978-1-4673-8851-1. DOI: 10.1109/CVPR.2016.90.
- [15] Wang Hechun and Zheng Xiaohong. “A Survey of Deep Learning-Based Object Detection”. In: *2nd International Conference on Big Data Technologies* (Jinan, China). Association for Computing Machinery, Aug. 2019, pp. 149–153. DOI: 10.1145/3358528.3358574.
- [16] Andras Horvath. *Which Java exceptions are the most frequent?* June 2018. URL: <https://blog.samebug.io/which-java-exceptions-are-the-most-frequent-f830b113c37f> (visited on 07/31/2021).
- [17] Andrew G. Howard et al. “MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications”. In: (Apr. 2017). URL: <https://arxiv.org/abs/1704.04861v1>.
- [18] Max Jokel. *TUM-Lens*. Version 1.0. July 31, 2021. URL: <https://play.google.com/store/apps/details?id=com.maxjokel.lens>.
- [19] Maximilian Jokel. “Implementing a TensorFlow-Slim based Android app for image classification”. Bachelor’s thesis. Garching: Lehrstuhl für Wissenschaftliches Rechnen, Technische Universität München, Nov. 2020. URL: <https://mediatum.ub.tum.de/1579885>.
- [20] Tao Kong et al. “HyperNet: Towards accurate region proposal generation and joint object detection”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Las Vegas, NV, USA). IEEE Computer Society, Dec. 2016, pp. 845–853. ISBN: 978-1-4673-8851-1. DOI: 10.1109/CVPR.2016.98.
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. “Imagenet classification with deep convolutional neural networks”. In: *Advances in neural information processing systems* 25 (2012), pp. 1097–1105.
- [22] Erik G Learned-Miller. “Introduction to supervised learning”. In: *I: Department of Computer Science, University of Massachusetts* (Feb. 2014).
- [23] Fei-Fei Li, Ranjay Krishna, and Danfei Xu. *CS231n: Convolutional Neural Networks for Visual Recognition*. 2021. URL: <https://cs231n.github.io/neural-networks-1/> (visited on 07/31/2021).
- [24] Fei-Fei Li, Ranjay Krishna, and Danfei Xu. *Detection and Segmentation*. May 2021. URL: [http://cs231n.stanford.edu/slides/2021/lecture\\_15.pdf](http://cs231n.stanford.edu/slides/2021/lecture_15.pdf) (visited on 07/31/2021).

## Bibliography

---

- [25] Tsung Yi Lin et al. “Focal Loss for Dense Object Detection”. In: *Proceedings of the IEEE International Conference on Computer Vision* (Venice, Italy). Vol. 2017-October. Institute of Electrical and Electronics Engineers Inc., Oct. 2017, pp. 2999–3007. DOI: 10.1109/ICCV.2017.324. URL: <http://arxiv.org/abs/1708.02002>.
- [26] Jie Liu et al. “Performance Analysis and Characterization of Training Deep Learning Models on Mobile Device”. In: *Proceedings of the International Conference on Parallel and Distributed Systems - ICPADS* (Dec. 2019), pp. 506–515. DOI: 10.1109/ICPADS47876.2019.00077.
- [27] Wei Liu et al. “SSD: Single Shot MultiBox Detector”. In: *Lecture Notes in Computer Science* (2016), pp. 21–37. ISSN: 1611-3349. DOI: 10.1007/978-3-319-46448-0\_2. URL: [http://dx.doi.org/10.1007/978-3-319-46448-0\\_2](http://dx.doi.org/10.1007/978-3-319-46448-0_2).
- [28] Tom M. Mitchell. *Machine Learning*. Ed. by Erick M. Munson. McGraw-Hill, Mar. 1997, p. xv. ISBN: 0-07-042807-7.
- [29] Andrew Ng. *Machine Learning Course by Stanford University*. 2021. URL: <https://www.coursera.org/learn/machine-learning> (visited on 07/31/2021).
- [30] *Proceedings of the IEEE International Conference on Computer Vision* (Venice, Italy). Vol. 2017-October. Institute of Electrical and Electronics Engineers Inc., Oct. 2017.
- [31] Team PyTorch. *PyTorch 1.3 adds mobile, privacy, quantization, and named tensors*. Oct. 2019. URL: <https://pytorch.org/blog/pytorch-1-dot-3-adds-mobile-privacy-quantization-and-named-tensors/> (visited on 07/31/2021).
- [32] Joseph Redmon et al. “You Only Look Once: Unified, Real-Time Object Detection”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Las Vegas, NV, USA). IEEE Computer Society, Dec. 2016, pp. 779–788. ISBN: 978-1-4673-8851-1. DOI: 10.1109/CVPR.2016.91.
- [33] Ahmad EL Sallab et al. “Deep reinforcement learning framework for autonomous driving”. In: *Electronic Imaging 2017* (2017), pp. 70–76.
- [34] A. L. Samuel. “Some Studies in Machine Learning Using the Game of Checkers”. In: *IBM Journal of Research and Development* 3 (3 July 1959), pp. 210–229. DOI: 10.1147/RD.33.0210. URL: <http://ieeexplore.ieee.org/document/5392560/>.
- [35] A. L. Samuel. “Some Studies in Machine Learning Using the Game of Checkers. II - Recent Progress”. In: *IBM Journal of Research and Development* 11 (6 Nov. 1967), pp. 601–617. DOI: 10.1147/RD.116.0601.
- [36] R. Sathya and Annamma Abraham. “Comparison of Supervised and Unsupervised Learning Algorithms for Pattern Classification”. In: *International Journal of Advanced Research in Artificial Intelligence* 2 (2 2013), pp. 34–38. ISSN: 2165-4069.
- [37] Laurent Sifre. “Rigid-Motion Scattering For Image Classification”. PhD thesis. École Polytechnique, CMAP, 2014. URL: [https://www.di.ens.fr/data/publications/papers/phd\\_sifre.pdf](https://www.di.ens.fr/data/publications/papers/phd_sifre.pdf).
- [38] Karen Simonyan and Andrew Zisserman. “Very Deep Convolutional Networks for Large-Scale Image Recognition”. In: *3rd International Conference on Learning Representations* (San Diego, CA, USA). ICLR, Sept. 2014. URL: <https://arxiv.org/abs/1409.1556v6>.

- [39] Petru Soviany and Radu Tudor Ionescu. “Optimizing the Trade-off between Single-Stage and Two-Stage Object Detectors using Image Difficulty Prediction”. In: *20th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing* (Timișoara, Romania). Institute of Electrical and Electronics Engineers Inc., Sept. 2018, pp. 209–214. URL: <https://arxiv.org/abs/1803.08707v3>.
- [40] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. Ed. by Frances Bach. second edition. The MIT Press, 2018. ISBN: 9780262039246.
- [41] Christian Szegedy et al. “Going deeper with convolutions”. In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Boston, MA, USA). IEEE Computer Society, June 2015, pp. 1–9. ISBN: 978-1-4673-6964-0. DOI: 10.1109/CVPR.2015.7298594.
- [42] James Vincent. *Google’s new machine learning framework is going to put more AI on your phone*. May 2017. URL: <https://www.theverge.com/2017/5/17/15645908/google-ai-tensorflow-lite-machine-learning-announcement-io-2017> (visited on 07/31/2021).
- [43] James Vincent. *The iPhone X’s new neural engine exemplifies Apple’s approach to AI*. Sept. 2017. URL: <https://www.theverge.com/2017/9/13/16300464/apple-iphone-x-ai-neural-engine> (visited on 07/31/2021).
- [44] Oriol Vinyals et al. “Grandmaster level in StarCraft II using multi-agent reinforcement learning”. In: *Nature* 575 (2019), pp. 350–354.
- [45] Mengwei Xu et al. “DeepType: On-Device Deep Learning for Input Personalization Service with Minimal Privacy Concern”. In: *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 2 (4 Dec. 2018), pp. 1–26. DOI: 10.1145/3287075. URL: <https://dl.acm.org/doi/abs/10.1145/3287075>.
- [46] Alex Zhitnitsky. *The Top 10 Exception Types in Production Java Applications - Based on 1B Events*. June 2016. URL: <https://www.overops.com/blog/the-top-10-exceptions-types-in-production-java-applications-based-on-1b-events/> (visited on 07/31/2021).
- [47] Xiaojin Zhu and Andrew B. Goldberg. “Introduction to Semi-Supervised Learning”. In: *Synthesis Lectures on Artificial Intelligence and Machine Learning* 3 (June 2009), pp. 1–130. DOI: 10.2200/S00196ED1V01Y200906AIM006. URL: <https://doi.org/10.2200/S00196ED1V01Y200906AIM006>.
- [48] Yousong Zhu et al. “CoupleNet: Coupling Global Structure with Local Parts for Object Detection”. In: *Proceedings of the IEEE International Conference on Computer Vision* (Venice, Italy). Vol. 2017-October. Institute of Electrical and Electronics Engineers Inc., Oct. 2017, pp. 4146–4154. DOI: 10.1109/ICCV.2017.444. URL: <https://doi.org/10.1109/ICCV.2017.444>.

## Acronyms

**ANN** Artificial Neural Network.

**API** application programming interface.

**CNN** Convolutional Neural Network.

**CV** Computer Vision.

**ML** Machine Learning.

**NNAPI** Neural Networks API.

**ResNet** Residual Neural Network.

# Glossary

**application programming interface** set of functions and procedures allowing the creation of applications that access the features or data of an operating system, application, or other service.

**differential privacy** protects an individual's information essentially as if her information were not used in the analysis at all, in the sense that the outcome of a differentially private algorithm is approximately the same whether the individual's information was used or not.