

# Energetsko učinkovit sistem za detekcijo slonov s pomočjo strojnega učenja

**Ključne besede:** strojno učenje, mikrokrmilnik, sklepanje na napravi, termalna kamera, sistem z majhno porabo, **UKD:** XXXXX

## Povzetek

### 1. Uvod

*Konflikti med ljudmi in sloni predstavljajo velik problem ohranjanja populacije slonov. Zaradi fragmentacije in pomanjkanja habitata sloni, v iskanju hrane, pogosto zaidejo na rizeva polja in plantaže, kjer pridejo v stik s človekom. Po podatkih skupnosti WILDLABS, zaradi konfliktov, samo v Indiji, letno umre povprečno 400 ljudi in 100 slonov. Sistemi zgodnje opozoritve nadomeščajo vlogo človeških stražarjev in opozarjajo bližnjo skupnost o bližini, potencialno nevarnih, slonov in tako pripomorejo k zmanjševanju konfliktov med ljudmi in sloni.*

*V tem magistrskem delu predstavljam strukturo sistema zgodnje opozoritve, ki je sestavljen iz večih, nizko porabnih, vgrajenih sistemov, ki so opremljeni s termalnimi kamerami in ene dostopne točke oz. prehoda (gateway). Vgrajeni sistemi so postavljeni na terenu, ob zaznavi slona pošljejo opozorilo preko brezžičnega omrežja do dostopne točke, ki nato lahko opozori lokalno skupnost. Za prepoznavo slonov iz zajetih termalnih slik smo uporabili metode strojnega učenja, bolj specifično konvolucijske nevronske mreže. Glavni cilji tega magistrskega dela so bili zasnova, izvedba in ovrednotenje modelov strojnega učenja, ki jih je možno poganjati na mikrokrmilnkih pod pogoji nizke porabe.*

## *2. Teoretični opis gradnikov sistema*

*V tem poglavju opisujemo osnovna znanja, ki jih bralec potrebuje za razumevanje tega magistrskega dela. Opišemo kako lahko strojno učenje pomaga reševati probleme, ki bi s klasičnimi tehnikami zahtevali kompleksne rešitve. Podrobno predstavimo izvajanje modelov strojnega učenja v vgrajenih sistemih. Ugotovimo, da kljub omejitvam kot so nizke procesorske moči in majhni spomini, prednosti kot so hitra odzivnost na dogodke, zmanjšanje porabe zaradi manjše potrebe po pošiljanju podatkov v oblak in povečane stopnje zasebnosti, hitro odtehtajo slabosti. Lotimo se opisa nevronskih mrež, aktivacijskih funkcij, konvolucijskih nevronskih mrež in tehnik prenosnega učenja. Predstavimo tudi platformo TensorFlow Lite for microcontrollers, ki nam je omogočila implementacijo nevronskih mrež na mikrokrmilnikih. Naredimo pregled možnih brezžičnih tehnologij in argumentiramo zakaj smo se odločili za LoRaWAN. Nazadnje opišemo delovanje termalnih kamer in predstavimo kako je potekala izbira optimalne termalne kamere. Izbrana kamera je bila FLIR Lepton.*

## *3. Zasnova modela nevronske mreže*

*V tem poglavju podrobno opišemo celoten proces zasnove modela, ki je sposoben klasificirati termalne slike in predvideti kateri objekt je na sliki. Pri zasnovi smo se omejili na prepoznavo 4 različnih razredov: sloni, ljudje, krave in narava oz. nakjučni objekt. Zadali smo si cilj, da klasificiranje termalne slike ne sme trajati več kot 1 sekundo in da mora biti model dovolj majhen, da ga lahko naložimo na mikrokrmilnik. Na začetku opišemo orodja in delovno okolje, ki smo jih uporabljali pri zasnovi modela (uporabljali smo platformo TensorFlow), nato pa se lotimo analize nabora termalnih slik, ki jih je zbral podjetje Arribada Initiative. Iz nabora termalnih slik smo izbrali slike, ki so ustrezale našim zahtevam. Nabor termalnih slik je vseboval veliko število slik slonov in ljudi, ampak ne veliko slik krav ali narave. Slednje smo posneli sami na terenu, s hitrim prototipom, ki smo ga izdelali sami.*

*Opisali smo kako smo so slike pripravili za učenje modela in predstavili smo osnovno strukturo modela. Ker je iskanje optimalnih hiperparametrov nehevristična naloga, smo določili možni razpon hiperparametrov in izvedli algoritem naključnega iskanja,*

*ki je naučil večje število modelov z različnimi hiperparameteri. Opisali smo tudi, kako poteka optimizacija modelov, ki bodo tekli na mikrokrmilnikih.*

*Nazadnje ponovno opišemo potek zasnove modela od začetka do konca, ampak tokrat to storimo s Edge Impulse Studijem.*

#### *4. Zasnova in izvedba sistema zgodnje opozoritve*

*V četrtem poglavju predstavimo sprva generalne gradnike sistema in njihove funkcije, nato pa predstavimo konkretne komponente. Odločili smo se za sistem z dvema mikrokrmilnikoma. Mikrokrmilnik NRF52840 kontrolira delovanje celotnega sistema in prezivi večino časa v režimu nizke porabe. Ob signalu iz PIR sensorja se zbudi iz spanja in vklopi drugi mikrokrmilnik, STM32F767ZI. STM32F767ZI je visoko zmogljiv mikrokrmilnik s Cortex-M7 jedrom. Povezan je s FLIR Lepton termalno kamero in kontrolira zajemanje slik ter njivovo procesiranje s nevronsko mrežo. STM32F767ZI sporoči rezultate klasifikacije NRF52840 mikrokrmilniku, ki jih obdela in nato pošlje preko LoRaWAN omrežja na dostopno točko. Za LoRa brezžični modul smo uporabili LR1110 čip.*

*Veliki del magistrskega dela se je ukvarjal s prenosom TensorFlow Lite for Microcontrollers platfrome na platformo naše izbire, libopenCM3. V procesu prenosa smo se podrobno spoznali s prevajanjem in povezovanjem kode, saj nismo uporabljali programskega okolja, ki bi to naredilo za nas. Tako smo ustvarili odprto-kodni projekt MicroML, ki omogoča uporabo TensorFlow lite kode na platformi libopenCM3. Sestava in uporabo MicroML-a smo podrobno opisali. MicroML smo uporabili pri pisanju kode za mikrokrmilnik STM32F767ZI, za NRF52840 pa smo uporabili operacijski sistem Zephyr.*

#### *5. Meritve in rezultati*

*TODO*

#### *6. Povzetek*

*TODO*

# Energy efficient system for detection of elephants with Machine Learning

**Key words:** machine learning, microcontroller, on-device inference, thermal camera, low-power system

**UKD:** XXXXX

## Abstract

*Human-Elephant Conflicts are a major problem in terms of elephant conservation. According to WILDLABS, an average of 400 people and 100 elephants are killed every year in India alone because of it. Early warning systems replace the role of human watchers and warn local communities of nearby, potentially life threatening, elephants, thus minimising the Human-Elephant Conflicts.*

*In this Master's thesis we present the structure of an early warning system, which consists of several, low-power embedded systems equipped with thermal cameras and a single gateway. To detect elephants from captured thermal images we used Machine Learning methods, specifically Convolutional Neural Networks. The main focus of this thesis was the design, implementation and evaluation of Machine Learning models running on microcontrollers under low-power conditions. We designed and trained several accurate image classification models, optimised them for on-device deployment and compared them against models trained with commercial software in terms of accuracy, inference speed and size. While writing firmware, we ported a part of TensorFlow library and created our own build system, suitable for libopencm3 platform. We also implemented reporting of inference results over LoRaWAN network and described possible server-size solution. We finally constructed fully functional embedded system from various development and evaluation boards, and evaluated its performance in terms of power consumption. We show that embedded systems with Machine Learning capabilities are a viable solution to many real life problems.*

## List of Abbreviations

WWF	World Wide Fund for Nature
HEC	Human-Elephant Conflicts
ML	Machine Learning
NN	Neural Networks
CNN	Convolutional Neural Networks
DNN	Deep Neural Networks
IoT	Internet of Things
RMS	Root Mean Square
ReLU	Rectified Linear Activation Unit
ISM	Industrial, Scientific and Medical
3GPP	The 3rd Generation Partnership Project
LoRa	Long Range
IR	Infrared
EM	Electromagnetic
LWIR	Long Wave Infrared Region
ROIC	Readout Integrated Circuit
VOx	Vanadium-Oxide
NETD	Noise Equivalent Temperature Difference
CPU	Central processing unit
FPA	Focal Point Array
TWI	Two Wire Interface

SPI	Serial Peripheral Interface Bus
PIR	Passive Infrared Sensor
I2C	Inter-Integrated Circuit
MOSI	Master Out Slave Input
MISO	Master Input Slave Out
FPA	Focal Point Array
AGC	Automatic Gain Control
SYS	System Information
VID	Video Processing Control
OEM	Original Equipment Manufacturer
RAD	Radiometry
GCC	the Gnu Compiler Collection
TTN	The Things Network
DK	Development Kit
EVK	Evaluation Kit
GNSS	Global Navigation Satellite System
DWT	Data Watchpoint Trigger
TTN	The Things Network
SQL	Structured Query Language

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Early warning system . . . . .	3
1.3	IRNAS and the Arribada Initiative . . . . .	4
1.4	Reasoning for the Machine Learning approach . . . . .	5
1.4.1	Implementation of Machine Learning algorithms . . . . .	6
1.4.2	Edge Impulse . . . . .	7
1.5	Objective . . . . .	7
1.6	Master's thesis outline . . . . .	8
<b>2</b>	<b>Theoretical description of system building blocks</b>	<b>9</b>
2.1	Machine Learning . . . . .	9
2.1.1	General Machine Learning workflow . . . . .	11
2.1.2	Machine Learning on embedded devices . . . . .	12
2.2	Neural Networks . . . . .	13
2.2.1	Activation functions . . . . .	14
2.2.2	Backpropagation . . . . .	15
2.2.3	Convolutional Neural Networks . . . . .	17
2.2.3.1	Convolutional layers . . . . .	18
2.2.3.2	Pooling layers . . . . .	20
2.3	Transfer Learning . . . . .	20
2.4	TensorFlow . . . . .	22
2.4.1	TensorFlow Lite for Microcontrollers . . . . .	23
2.4.1.1	Post-training quantization . . . . .	24

2.5	IoT and wireless technologies . . . . .	25
2.5.1	LoRa and LoRaWAN . . . . .	26
2.6	Thermal cameras . . . . .	28
2.6.1	Choosing the thermal camera . . . . .	30
<b>3</b>	<b>Neural Network model design</b>	<b>34</b>
3.1	Model objectives . . . . .	34
3.2	Tools and development environment . . . . .	36
3.3	Creating the dataset . . . . .	37
3.4	Exploring the dataset . . . . .	39
3.5	Image preprocessing . . . . .	42
3.6	Model creation and training . . . . .	44
3.7	Model optimisation . . . . .	46
3.8	Neural Network model design in Edge Impulse Studio . . . . .	47
<b>4</b>	<b>Design and implementation of the early warning system</b>	<b>50</b>
4.1	Hardware . . . . .	52
4.1.1	Nucleo-F767ZI . . . . .	52
4.1.2	nRF52840 DK . . . . .	53
4.1.3	LR1110 development kit . . . . .	54
4.1.4	Boost converter evaluation kit . . . . .	54
4.1.5	FLIR Lepton 2.5 camera module and Lepton breakout board	55
4.1.6	PIR Sensor . . . . .	57
4.2	Firmware . . . . .	58
4.2.1	Tools and development environment . . . . .	58
4.2.1.1	Development environment for STM32f767ZI . . . . .	58
4.2.1.2	Development environment for nRF52840 . . . . .	59
4.2.2	Architecture design . . . . .	59
4.2.3	MicroML and build system . . . . .	61
4.2.4	Running inference on a microcontroller . . . . .	64
4.3	Server-side components and software . . . . .	67

<b>5 Measurements and results</b>	<b>70</b>
5.1 Comparison of models . . . . .	70
5.1.1 Hyperparameter search space and results analysis . . . . .	70
5.1.2 Comparison of selected, re-trained models . . . . .	75
5.1.3 Comparison of Edge Impulse models . . . . .	78
5.2 On-device performance testing . . . . .	80
5.2.1 Comparison of code sizes . . . . .	82
5.2.2 Comparison of different optimisation options . . . . .	83
5.2.3 Scoring trained models . . . . .	84
5.3 Power profiling of an embedded early warning system . . . . .	87

# List of Figures

1.1	Various Human-Elephant Conflicts. . . . .	2
1.2	Diagram of an early warning system. . . . .	3
2.1	Workflow diagram of solving a generic Machine Learning problem. . . . .	11
2.2	Mathematical model of artificial neuron and fully connected 3-layer neural network. . . . .	13
2.3	Different activation functions and their equations. . . . .	15
2.4	Structure of a Convolutional Neural Network. . . . .	17
2.5	Dot product operation between filter and zero-padded input matrix. .	18
2.6	Convolutional layer with five different filters. . . . .	19
2.7	Polling layer. . . . .	20
2.8	Transfer Learning. . . . .	21
2.9	Workflow of preparing an ML model for an inference on a microcontroller.	24
2.10	Properties of Lora signal. . . . .	27
2.11	Comparison between a normal image and thermal image . . . . .	28
2.12	Focal point array and bolometer. . . . .	29
2.13	Comparison of images of the same object taken with cameras with different NETD values. . . . .	31
2.14	Comparison of image quality made by different thermal cameras. . . .	32
3.1	Camera trap used in Assam, India. . . . .	37
3.2	Thermal images of elephants from the dataset. . . . .	38
3.3	Camera setup used for taking thermal images with FLIR Lepton 2.5.	39
3.4	Distribution of used images from thermal dataset depending on image location and type of sensor. . . . .	41

3.5	Class distribution of thermal images. . . . .	41
3.6	Image preprocessing pipeline. . . . .	42
3.7	Distribution of training images before and after resampling. . . . .	44
3.8	CNN architecture written in Python using Keras Sequential API. . .	46
3.9	Creating a Neural Network in Edge Impulse Studio. . . . .	48
3.10	Training performance report. . . . .	49
4.1	General block diagram of an embedded system . . . . .	51
4.2	Diagram describing behavior of embedded early warning system . .	51
4.3	Hardware diagram of embedded early warning system . . . . .	52
4.4	Nucleo-F767ZI development board . . . . .	53
4.5	nRF52840DK development board . . . . .	53
4.6	MAX17225ENT+T boost converter breakout board . . . . .	54
4.7	MAX17225ENT+T boost converter breakout board . . . . .	55
4.8	Front and back side of FLIR Lepton breakout board with thermal camera module inserted. . . . .	56
4.9	Command and control interface of FLIR Lepton camera. . . . .	57
4.10	Front and back side of FLIR Lepton breakout board with thermal camera module inserted. . . . .	58
4.11	Architecture diagram of the firmware that is running on a STM32F767ZI microcontroller. . . . .	59
4.12	Architecture diagram of the firmware that is running on a nRF52840 microcontroller. . . . .	60
4.13	Directory structure of MicroML project. . . . .	62
4.14	Build system of MicroML project. . . . .	63
4.15	Example of TensorFlow Lite inference code in C++. . . . .	65
4.16	Server side flow of information. . . . .	67
4.17	Node-RED flow . . . . .	68
4.18	Example of Grafana graph. . . . .	69
5.1	Confusion matrices of <i>0a</i> model (left) and <i>460b</i> model (right). . . . .	77
5.2	Comparison of time of inference of different models. . . . .	81

5.3	Example output of arm-none-eabi-size command.	82
5.4	Comparison of Flash and RAM size of compiled example models.	82
5.5	Inference time of <i>Ob</i> model using different optimisations.	83
5.6	Inference time of <i>Ob</i> model using different optimisations.	86

## List of Tables

5.1	First hyperparameter search space . . . . .	71
5.2	Partial results of first random search of hyperparameters . . . . .	72
5.3	Second hyperparameter search space . . . . .	74
5.4	Partial results of second random search of hyperparameters . . . . .	74
5.5	Properties of selected models . . . . .	76
5.6	Precision and recall metrics of trained models . . . . .	76
5.7	Properties of Edge Impulse models using CNN architecture. . . . .	79
5.8	Properties of Edge Impulse models using Transfer Learning technique	79
5.9	Precision and recall metrics of trained Edge Impulse models . . . . .	79

# 1 Introduction

## 1.1 Motivation

As a result of increasing human population and habitat loss, human-wildlife conflicts have become increasingly common in recent decades [1]. According to the organisation The World Wide Fund for Nature (WWF), human-wildlife conflicts are defined as: "any interaction between humans and wild animals, that results in negative impacts on human social, economic or cultural life, on the conservation of wildlife populations, or on the environment" [2]. These conflicts range from mostly harmless, non-violent contacts, such as sightings of wildlife animals in urban areas, to the destruction of crops and infrastructure, killing of livestock, and, in the worst cases, loss of human lives. In more severe cases these conflicts end in defensive or retaliatory killings of wild animals, which can drive an already endangered species to extinction.

Polar bears, tigers, and elephants are generally considered to be the most problematic [1]. In the Arctic, as a consequence of the reduction of their natural habitat, polar bears are drawn to human settlements and food dumps while searching for food [3]. Unexpected encounters can become deadly for both sides. As wide-ranging animals, tigers need large areas where they can roam, hunt, and breed [4]. When their natural prey population is depleted, they often turn their attention to poorly protected livestock. Their attacks often have economic, social, and psychological consequences. According to WILDLABS, tigers killed 101 people between the years 2013 and 2016, in India alone [4].

As herbivores, elephants might be seen as less problematic when compared to polar bears or tigers, but this assumption could not be further from the truth.

Although exact numbers vary between sources, casualties from Human-Elephant Conflicts (HEC) are much higher compared to conflicts involving polar bears or tigers. According to WILDLABS, an average of 400 people and 100 elephants are killed every year in India [5]. The leading cause of death of elephants is electrocution (by electric fences, unprotected power lines), followed by train accidents, poaching, and poisoning [6]. Reasons for HEC are similar to the conflicts with polar bears and tigers. Their habitat is shrunk continuously and replaced with crop fields and plantations. As their food options decrease, surrounding agricultural landscapes become inviting. Various HECs can be seen in figure 1.1.



Figure 1.1: Various Human-Elephant Conflicts. Top left: two elephants running from a mob hurling flaming tar balls, top-right: palm plantation owner inspecting damage done by elephants to the crops, bottom-left: an elephant crossing the protective barrier, bottom-right: An elephant that died because of HEC. Image sources: [5] [7] [8] [9]

One of the HEC hotspots is in the Sonitpur District, Assam Province, India. In a  $5300 \text{ km}^2$  large area around 200,000 people and 200 elephants share the same space [5]. Elephants often venture into paddy fields which represent livelihood for local communities. A single elephant can quickly trample fields of rice crops in a few hours, causing big financial problems to already impoverished farmers [5].

Several measures have been taken to minimise HEC: Electrical fences, watch towers, trenches, chilli-based deterrents, regular patrols, usage of trained captive elephants and camera traps with motion sensors.

Although the above mentioned measures function to some degree, they are not effective enough, since they are unreliable, or come into effect when the damage has already been done [10].

## 1.2 Early warning system

One important component of minimizing Human-Elephant Conflicts is a reliable early warning system. A system capable of detecting the presence of nearby elephants would warn nearby communities and give them enough time to prepare and respond non-violently. The same kind of system would also provide information about common elephant paths, thus giving farmers knowledge on how to construct and place their fences better to minimise HEC. The system (Figure 1.2) should consist of several small, deployed devices with mounted cameras that will detect elephants, and one server which will aggregate alerts and forward them to mobile phones and computers, where the local community will see them.

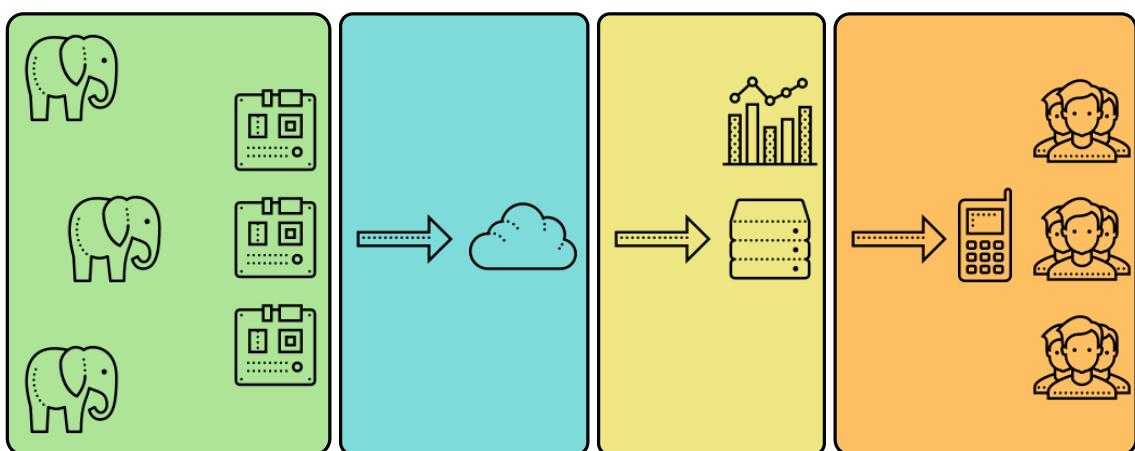


Figure 1.2: Diagram of an early warning system. Icons source: [11]

Although most of the villages in the Sonitpur district have access to cell phones and the internet, the connectivity can be unreliable [5]. Furthermore, devices will be placed quite far from the main server, which makes sending a large amount of data a

problem. This limits the choice of wireless networks to long range, low bandwidth technologies. It is, therefore, preferable that elephant detection is done on deployed devices, and only results (which can be few bytes big) are sent over some radio network to the main server. Deployed devices will be placed in forests, fields, with no access to electricity, therefore, they need to be battery powered. Low maintenance of deployed devices is a desirable quality, which means that they should be functional for longer periods without any human interactions. To achieve that with a limited power source, they should be energy efficient, equipped with solar panels, and a low power radio. Devices should spend most of their time in sleep mode, conserving energy, only waking up to take a photo, processing it, and sending results to the main server. A thermal camera is needed, as most of the Human-Elephant Conflicts happen during the night [5].

Elephant recognition can be done with the help of a convolutional neural network (CNN) running inference on a microcontroller. Making this possible and evaluating the solution is the focus of this master's thesis.

### 1.3 IRNAS and the Arribada Initiative

The system described above is currently in development at IRNAS in the collaboration with Arribada Initiative. The Slovenia based Institute IRNAS offers a complete development service, starting with an idea on paper and ending with a finished product. Its previous projects cover a wide range of different fields, from free space optical systems, bio-printing, to Internet of things (IoT) solutions that cover various industrial and nature conservative use cases. The Arribada Initiative is a London based team, that uses open source solutions for purposes of nature conservation. As the winner of WWF and WILDLABS Human Wildlife Conflict Tech Challenge [12], Arribada received funding to develop an early warning system. They spent some time in Assam, India, where they tested a proof-of-concept design [13]. They decided on devices with thermal cameras, as Human-Elephant conflicts often happen during the night. The sensor of choice was an FLIR Lepton 2.5 and/or 3.5. They also created a large dataset of elephant images while filming elephants in Whipsnade Zoo in the

United Kingdom. This dataset will be used for training the Neural Networks and will be discussed in Section 3.4. To create a final embedded system with on device Machine Learning, Arribada chose to work with IRNAS.

## 1.4 Reasoning for the Machine Learning approach

Today Machine Learning (ML) is present in many products that we use on a daily basis. It can be found in email spam detection, recommendation algorithms on Facebook and YouTube, speech recognition on smartphones and medical applications. ML can help us solve problems that are hard to solve by conventional methods. For example, to develop an email spam detector, a programmer would have to write a program that would scan the content of an email while checking for the common words and phrases that appear in spam emails and flag the email as spam if they would be found. This would take several iterative cycles of writing the rules, evaluating the solution and analysing the possible mistakes. Even if a possible deterministic solution would be made, it would not stand the test of time, as new forms of spam emails would emerge, tricking the system.

Compare that to a Machine Learning approach. Given enough examples of spam and normal emails, we can train an ML algorithm and let it to discover by itself the rules that mandate what is a spam and what is a normal email. The program would be much smaller compared to the one made by the conventional approach. After the program is deployed into the real world, we can use it to store new data and relearn, thus always adapting to new changes. This process can be automated.

Same parallels can be drawn to recognising elephants from thermal images. It is an difficult task to write a deterministic algorithm that could identify an elephant successfully from an image and not confuse it for a human or livestock. Using an ML approach we can train the algorithm on a image dataset and let it figure out the connections between the images and correct labels.

### 1.4.1 Implementation of Machine Learning algorithms

Since ML algorithms are at their core basic math operations, they can be implemented (although maybe not efficiently) in any programming language and on any hardware platform from scratch. However, to avoid reinventing the wheel and wasting the time on algorithm optimisation, it is more logical to use one of the available ML frameworks. Frameworks such as scikit-learn, Keras, Caffe and TensorFlow enable programmers to write application code in a high level language such as Python, which, at run time, translates to efficient C/C++ code. These frameworks take away low-level details of ML algorithms, enabling programmers to deal only with application code and not its underlying implementation.

Over the past several years there has been a growing desire to expand ML applications to embedded devices. Executing ML algorithms directly on microcontrollers has its benefits and presents new challenges compared to executing those algorithms on PC's or servers. These comparisons are described further in Section 2.1.2. There are several frameworks, proprietary and open-sourced, that can be used to develop ML applications on microcontrollers. STMicroelectronics created X-CUBE-AI, a tool that converts the pre-trained model created by one of the various Deep Learning frameworks into an optimised library. X-CUBE-AI works only with STM32 microcontrollers, and is proprietary. Another framework, TensorFlow Lite for microcontrollers, was created as an extension of TensorFlow Lite, which is used to develop ML models for mobile applications. It provides converter tools and C++ implementations of common ML operations. In 2019 another framework,  $\mu$ Tensor, was merged with TensorFlow Lite, providing it with support for an efficient CMSIS-NN library developed by ARM.

For this thesis we used TensorFlow Lite for microcontrollers. It can be used with any family of microcontrollers, and is open-source, so we can study its internal code.

### 1.4.2 Edge Impulse

Regardless of the many ML frameworks on the market, companies that specialise in ML on embedded devices are scarce. One of them is Edge Impulse, which is a recently founded company from San Jose, USA. They provide users with an end to end web solution for developing ML applications for embedded devices. Instead of designing and writing specific programs that deal with preprocessing of training data and creation and training of ML models, Edge Impulse does this automatically for their customers. After the ML model is created and converted into an optimised format for embedded systems, customers get the immediate first approximation of how much RAM and FLASH memory the model will take and how fast it will run. We can then run the model on the local machine, or deploy to a variety of different platforms, such as Arduino, STM32, OpenMV, Zeyphr and others.

Their solution will be used as a benchmark for our work with TensorFlow Lite.

## 1.5 Objective

The objective of this Master's thesis is to evaluate the feasibility to recognize animals, especially elephants, from thermal images, with Machine Learning algorithms, running directly on a microcontroller.

Additionally we will design and build a system capable of detecting an elephant with the help of Machine Learning algorithms in the day or at night. Detection of an elephant needs to be reported over a wireless network to an application server.

For that we will:

- Train a Neural Network model capable of classifying elephants, humans and other animals from thermal images.
- Optimise Neural Network model for on device inference.
- Implement on device inference on an STM32 microcontroller using TensorFlow Lite.

- Compare the performance of our implementation against the Edge Impulse implementation.
- Build a system around the STM32 microcontroller with a thermal camera and wireless network support.
- Profile power consumption of the built system.
- Establish system requirements for different ML applications.

## 1.6 Master's thesis outline

This chapter provided an overview of motivation and the companies involved, some reasoning for choosing the Machine Learning approach and the objectives of this thesis. Chapter 2 provides a theoretical description of the system building blocks. Machine Learning, Neural Networks, thermal cameras, TensorFlow Lite, and others topics are discussed there. Chapter 3 revolves around the design of an image classification Neural Network model. Chapter 4 describes the design and implementation of an early warning system, from the hardware, firmware and software points of view. In Chapter 5 we describe the measurement procedure and results. Chapter 6 presents our findings, describes the limitations of our project, and suggest paths for further research.

## 2 Theoretical description of system building blocks

### 2.1 Machine Learning

According to Arthur Samuel (qtd. in Geron [14]) Machine Learning is a field of study that gives computers the ability to learn without being programmed explicitly. This ability to learn is the property of various Machine Learning algorithms. We will be using the terms "Machine Learning" and "learning" interchangeably. To learn, these learning algorithms need to be trained on a collection of examples of some phenomenon [15]. These collections are called **datasets**, and can be generated artificially or collected in nature.

To understand how the ML approach can solve problems better, we can examine an example application. Let us say that we would like to build a system that can predict a type of animal movement based on accelerometer data. To train its learning algorithm, also known as a **model**, we need to train it on a dataset that contains accelerometer measurements of different types of movement, such as walking, running, jumping and standing still. Input to the system could be either raw measurements from all three axes, or components extracted from raw measurements such as RMS, spectral power, peak frequency and/or peak amplitude. These inputs are also known as **features**, they are values that describe the phenomenon being observed [15]. The output of the system would be a predicted type of movement. Although we would mark each example of measurement data with what type of movement it represents, we would not define the relationship between the two directly. Instead, we would let the model figure out the connection by itself, through the process of training. The trained model should be general enough so that it can predict the type of movement on unseen accelerometer data correctly. Performance of the model is

measured in **accuracy**. It tells us what percentage of input data has been predicted correctly.

There exists a large variety of different learning algorithms. We can categorise them broadly in several ways, and one of them depends on how much supervision the learning algorithm needs in the training process. Algorithms like K-nearest neighbours, linear and logistic regression, Support Vector Machines fall into the category of supervised learning algorithms. Training data that is fed into them includes solutions, also known as **labels** [14]. The above described example is an example of a supervised learning problem.

Algorithms like k-Means and Expectation Maximization, fall into the category of unsupervised learning algorithms. Here, training data is unlabelled, algorithms are trying to find similarities in data by themselves [14]. Other categories exist, such as semi-supervised learning which is a combination of the previous two and reinforcement learning, where the model acts inside the environment according to learned policies [14].

Neural Networks, algorithms inspired by neurons in human brains [14] [16], can fall into either of the categories. They are appropriate for solving complex problems like image classification, speech recognition, and autonomous driving, but they require a large amount of data and computing power for training. They fall into the field of Deep Learning, which is a sub-field of Machine Learning.

Training of ML models is computationally demanding, and is usually done on powerful servers or computers with dedicated Graphic Processing Units to speed up training time. After a model has been trained, data can be fed in and a prediction is computed. This process is also known as **inference**. The inference is computationally less intensive compared to the training process, so, with properly optimised models, we can run inference on personal computers, smartphones, tablets, and even directly in internet browsers.

### 2.1.1 General Machine Learning workflow

There are several steps in ML workflow that need to happen to get from an idea to a working ML based system, and this is represented in Figure 2.1.

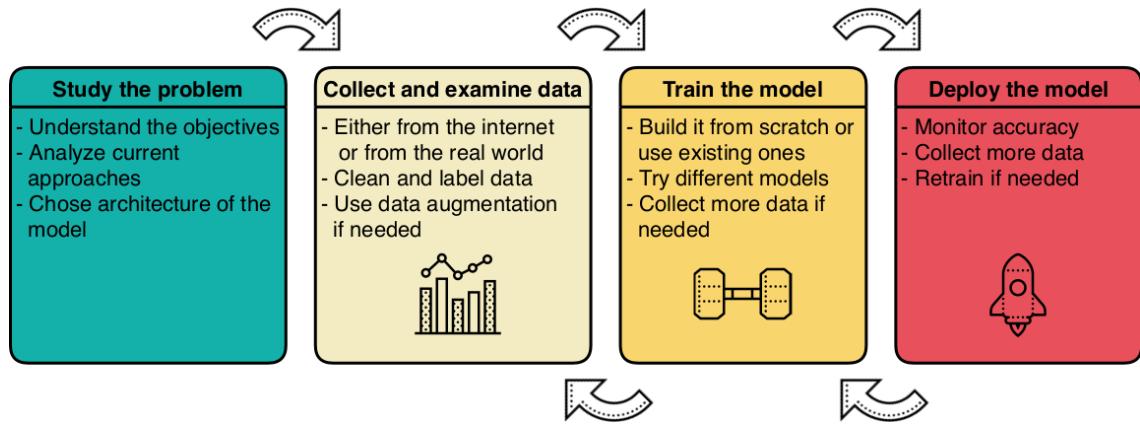


Figure 2.1: Workflow diagram of solving a generic Machine Learning problem. Icons source: [11]

First, the problem has to be studied, it has to be understood what are the objectives, what are current solutions and which approach should be used. Here, we decide on the rough type of ML model that we will use, based on the problem. In the second step we collect and clean up data. We should always strive to collect a large amount of quality and diverse data that represent a real world phenomenon. Collecting that kind of data can be hard and expensive, but we can use various tools, such as data augmentation or data synthesis, thus increasing data size and variety. Sometimes data are not collected by us, in that case we should examine them and extract information that we need. Third, we train the ML model. We might create something from scratch or use an existing model. We can train several different types of models and choose the one that performs the best. To achieve the desired accuracy, steps two and three can be repeated many times. In step four we deploy our model and monitor its accuracy. If accuracy drops we can always collect new data and retrain the model.

### 2.1.2 Machine Learning on embedded devices

Machine Learning on embedded devices is an emerging field, which coincides nicely with the Internet of Things. Resources about it are limited, especially when compared to the vast number of resources connected with Machine Learning on computers or servers. Most of the information about it can be found in the form of scientific papers, blog posts and Machine Learning framework documentation [17] [18] [19].

Running learning algorithms directly on embedded devices comes with many benefits.

**Reduced power consumption** is one of them. In most IoT applications devices send raw sensor data over a wireless network to the server, which processes it either for visualisation or for making informed decisions about the system as a whole. Wireless communication is one of the more power hungry operations that embedded devices can do, while computation is one of more energy efficient [19]. For example, a Bluetooth communication might use up to 100 milliwatts, while a MobileNetV2 image classification network running 1 inference per second would use up to 110 microwatts [19]. As deployed devices are usually battery powered, it is important to keep any wireless communication to a minimum, so minimising the amount of data that we send is paramount. Instead of sending everything we capture, is much more efficient to process raw data on the devices and only send the results.

Another benefit of using ML on embedded devices is **decreased latency time**. If the devices can extract high-level information from raw data, they can act on it immediately, instead of sending it to the cloud and waiting for a response. Getting a result now takes milliseconds, instead of seconds.

Such benefits do come with some drawbacks. Embedded devices are a more resource constrained environment when compared to personal computers or servers. Because of limited processing power, it is not feasible to train ML models directly on microcontrollers. Also it is not feasible to do online learning with microcontrollers, meaning that they would learn while being deployed. Models also need to be small enough to fit on a device. Most general purpose microcontrollers only offer several hundred kilobytes of flash, up to 2 megabytes. For comparison, the MobileNet v1

image classification model, optimised for mobile phones, is 16.9 MB in size [20]. To make it fit on a microcontroller and still have space for our application, it would have to be simplified.

The usual workflow while developing Machine Learning models for microcontrollers, is to train a model on training data on a computer. When we are satisfied with the accuracy of the model we quantize it, and convert it into a format understandable by our microcontroller. This is described further in Section 2.4.1.

## 2.2 Neural Networks

Although the first models of Neural Networks (NN) were presented in 1943 (by McCulloch and Pitts) [14] and hailed as the starting markers of the Artificial Intelligence era, several decades of research and technological progress had to pass before they could be applied to practical, everyday problems. Early models of NNs, such as the one proposed by McCulloch and Pitts, were inspired by how real biological neural systems work. They proved that a very simple model of an artificial neuron, with one or more binary inputs and one binary output, is capable of computing any logical proposition when used as a part of a larger network [14].

To learn how NNs work we can refer to Figure 2.2a, which shows a generic version of an artificial neuron.

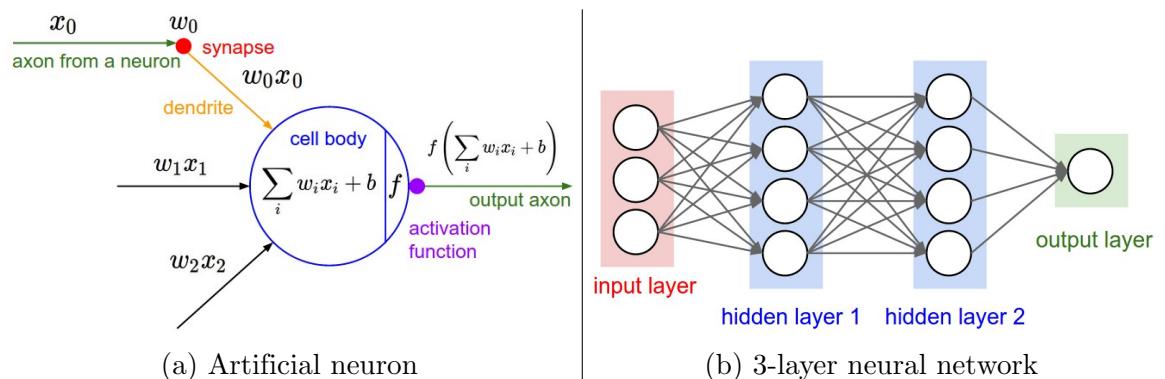


Figure 2.2: (a) Mathematical model of an artificial neuron, similarities with biological neurons can be seen. (b) Fully connected 3-layer neural network. Image source: [16]

A Neuron takes several inputs, multiplies each input with its **weight** and sums them

up. A **bias** parameter is added to the sum, which is then passed to an activation function.

NNs consist of many neurons, which are organised into **layers**. Neurons inside the same layer do not share any connections, but they connect to the layers before and after them. The first layer is known as the **input** layer and last one is known as the **output** layer. Any layers between are said to be **hidden**. In Figure 2.2b we can see a neural network with an input layer with three inputs, two hidden layers with four neurons each, and an output layer with just one neuron. If all inputs of neurons in one layer are connected to all outputs from the previous layer, we say that a layer is **fully connected** or **dense**, Figure 2.2b is an example of one. NNs with many hidden layers fall into the category of Deep Neural Networks (DNN).

### 2.2.1 Activation functions

Activation functions introduce non-linearity to a chain of otherwise linear transformations, which enables ANNs to approximate any continuous function [14]. There are many different kinds of activation functions, as seen on Figure 2.3, such as sigmoid function and rectified linear activation function (ReLU). A sigmoid function was used commonly in the past, as it was seen as a good model for a firing rate of a biological neuron: 0 when not firing at all, and 1 when fully saturated and firing at maximal frequency [16]. It takes a real number and squeezes it into a range between 0 and 1. It was later shown that training NNs with sigmoid activation function often hinders the training process, as saturated outputs cut off parts of networks, thus preventing the training algorithm from reaching all neurons and configuring the weights correctly [16]. It has since fallen out of practice, and is nowadays replaced by ReLu or some other activation function.

Another commonly used activation function is a softmax function (seen in 2.1, which takes a vector as an input, computes the exponential of every element and divides that with the sum of exponentials of all elements [14] The input vector of values becomes a vector of probabilities. Softmax is usually used as an activation in the

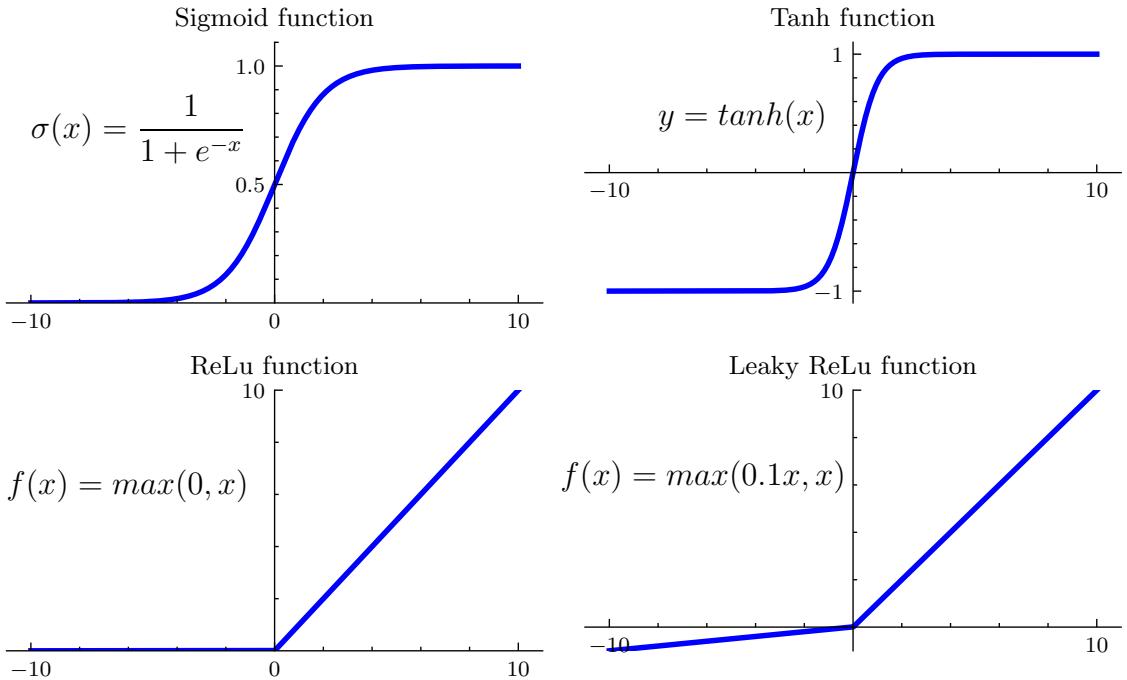


Figure 2.3: Different activation functions and their equations.

last layer of a classifier network.

$$\sigma(y_i) = \frac{e^{y_i}}{\sum_{j=1}^K e^{y_j}} \quad \text{for } i = 1, \dots, K \text{ and } y = (y_1, \dots, y_k) \in \mathbb{R}^K \quad (2.1)$$

Where:

$y$  - Input vector

$K$  - Number of elements in the input vector

$\sigma(y_i)$  - Computed probability of the  $i$ -th element in the input vector

### 2.2.2 Backpropagation

Training of Neural Networks is done with a training algorithm, known as **backpropagation**. As mentioned before, we train the Neural Network by showing it a large amount of training data with labels. At the start of the training phase, all weights and biases are set to randomly small values. During each training step, a Neural Network is shown a small batch of training data. Each instance is fed into the NN and the final output label is calculated. This is known as **forward pass**, which is

the same as making predictions, except that intermediate results are stored from each neuron from every layer. Calculated output is compared to an expected one using a **loss** (also known as **cost**) function. The loss function returns a single value, which tells us how badly our NN performing: the higher it is, the worse is our NN performing. The goal is to minimise the loss function, thus increasing the accuracy of our NN. In the context of multivariable calculus, this means that we have to calculate the negative gradient of weights and biases, which will tell us in which direction we have to change each weight and bias so that the value of loss function decreases.

Doing this for all weights and biases at the same time would be complicated, so the backpropagation algorithm does this in steps. After computing the loss function, the algorithm calculates analytically how much each output connection contributed to the loss function (essentially the local gradient) with the help of previously stored intermediate values. This step is done recursively for each layer until the first input layer is reached. At that moment the algorithm knows in which direction each weight and bias should change, so that the value of the loss function lowers. A procedure is then performed, known as a **Gradient Descent**. All local gradients are multiplied with a small number known as a **learning rate**, and then subtracted from all weights and biases. This way, in each step we change weights and biases slowly in the right direction, while minimising the loss function. Gradient Descent is not only used when training neural networks but also when training other ML algorithms.

We do not have to execute a backpropagation algorithm for each training instance, instead, we can calculate predictions for a small set of training data, calculate the average loss function and then apply backpropagation.

---

<sup>1</sup>Scientists Weisel and Hubel showed that different cells in the primary visual cortex of a cat responded differently to different visual stimuli [16]. Some were activated when shown a horizontal line in a specific location, some were activated by vertical lines. More complex cells responded to boxes, circles and so on. CNNs also detect simpler shapes first and use them to detect more complex ones later.

### 2.2.3 Convolutional Neural Networks

Convolutional Neural Networks (CNN) are a kind of Neural Networks that work especially well with image data. Like NNs they have found inspiration in nature, in their case the visual cortex of the brain.<sup>1</sup>

In Figure 2.4 we can see an example of CNN which takes an image of a car as an input and outputs probability results in five different classes.

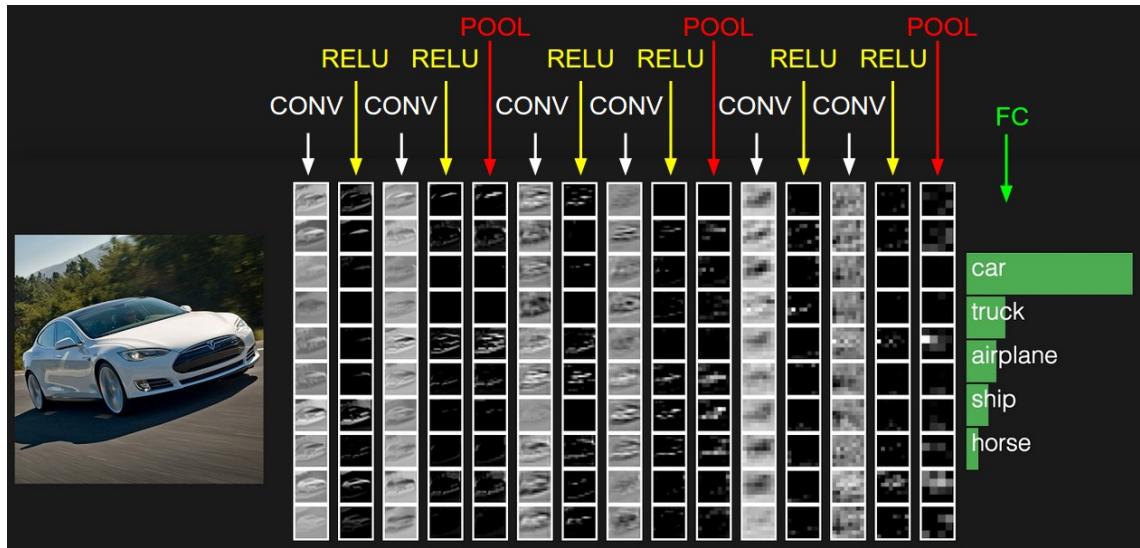


Figure 2.4: Structure of a Convolutional Neural Network. Image source: [16]

Specific to CNNs are two different types of layers, **convolutional** layers and **pooling** layers. Each convolutional layer detects some sort of shapes: the first ones detect different kinds of edges, while later ones detect more complex shapes and objects, like wheels, legs, eyes, ears. Pooling layers downsample the data in the spatial dimension, thus decreasing the number of parameters and operations needed in CNN. After a few alternating pairs of convolutional and pooling layers, the output of the last pooling layer is flattened out into one dimensional vector and fed into a fully connected NN which produces probability results in given classes.

It makes sense to explain how convolutional and pooling layers work in greater detail, as this will be important later when we will be designing our CNN models in Section 3.6.

### 2.2.3.1 Convolutional layers

Data that CNNs operate on are 3 dimensional matrices, where width and height correspond to an image resolution, and depth corresponds to the number of colour channels, 3 for colourful images (red, green, blue) and 1 for greyscale. When speaking about these matrices we will refer to them as volumes.

Convolution layers perform dot products between input volume and several **filters** or **kernels** to produce output volume. In these layers, filters are configured through the training phase. We can see a concrete example in Figure 2.5. 2D filter with size  $2 \times 2$  covers a part of the input volume, over which element-wise multiplication is computed, elements are summed and the result is written into the first element of output volume.

Input	Kernel	Output
$\begin{matrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 0 \\ 0 & 3 & 4 & 5 & 0 \\ 0 & 6 & 7 & 8 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{matrix}$	$*$	$\begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix}$
	=	$\begin{matrix} 0 & 3 & 8 & 4 \\ 9 & 19 & 25 & 10 \\ 21 & 37 & 43 & 16 \\ 6 & 7 & 8 & 0 \end{matrix}$

Figure 2.5: Dot product operation between filter and zero-padded input matrix.  
Image sources: [21]

The filter then moves a fixed distance or **stride** and the process is repeated. It is important to note that, although we can choose the width and height of the filter, the depth of the filter is always equal to the depth of the input volume. If the depth is larger than one, then dot products are done for each 2D matrix in the depth dimension separately, and then an element-wise sum operation is performed between these matrices. To avoid losing information from the image pixels that are on the edges (as they would be included in dot products fewer times compared to central ones) we often pad input images with zeros.

The size of output volume depends on several factors as seen in 2.2.

$$V_o = (V_i - F + 2P)/S + 1 \quad (2.2)$$

Where:

$V_i$  - Input volume size, only width or height

$V_o$  - Output volume size, only width or height

$F$  - Filter or receptive field size

$P$  - Amount of zero padding used on the border

$S$  - Stride length

If we examine the example in Figure 2.5 we can see that input with a size 3 x 3, stride 1, padding 1 and filter with a size 2 x 2, produces an output with size 4 x 4.

The depth of output volume is equal to the number of filters used in the convolutional layer as seen in Figure 2.6, it is a norm that a single convolutional layer uses a large number of filters to produce a deep output volume [16]. It is also common to set padding, stride and filter size so that the width and height of the input volume are preserved. This prevents the information at the edges from being lost too quickly [16].

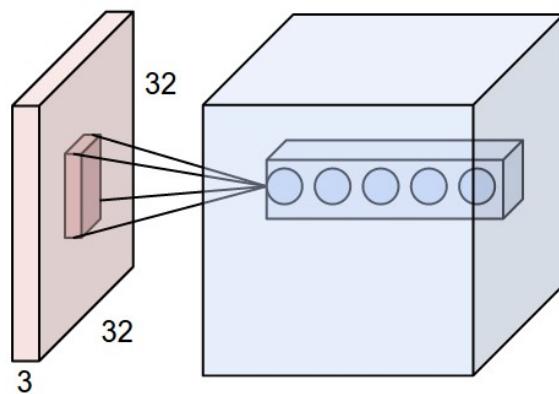


Figure 2.6: Convolutional layer with five different filters. Image sources: [21] [16]

At the end of the convolutional layer the output volume is fed into neurons similar to the one described in Section 2.2. All elements in the same depth are affected by

the same bias term and fed into the activation function. The size of the volume is preserved in this step.

### 2.2.3.2 Pooling layers

Pooling layers perform the downsampling of input volumes in both width and height dimensions. This is done by sliding a filter of fixed size over the input and doing a MAX operation on elements that the filter covers, and only the largest value element is copied into the output (Figure 2.7). Pooling is done on each depth slice separately from other slices, so depth size is preserved through the layer.

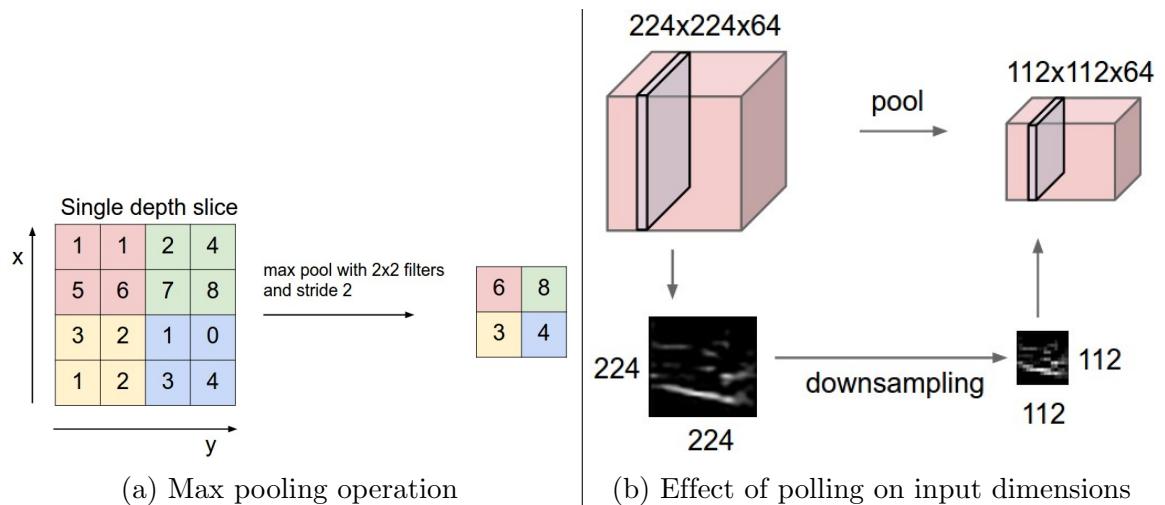


Figure 2.7: Pooling layer. Image source: [16]

It is common to select pool size  $2 \times 2$  and stride 2. Like this, inputs are downsampled by two in height and width dimensions, discarding 75 % of activations. Pooling layers therefore reduce the number of activations, and prepare them to be flattened out and fed into a fully connected layer.

## 2.3 Transfer Learning

Transfer Learning is a ML technique, which transfers knowledge that was learned on one problem, to another, similar problem.

We mentioned earlier in Section 2.2.3 that CNNs learn in first few layers how to detect simple shapes, like edges and lines, which are in later layers used to detect

boxes, circles, faces and so on. This behavior is not limited only to CNNs, but to all DNNs.

The task of detecting simple shapes in first few layers of DNNs can be common between models that are trying to solve a similar problem. Example would be a model classifying images depicting different types of objects and a model classifying images with different types of animals. First model is trying to solve a more general problem than the second one, however, the learned weights in first few layers from first model can be reused in the second model.

And this is exactly what Transfer Learning is all about, it is possible to train a big, complex model on a large dataset and then reuse first few layers in a different model. Example Transfer learning is shown in Figure 2.8. It is important to note that the weights of transferred layers are fixed, they are not being trained.

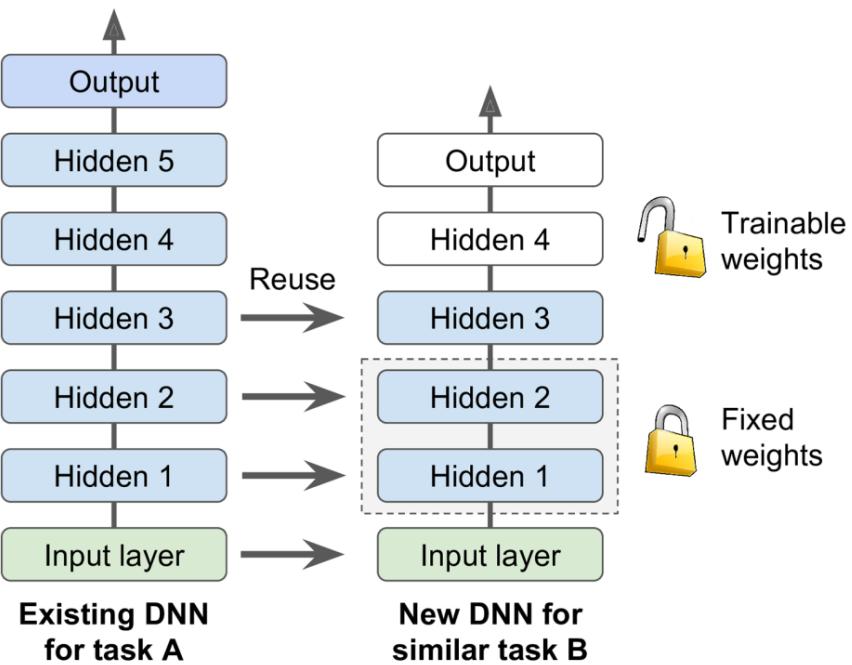


Figure 2.8: Transfer Learning. Image source: [14]

Benefits of using Transfer learning are:

- **Shorter training time.** Using a pre-trained model with fixed weights and adding few trainable layers on its top greatly lower training time, simply because there are less parameters to train.

- **No need for large datasets.** Training DNN usually requires large amount of data, which we might not have. Using a pre-trained model, that was trained on a large dataset, and applying it to a specific small dataset yielded better performance than just training a model from scratch [14].

## 2.4 TensorFlow

TensorFlow is a free and open-source framework for numerical computation. It is particularly suited for large-scale Machine Learning applications [14]. It started as a proprietary project developed by a Google Brain team at Google in 2011, and became open-source in late 2015. It is used in many of Google's products such as Gmail, Google Cloud Speech and Google Search.

TensorFlow gives programmers tools for creating and training ML models, without needlessly diving into the specifics of computing Neural Networks. Programmers can write high level code in Python API, which calls a highly efficient C++ code. When using TensorFlow, the hardest part of an ML project is usually data preparation. After that is done, the creation of an ML model, its training and evaluation can be done in a few lines of Python code.

TensorFlow also supports Keras high level API for building ML models. Keras is a Python library that functions as a wrapper for TensorFlow. When building ML models developers can use Keras Sequential API, where each layer in a model is represented as one line of code. Users do not need to care about connections between the layers, they only need to choose the type of layer (convolutional, max pool, fully connected), its size and a few other specific parameters. Sequential API is used most of the time, but if a finer level of control is needed TensorFlow provides low level math operations as well.

Finally, TensorFlow's trained model is portable [14]. Models can be trained in one environment and executed in another. This means that we can train our model by writing Python code on a Linux machine and execute it with Java on an Android device. This last functionality is important for running ML models on microcontrollers.

### 2.4.1 TensorFlow Lite for Microcontrollers

TensorFlow Lite (TFLite) is a set of tools and libraries that enable running ML inferences on constrained devices [22]. It provides support for Android and iOS devices, and embedded Linux. TensorFlow Lite for Microcontrollers (TFLite Micro) is a recent port of TFLite (as of mid 2019), dedicated to running ML models on microcontrollers. TFLite itself provides API in different languages, such as Java, Swift, Python and C++. TFLite Micro uses C++ API, specifically C++11, which reuses a large part of the general TensorFlow codebase.

TFLite Micro library does not require any specific hardware peripherals, which means that the same C++ code can be compiled to run on a microcontroller or a personal computer with minimal changes. Users are only expected to implement their version of the `printf()` function. As microcontroller binaries are usually quite big, flashing firmware to a microcontroller is a time consuming process. It makes sense to first test and debug a program that includes only ML inference specific code on a personal computer, before moving on to a microcontroller, to save time. Implementation of the test setup is described in 4.2.3.

The TFLite Micro library is available publicly as a part of a much larger TensorFlow project on GitHub [22]. To use the library for embedded development the whole project has to be cloned from the GitHub. The TensorFlow team provides users with several example projects that have been ported to several different platforms, such as Mbed, Arduino, OpenMV and ESP32. Example projects show how to use TFLite API, while showcasing different ML applications: Motion detection, wake word detection and person detection.

Is important to know that TFLite is just an extension of the existing TensorFlow project. General steps for creating a trained ML model are still the same as seen in Figure 2.1, although we have to be aware of some details. Figure 2.9 shows all steps that are needed to prepare an ML model for running on a microcontroller.

We start with a small but inaccurate model that can still accomplish the basic criteria that our objective demands. When the end of this workflow process is reached and we

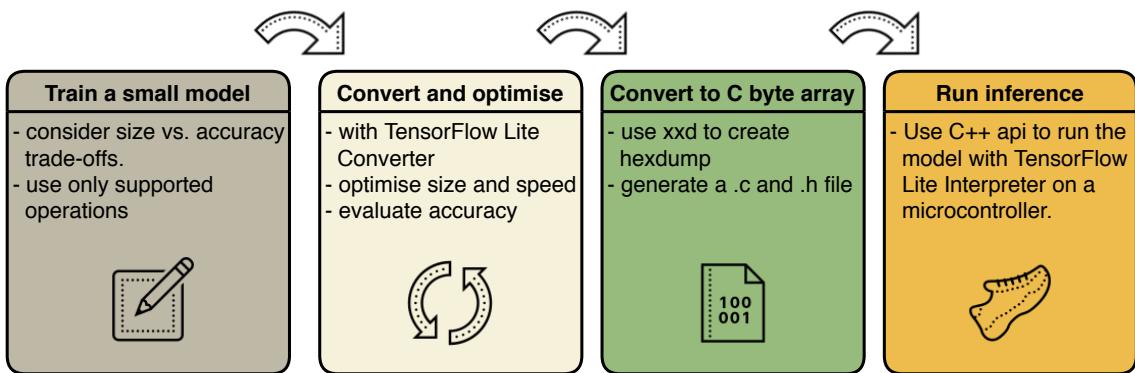


Figure 2.9: Workflow of preparing a ML model for an inference on a microcontroller. Icons source: [11]

made sure that the model can fit into a flash memory area of our target microcontroller, we can start training more complex models to increase accuracy. We are allowed to use only operations that have supported implementations on microcontrollers. This is usually not a restriction, as many of them are supported.

The model that is created is usually quite big, and needs to be converted with the TensorFlow Lite Converter tool. This tool provides a non-optimised conversion and several different optimised conversions.

To import and use the optimised model, we need to convert it into binary format, which is done with the command line tool `xxd`. The model is then ready to be executed on a microcontroller, we can run it and process the results. Accuracy will be the same as compared to running the same `.tflite` model on a personal computer, but execution time will naturally be different. If needed, we can tweak the model parameters, train a new model and repeat the described workflow again.

#### 2.4.1.1 Post-training quantization

By using quantization optimisation we approximate floating-point numbers in a different format, usually with 8-bit integers. When computing Neural Networks we can quantize weights, biases and intermediate values output by separate neurons. Quantization has a dramatic effect on the size of the model and its execution speed. By changing 32-bit floating-point numbers with 8-bit integers size decreases by a

factor of 4. Floating-point math is by nature slow to compute, many microcontrollers do not even have a floating-point unit. In comparison integer math is faster to compute, therefore quantized models are executed faster. Model accuracy decreases after using quantization, but usually by less than a percent.

## 2.5 IoT and wireless technologies

The Internet of Things, or IoT, is a system of uniquely identifiable devices, which communicate with each other or other systems over wireless networks [23]. A device or a thing is a battery powered embedded system such as smart watch, heart monitor, or animal tracker which would transmit collected sensor data to an IoT gateway, which would relay the data over to the cloud. These data can then be analysed and displayed in such a fashion which would provide businesses or users with valuable information. Examples of this would be tracking the energy consumption of machines in factories, monitoring conditions of crops in agriculture, or monitoring locations of endangered species in African conservation parks.

An important part of the IoT system is a wireless network that is used to transport data from edge devices to gateways, or directly to the Internet. The choice of a wireless network is highly dependent on a type of problem an IoT solution is trying to solve. Factors such as required battery life, amount of data being sent, the distance that data have to travel and environment conditions of the edge device itself are important.

Because our early warning system demands a decent battery life of several months and needs to send a small amount of data over one or two kilometres, we will focus on wireless technologies such as NB-IoT, Sigfox and LoRa.

Narrowband IoT or NB-IoT is a radio technology standard developed by the 3GPP standard organisation [24]. NB-IoT was made specifically with embedded devices in mind, it has a range of up to 15 km and it has deep indoor penetration [24]. Compared to Sigfox and LoRa it has better latency and a higher data rate, but also higher power consumption [25]. However it is unsuitable for our use case as it

operates on the network provided by the cellular base towers, which is inconvenient as the mobile connection in Assam, India can be inconsistent [5].

Sigfox is a radio technology developed by the company of the same name that operates on an unlicensed Industrial, Scientific and Medical (ISM) radio band. In many views it is similar to LoRa, as it has the comparative range and power consumption [25]. However, there are a few important differences. Although Sigfox modules are a bit cheaper when compared to Lora modules, each message is paid, devices are limited to 12 bytes per uplink, 140 uplinks per day and only 4 downlinks are available per day. Sigfox devices can also only communicate with base stations, installed by the Sigfox company [25]. This means that users can not build their own network and are dependent on the coverage provided by Sigfox.

This leaves us with the Lora protocol, which covers our use case from points view of long range, low power consumption and the ability to set up our own network.

### 2.5.1 LoRa and LoRaWAN

LoRa (Long Range) is a physical layer protocol that defines how information is modulated and transmitted over the air [26] [24]. The protocol is proprietary and owned by a semiconductor company, Semtech, who is the sole designer and manufacturer of Lora radio chips in the world. The LoRa protocol uses a modulation similar to chirp spread spectrum modulation [26]. As the protocol is proprietary, exact details of it are not known, although it was reverse engineered by a radio frequency specialist [27]. An example of a LoRa signal that was captured with a software defined radio can be seen in Figure 2.10a. Each symbol is modulated into a radio signal whose frequency is either increasing or decreasing with a constant rate inside of a specified bandwidth. When the bandwidth boundary is reached, the signal "wraps around" and appears at the other boundary. Although the frequency is always changing with a constant rate, it is not continuous inside the bandwidth window, but it can change to a different frequency immediately and continue from there.

This kind of modulation gives LoRa extreme resiliency against the interference of other radio frequency signals that might be using the same frequency band [26] [28].

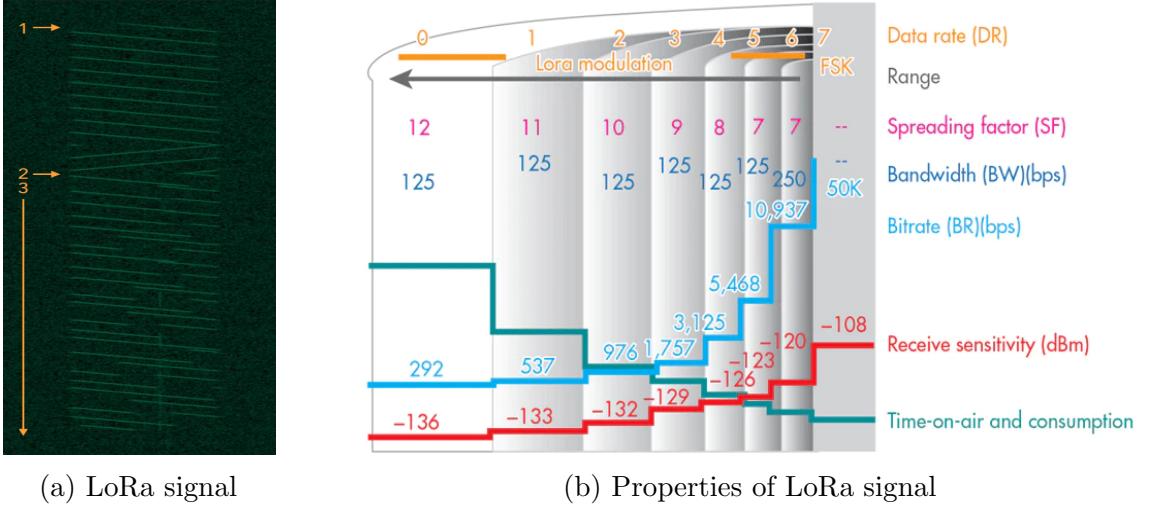


Figure 2.10: LoRa signal (left) and different properties of LoRa with their effects on range, bit rate, receiver sensitivity, time on air and consumption (right). Image sources: [27] [28]

For example, on a lower part of Figure 2.10a we can see a signal with constant frequency transmitting inside the bandwidth window that the LoRa signal is using. This kind of interference is filtered out easily by a LoRa receiver.

The size of a bandwidth window, rate of frequency change (also known as a spreading factor) and transmitting power further define the LoRa signal. With these factors, we can influence the range, power consumption and bit rate of a LoRa signal. For example, as seen in Figure 2.10b, by increasing the spreading factor we increase the time on air, thus giving the receiver more time to sample the signal, which leads to better sensitivity, but increases power consumption.

While LoRa defines the physical layer, LoRaWAN defines the media access control protocol for wide area networks, which are built on top of LoRa [26]. Its specification is open, so anyone can implement it. LoRaWAN takes care of communication between end-devices and gateways and manages communication frequency bands, data rates and transmitting power.

LoRaWAN has a star of stars topology [26]. Devices deployed in the field transmit messages on frequency bands that differ from region to region. Messages are received by gateways which relay them to the network server. The network server displays relayed messages, decodes them and sends them to various applications. If the same

message is heard by several gateways, the server drops all duplicates. The server also decides which gateway will send a downlink message to a specific device.

Because LoRaWAN operates on an unlicensed ISM band, anyone can setup up their network without any licensing fees. For some use cases, a single gateway with an internet connection is enough to provide coverage to a large number of devices.

## 2.6 Thermal cameras

Thermal cameras are transducers that convert infrared (IR) radiation into electrical signals, which can be used to form a thermal image. A comparison between a normal and a thermal image can be seen in Figure 2.11. IR is an electromagnetic (EM) radiation, and covers part of the EM spectrum that is invisible to the human eye. The IR spectrum covers wavelengths from  $780 \mu\text{m}$  to  $1 \text{ mm}$ , but only a small part of that spectrum is used for IR imaging (from  $0.9 \mu\text{m}$  to  $14 \mu\text{m}$ ) [29]. We can classify IR cameras broadly into two categories: photon detectors or thermal detectors [29]. Photon detectors convert absorbed EM radiation directly into electric signals by the change of concentration of free charge carriers [29]. Thermal detectors convert absorbed EM radiation into thermal energy, raising the detector temperature [29]. The change of the detector's temperature is then converted into an electrical signal. Since photon detectors are expensive, large and therefore unsuitable for our use case, we will not describe them in greater detail.



Figure 2.11: Comparison between a normal image and thermal image. Image source: Arribada Initiative [30]

Common examples of thermal detectors are thermopiles and microbolometers. Thermopiles are composed of several thermocouples. Thermocouples consists of two different metals joined at one end, which is known as the hot junction. The other two ends of the metals are known as cold junctions. When there is a temperature difference between the hot and cold junctions, a voltage proportional to that difference is generated on the open ends of the metals. To increase voltage responsivity, several thermocouples are connected in series to form a thermopile [29]. Thermopiles have lower responsivity when compared to microbolometers, but they do not require temperature stabilisation [29].

Microbolometers can be found in most IR cameras today [29]. They are sensitive to IR wavelengths of 8 to 14  $\mu\text{m}$ , which is a part of the long wave infrared region (LWIR) [29]. Measuring part of a microbolometer is known as Focal Point Array (FPA) (Figure 2.12a). FPA consists of IR thermal detectors, bolometers (Figure 2.12b), that convert IR radiation into an electric signal. Each bolometer consists of an absorber material connected to a Readout Integrated Circuit (ROIC) over thermally insulated, but electrically conductive legs [31].

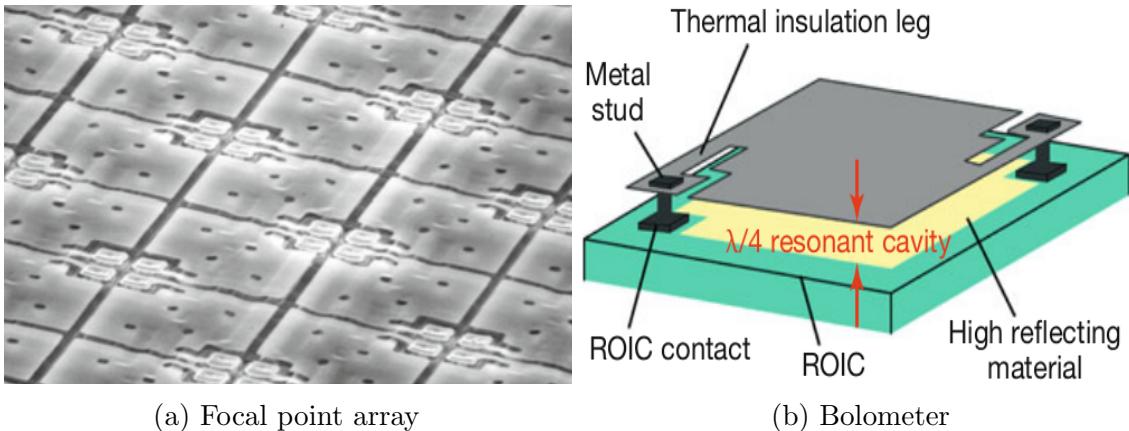


Figure 2.12: (a) Focal point array under electronic microscope. (b) Bolometer with  $\lambda/4$  resonant cavity. Image source: Vollmer, Möllmann [29]

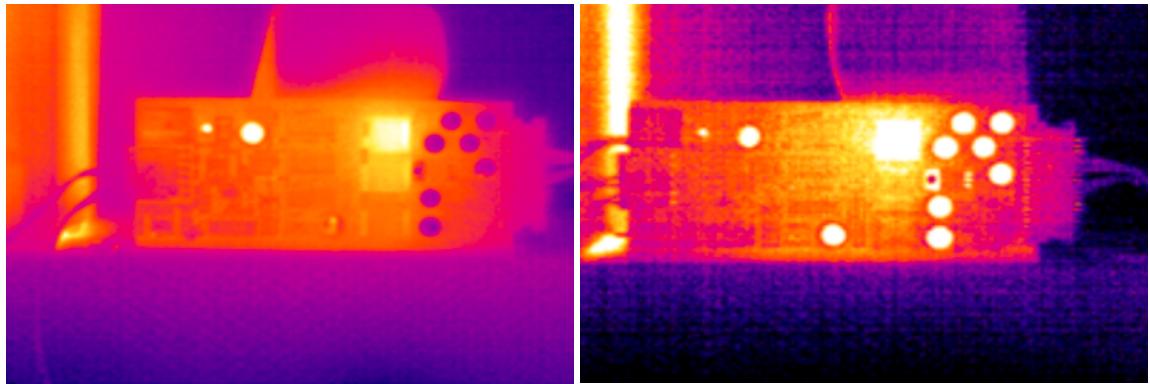
Absorber material is either made out of metals such as gold, platinum, titanium, or, more commonly, out of semiconductors such as vanadium-oxide (VO<sub>x</sub>) [31]. The important property of absorber materials is that electrical resistance changes

proportionally with the material's temperature [29]. When IR radiation hits absorber material, it is converted into thermal energy, which raises the absorber's temperature, thus changing its resistance. To detect the change in resistance, ROIC applies a steady-state bias current to the absorber material, while measuring voltage over the conductive legs [29].

When deciding between different types of thermal cameras we are often comparing them in terms of the cost, size and image resolution. One important property that also has to be taken into account is temperature sensitivity, also known as Noise Equivalent Temperature Difference (NETD). NETD is measured in mK, and tells us the minimum temperature difference that can still be detected by a thermal camera. In microbolometers, NETD is proportional to the thermal conductance of the absorber material, among other factors [29]. The thermal conductance of bolometers is minimised by enclosing FPA into the vacuum chamber, thus excluding thermal convection and conduction due to the surrounding gases. The only means of heat transfer that remain are radiant heat exchange (highly reflective material below the absorber is minimising its radiative losses), and conductive heat exchange through the supportive legs. NETD also depends on the temperature inside the camera, as higher ambient temperatures can raise the internal temperature, thus increasing the NETD and noise present in the thermal image. Today's thermopiles can achieve NETD of 100 mK, microbolometers 45 mK, while photon detectors can have NETD of 10 mK. Although tens of mK does not seem a lot, we can see in Figure 2.13 what a difference of 20 mK means for image resolution and noise.

### 2.6.1 Choosing the thermal camera

The choice of thermal camera was made by the Arribada Initiative [30]. They tested several different thermopiles and microbolometers while searching for the desired properties. The camera had to be relatively inexpensive, and small enough so that it could be integrated into a relatively small housing. The main property that they searched for was that elephants could be recognised easily from thermal images. That



(a) NETD is 60 mK

(b) NETD is 80 mK

Figure 2.13: Comparison of images of the same object taken with cameras with different NETD values. Low NETD values are more appropriate for object recognition. Image source: MoviTherm [32]

meant that the camera needed to have decent resolution and low NETD. Cameras were tested in Whipsnade Zoo and the Yorkshire Wildlife Park where images of elephants and polar bears could be made.

They tested two thermopile cameras (Heimann 80x64, MELEXIS MLX90640) and two microbolometer cameras (ULIS Micro80 Gen2, FLIR Lepton 2.5). Although thermopile cameras were cheaper than microbolometer cameras, the quality of images they produced was inferior, as can be seen in Figure 2.14.

The MELEXIS MLX90640 camera had a resolution of 32 x 24 pixels and NETD of 100 mK, while the Heimann camera had a resolution of 80 x 64 pixels and NETD of 400 mK. It was concluded that images taken by either one of the thermopile cameras could not be used for object recognition, merely only if the object was present or not [30].

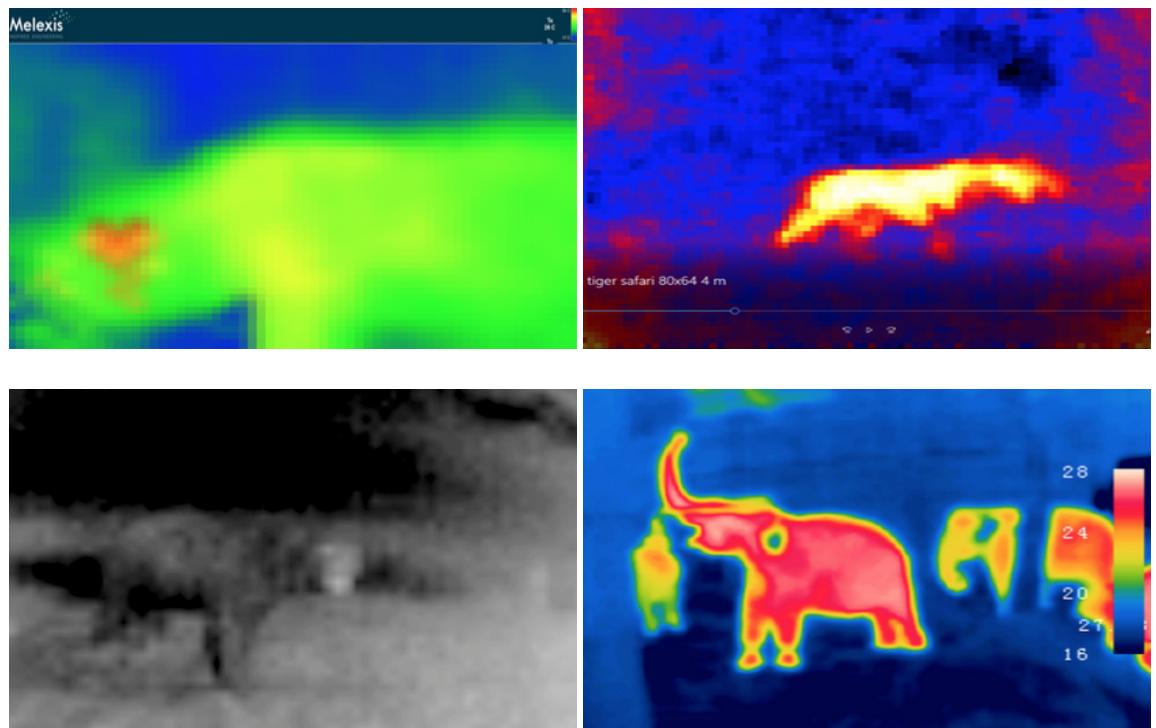


Figure 2.14: Comparison of image quality made by different thermal cameras, MELEXIS MLX90640 (top left), Heimann 80x64 (top right), ULIS Micro80 Gen2 (bottom left) and FLIR Lepton 2.5 (bottom right). Image source: Arribada Initiative [30]

Microbolometers produced better results. Both the Ulis Micro80 and FLIR Lepton had a similar resolution, 80 x 80 and 80 x 60 respectively, but the Ulis Micro80 had two times bigger NETD compared to the FLIR Lepton camera, 100 mK and 50 mK, respectively. Images produced by the FLIR Lepton were much cleaner, so it was chosen as an appropriate camera for the task.

It is important to note that the FLIR Lepton, like all microbolometers, requires frequent calibration to function properly. In temperature non-stabilised cameras small temperature drifts can have a major impact on image quality [29]. Calibration is done either by internal algorithms of the camera or by exposing the camera to a uniform thermal scene. The FLIR Lepton camera comes with a shutter, which acts as a uniform thermal signal and enables regular calibration. Calibration in the FLIR Lepton is by default automatic, triggering at startup, and every 3 minutes afterwards or if camera temperature drifts by more than 1.5 °C.

The FLIR Lepton camera comes in two versions, 2.5 and 3.5. Both cameras function the same and have exactly the specifications, they only differ in resolution: the 3.5 has a resolution of 120 x 160, while the 2.5 has 60 x 80. Both were used in the process of image collection.

## 3 Neural Network model design

In this chapter we describe the design of a Convolutional Neural Network that can process thermal images and predict what object they contain. The workflow that we followed is largely a combination of the workflows shown in Figures 2.1 and 2.9.

We first had to set concrete objectives, while keeping in consideration various constraints. The tools and development environment that were used in the process are then described. The methods of dataset creation are described afterwards: first, the dataset that was created by Arribada Initiative, then the dataset provided by us.

We then explored both datasets, analysed different class representations, and decided if they were appropriate for accomplishing the objectives that we set earlier.

In the image preprocessing phase we imported images and connected them with metadata that was parsed from the Excel database. We analysed the dataset, split it into different sets, and applied image correction procedures. We then decided on a rough CNN architecture with variable hyperparameters, and ran a random search algorithm, which searched for the best performing models based on accuracy.

We finish this chapter by going through the same design process again, but this time using tools provided by Edge Impulse.

### 3.1 Model objectives

The accuracy of our early warning system should be equal or similar to the one of the human observers, no matter if it is operating in daytime or night time. Although the system will be placed on paths that are traversed regularly by elephants, they

are not the only possible objects that can appear on the taken thermal images. Humans and various livestock, such as goats and cows, could also be photographed. Reporting false positives should be avoided, which means that the system should not label a human or a livestock animal incorrectly as an elephant. At the same time, false negatives also need to be avoided, as an elephant could pass the system undetected. These kinds of mistakes could undermine the community's confidence in the early warning system and defeat the purpose. This means that, besides elephant detection, we should also focus on classifying humans and livestock correctly, while providing a nature/random class for all other unexpected objects or simply images of nature.

It would be beneficial if the thermal camera could take several images of the same object in a short time, thus increasing the confidence of the computed label of the object. However, this is constrained by the image processing time and the camera's field of view. The thermal camera FLIR Lepton has a horizontal field of view of 57 degrees. The closer an object passes by a thermal camera, the quicker it traverses the camera's field of view, thus giving the camera less time for capture. This problem can be solved by minimising the execution time of the ML model, or by placing the early warning system in a position that is several metres away from the expected elephant's path. As the latter option might not always be possible, we should strive to keep the whole image processing time as short as possible.

Finally, as our Neural Network has to run on a microcontroller and not on a computer or a server, we have to keep it lightweight in terms of memory. Extra model complexity that brings few percents of accuracy does not matter much if the model is too large to fit on a microcontroller or takes too long to run.

To summarise:

- We will create an image classification ML model that will be capable of processing a thermal image and sorting it into one of 4 possible categories: Elephant, human, livestock, and nature/random.
- Total image processing time should be as short as possible: we should try to

keep it under 1 second.

- The model should be small enough to fit on a microcontroller of our choice, while still leaving some space for the application code. The microcontroller of our choice (STM32F767) has 2 MB of flash memory, so the model size should be smaller than that.

### 3.2 Tools and development environment

All of the work connected with image preparation and ML model creation was done in Python 3.6, Numpy was used for image preprocessing, Pandas for Excel database manipulation, and Matplotlib for plot generation. The Neural Networks were designed in TensorFlow 2.4, using Keras high-level API and a Keras Tuner model was used for the hyperparameter search.

As training Neural Networks is a computationally demanding process, it would not be feasible to do it on a personal laptop. Amazon's Elastic Compute Cloud web service was used instead. Elastic Compute Cloud, or EC2, enables users to create an instance of a server in a cloud with a specified amount of processing power and memory. Some instances come with dedicated software modules and dedicated graphics cards for an extra boost in performance. We created an instance of a Linux server that came with TensorFlow, Numpy, and other libraries pre-installed. Interaction with servers was done one command line through the SSH protocol.

Instead of writing Python scripts and executing them through the command line, we used Juptyer Notebook. Juptyer Notebook is a web-based application that can run programs that are a mix of code, explanatory text, and computer output. Users can divide code into segments, which can be executed separately, and visual output from modules such as Matplotlib is also supported. To use Juptyer Notebook on our cloud instance, we had to install and run it. We could then access the web service simply through a web browser by writing the IP address of the server, followed by the default Juptyer Notebook server port, 8888.

### 3.3 Creating the dataset

As mentioned in Section 1.3, the major part of the thermal image dataset was provided by the Arribada Initiative [12] [13]. Images in the dataset came from two different locations: Assam, India, and ZSL Whipsnade Zoo, in the United Kingdom.

Assam served as a testing ground. The Arribada team positioned two camera traps at two locations that overlook paths commonly used by elephants. Cameras were built out of Raspberry Pi, a Passive Infrared Sensor (PIR) sensor, an FLIR Lepton 2.5 camera, and batteries, all of which were enclosed in a plastic housing. The insides of the camera and an example of a deployed camera can be seen in Figure 3.1.



Figure 3.1: Camera trap used in Assam, India. Image source: Arribada Initiative [13]

The PIR sensor functioned as a photo trigger: whenever an object passed in front of it, the camera made an image. This setup provided Arribada with elephant images in real-life scenarios, however, they could not capture elephants in a variety of different conditions. It is important to create an image dataset, where the object can be seen in different orientations, distances, angles, and temperature conditions. Models that were trained on diverse datasets end up being much more robust and, therefore, perform better on never before seen image data, when deployed in real life.

This was accomplished in ZSL Whipsnade Zoo, where they took many images of elephants in a variety of different conditions [33]. With elephants in the enclosure, researchers could move cameras around and get images that were needed. The PIR sensor trigger approach was dropped in favour of a 5 second time-lapse trigger. Two cameras were used again, although, one of them now used an FLIR Lepton 3.5

camera with better resolution.

Images of elephants that came from both locations can be seen in Figure 3.2.

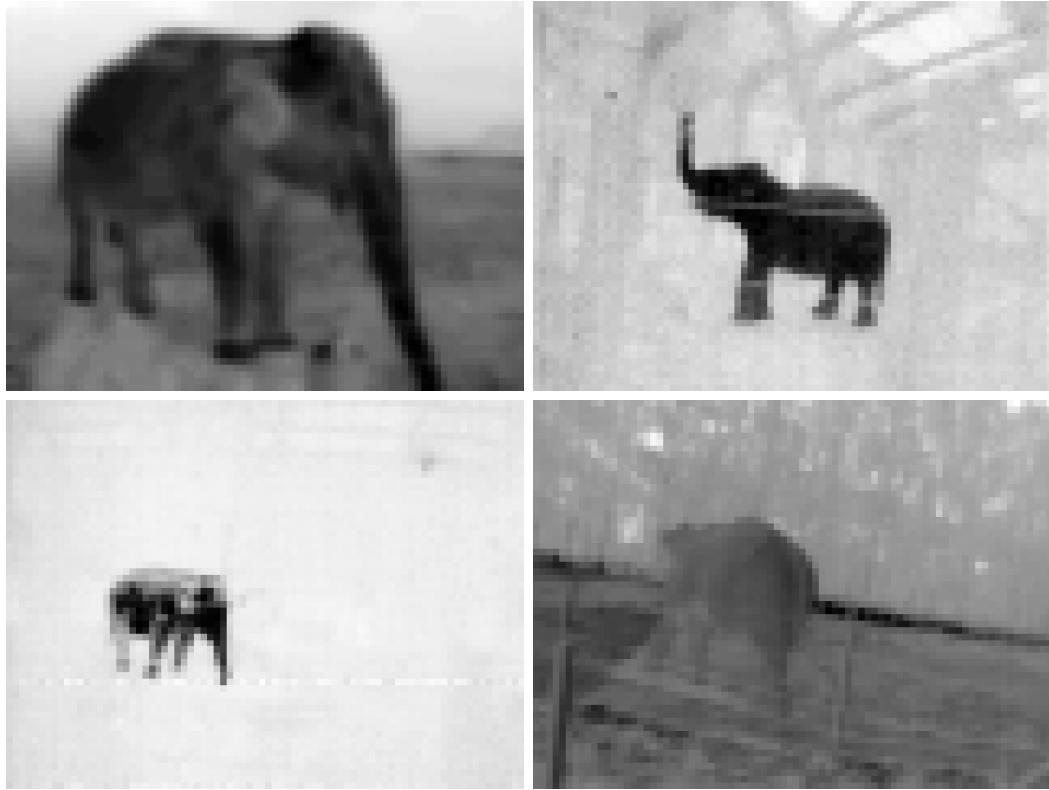


Figure 3.2: Thermal images of elephants from the dataset.

A small part of the thermal image dataset was provided by us. This was done because the number of images of cows was low compared to the number of human and elephant images, and because we also did not have any images that could be used for the nature/random class. We wanted to gather images as quickly and efficiently as possible, so we built a prototype camera made out of an FLIR Lepton 2.5 breakout board, a Raspberry Pi Zero, and a power bank. We used an open-source library [34] for the FLIR Lepton module, which used a simple C program to take a single image with a thermal camera and save it to a Raspberry Pi. The image of the setup can be seen in Figure 3.3.

We wrote a simple Python script that executed the C program every time we pressed the trigger push-button. An additional shutdown button was added to call the Raspberry Pi shutdown routine, as removing power from it forcibly would corrupt

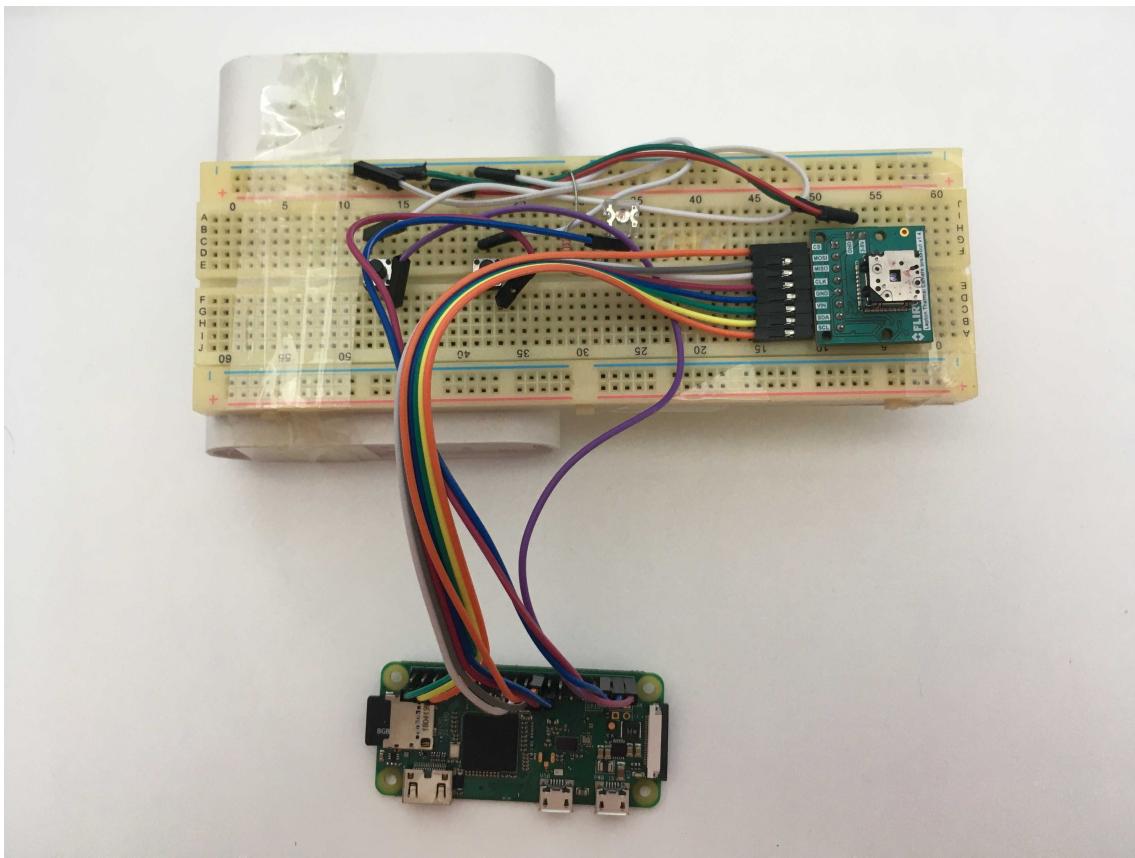


Figure 3.3: Camera setup used for taking thermal images with FLIR Lepton 2.5.

freshly taken thermal images on the Raspberry Pi's SD card.

With this setup, we made 365 images of cows in varying conditions, 308 images of nature, and 124 images of humans that were made on the go. We then sorted the images manually into appropriate folders and added them to the dataset.

### 3.4 Exploring the dataset

A thermal image dataset created by Arribada was given to us in the form of a Google Drive folder, which we downloaded to our computer.

After examining the folder, we came to several conclusions.

1. We saw that the primary focus of the Arribada team was to build an object localization model, not an image classification model. In object localisation, the Neural Network draws bounding boxes around objects that it recognises and assigns them labels, while the image classification model only labels the image

as a whole. Object localization produces a bigger and more complex model than image classification, and it is unsuitable for running on a microcontroller. All major work that was done by the Arribada team was contained in one folder where each image had an accompanying text file of the same name. Text files were produced by a DeepLabel software, which is used for preparing images for training object localisation models. Each line in a text file described the location of the bounding box and its label. This dataset format was not suitable for us, as many images contained more bounding boxes, which would be troublesome to sort into a distinct label.

We later saw that there were a few folders with names such as "Human", "Single Elephant", "Multiple Separate Elephants", "Multiple obstructing Elephants", "Cows", "Goats" and so on, which contained sorted images that we could use. We merged all folders with elephant pictures into one folder, as we did not care if the model can differentiate how many elephants are on a taken image, as we only wanted to know if there were any elephants on it or not.

2. We found out that all images were documented in a large Excel database. For each image, there was a row in a database that connected the image file name with the information on where the image was taken and with what sensor. This enabled us to generate the graph seen in Figure 3.4.

We used a total of 13,667 images from the thermal image dataset: almost 88 % of them were made in Whipsnade Zoo, while the rest of them were made in Assam. All images from Assam were made with the FLIR Lepton 2.5, while both cameras were used in Whipsnade zoo. However, more photos were made with the 2.5 version of the thermal camera.

3. After inspecting the folder with goat images manually, we saw that it contained mostly images of a herd of goats standing around a single elephant. This folder was usable only for object localisation ML models, where each goat could be tagged with a bounding box. In the case of an image classification model, this sort of training data is not desirable, as it would be too similar to another separate class, in our case the elephant class. We therefore dropped goat images

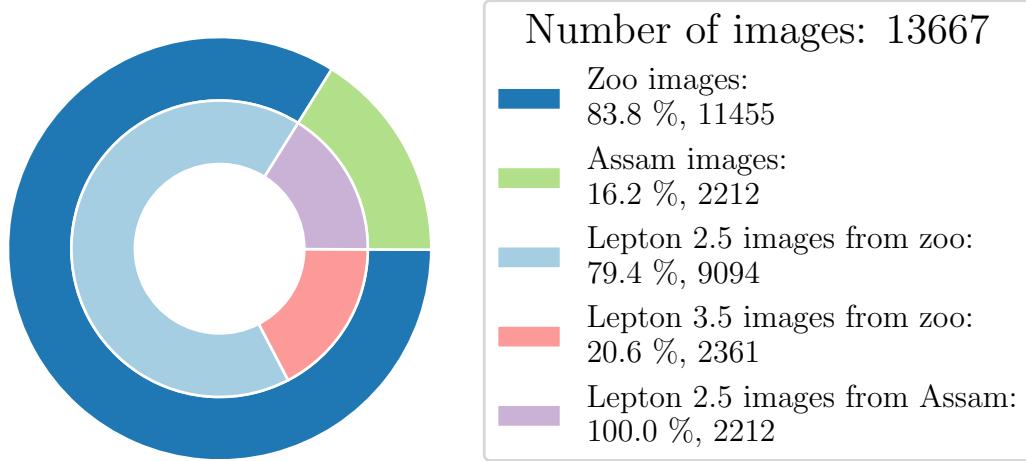


Figure 3.4: Distribution of used images from thermal dataset depending on image location and type of sensor.

out of our training data entirely. Livestock class was replaced with cow class.

4. We also realised that there was a large class imbalance, as seen in Figure 3.5 in favour of the elephant class.

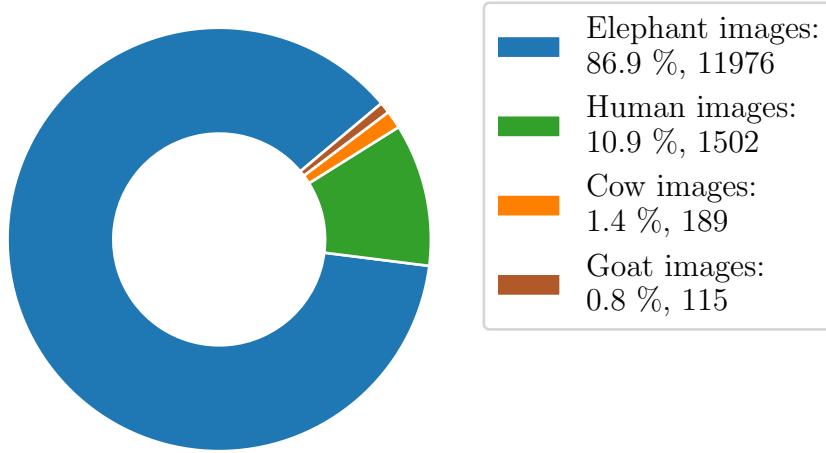


Figure 3.5: Class distribution of thermal images.

The number of elephant images was more than 4 times larger than the number of images of the all other classes combined. We solved this issue by acquiring more images of the minority class and oversampling the minority class.

### 3.5 Image preprocessing

The image preprocessing phase is a pipeline process that differs from project to project. Our process can be seen in Figure 3.6.

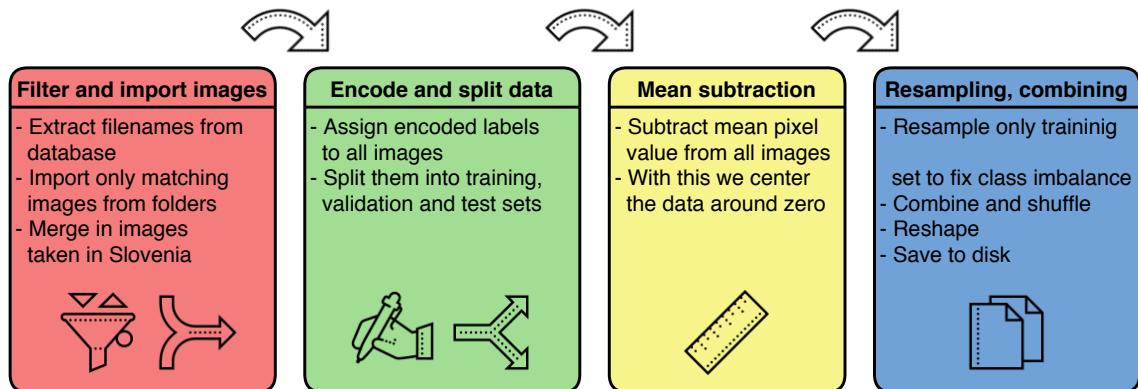


Figure 3.6: Image preprocessing pipeline. Icons source: [11]

At the start of the process, we compared the filenames of each separate folder to the list of filenames found in the Excel database. We imported only the images found in both sources, as lists were not identical, and we wanted to keep track of the different metadata information. As some images were made with two different FLIR Lepton cameras with different resolutions (60 x 80 and 120 x 160), we downscaled higher resolution images directly in the importing process. After this, we added images that were taken by us in Slovenia. At this point, we had four separate Numpy arrays, one for each class, with 3 dimensions: The first dimension stored a number of different images in that class, the second and third dimensions stored images pixel values (60 and 80 pixels respectively).

The next step was assigning labels to each image. As the output of NNs are numbers, we cannot just assign labels in strings format to data. Instead, we assigned every image a single number that represented that class, 0 for an elephant, 1 for a human, 2 for a cow, and 3 for a nature/random class. We shuffled images inside of each class and then split them into training, validation and test sets.

The training set was used for model training, while the validation set helped to choose the best model based on accuracy. The test set is normally set aside and used

only at the end, after the model is chosen, to assess how the model performs on never seen data. If we did not use the validation set and only chose the best model according to the test set, we could be overfitting<sup>1</sup>a model and we would have no accurate measure of how well our model would perform on unseen data.

At end of this step we had 4 different Python dictionaries for each class. Each dictionary had 3 key-value pairs for every training, validation, and test set, which held image data and encoded labels.

We next applied the simplest form of normalisation to all images, a mean subtraction. We calculated a two-dimensional matrix that held mean values of pixels averaged over the whole training set, which we subtracted from all images, essentially zero centring the data. This is a common preprocessing step in every ML image preprocessing pipeline, which is usually combined with standardisation<sup>2</sup>. We achieved this by resampling the human, cow, and nature/random classes. The human class was resampled 5 times, while both cow and nature/random classes were resampled 8 times. Figure 3.7 shows the distribution of training images before and after resampling.

We only resampled training sets, not validation or test sets. If we resampled everything, the model would be seeing the same image several times during testing, thus reporting incorrect accuracy in the validation and test phases.

After resampling we merged and shuffled all data, and saved it to the disk for later use.

---

<sup>1</sup>Overfitting means that the model performs well on the training data, but it does not generalise well to real-world examples [14].

<sup>2</sup>Standardisation scales the whole range of input pixel values into -1 and 1 interval. This is only needed if different input values have widely different ranges [16]. Because images that were created with the FLIR camera were all 8-bit encoded and therefore had the same range, this was not needed.

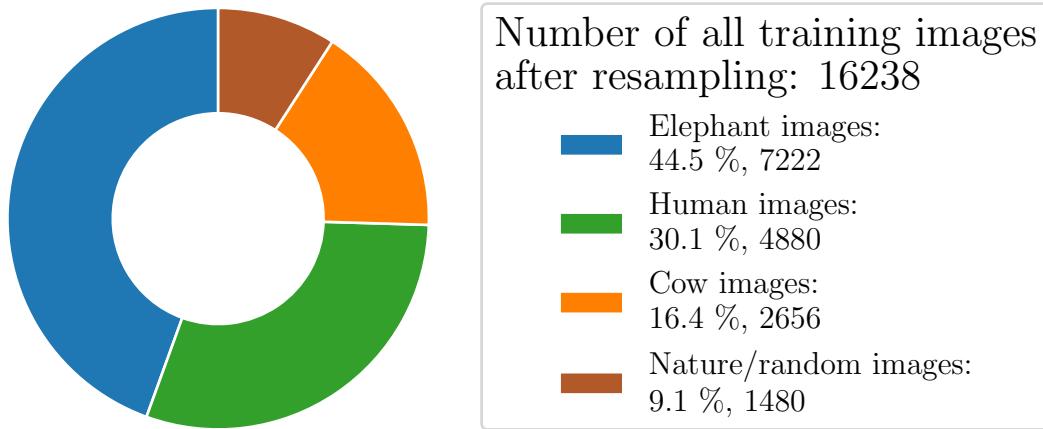
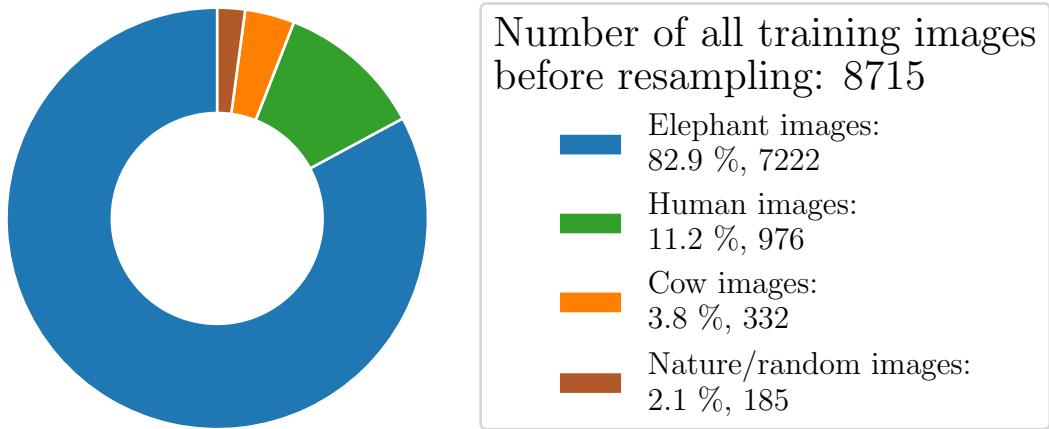


Figure 3.7: Distribution of training images before and after resampling.

### 3.6 Model creation and training

For the creation of CNN models, we used Keras Sequential API and a Keras Tuner module. Sequential API abstracted many low-level details of the model's design. While specifying layers, we only had to specify what type of layer we wanted, its size and layer-specific features. We did not have to keep track of any connections between or in layers, as this was done automatically by Keras.

For a model architecture we decided to use a simplified version of the common CNN architecture that was shown in Figure 2.4. The best way to present the model is by

inspecting the Sequential API code that creates it, and the code is shown in Figure 3.8.

The model consisted of two pairs of convolutional and max-pooling layers, followed by a final convolutional layer. For activation the function ReLu was chosen, as it is currently the most effective and popular option [16] [14]. The padding option was set to same, which meant that a spatial dimension of a volume would not change before and after a convolutional layer. Pooling layer kernel size was set to 2 x 2, with a default stride of 2.

The output volume of the last convolutional layer was flattened out into a single vector and fed into a dense layer, which was followed by a dropout layer<sup>3</sup>.

The last dense layer was a final output layer with only 4 neurons, each one representing one class. Softmax activation was used to calculate class probabilities. The model was set to use the Adam optimizer and Sparse Categorical Crossentropy loss function. Adam is an upgraded version of the gradient descent method, which adapts the learning rate automatically to decaying gradients [14]. It is generally easier to use than gradient descent, as it requires less tuning or learning rate hyperparameters. Sparse Categorical Crossentropy loss function is used when building a multi-class classifier.

The above set hyperparameters follow the general rules of thumb and serve as a good starting point when building CNNs [16]. However, hyperparameters such as the number of filters, filter size, size of a hidden dense layer, dropout rate, and learning rate, are specific to each dataset, and cannot be chosen heuristically.

To find hyperparameters that would yield the highest accuracy we used the Keras Tuner module. The exact configuration of a Keras Tuner module and comparison of trained models is presented and discussed in Section 5.1.

---

<sup>3</sup>The Dropout layer decides with probability  $p$  in each training step how many activations from the previous layer will be passed on to the next layer. It is active only during the training phase, during the testing phase activations are multiplied with  $(1 - p)$  factor to compensate. It is a very popular type of regularisation technique, which makes models more robust to the input data [14].

```

1 model = models.Sequential()
2
3 model.add(Conv2D(FilterNum1, FilterSize,
4                  activation='relu',
5                  padding="same",
6                  input_shape=(60,80, 1)))
7 model.add(MaxPooling2D((2, 2)))
8
9 model.add(Conv2D(FilterNum2, FilterSize,
10                 activation='relu',
11                 padding="same"))
12 model.add(MaxPooling2D((2, 2)))
13
14 model.add(Conv2D(FilterNum3, FilterSize,
15                 activation='relu',
16                 padding="same"))
17 model.add(Flatten())
18
19 model.add(Dense(DenseSize, activation='relu'))
20 model.add(Dropout(DropoutRate))
21 model.add(Dense(4), activation='softmax')

```

Figure 3.8: CNN architecture written in Python using Keras Sequential API.

### 3.7 Model optimisation

Keras supports saving models in h5 format, which model's architecture, values of weights, and information used while compiling the model. h5 format cannot be used directly for running trained models on mobile devices and microcontrollers: conversion to a .tflite format has to be done with the TFLite Converter tool.

The TFLite converter can convert a model in .h5 format into four differently optimised tflite models:

- **Non-quantized tflite model**, no quantization, just basic conversion from .h5 to .tflite format is done.
- **Float16 model**, weights are quantized from 32-bit to 16-bit floating-point values. The model size is split in half, and the accuracy decrease is minimal, but there is no boost in execution speed.
- **Dynamic model**, weights are quantized as 8-bit values, but operations are still done in floating-point math. Models are 4 times smaller and execution speed is faster when compared to float16 optimisation but slower than full integer optimisation.

- **Full integer model**, weights, biases, and math operations are quantized, execution speed is increased. It requires a representative dataset at conversion time.

A full integer model is an ideal choice for running models on microcontrollers, although, it should be noted that not all operations have full integer math support in TFLite Micro.

Furthermore, created tflite models need to be converted into a format that is understandable to the C++ TFlite API running on a microcontroller. This is done with the **xxd**, a Linux command-line tool that creates a hex dump out of any input file. By setting **-i** flag, the xxd tool creates a hex dump of our model, and formats it as a char array in the C programming language.

To automate the optimisation process we wrote a Python script that took the model in raw .h5 format and converted it into every possible version of the optimised tflite model. Each model was then processed with the xxd tool and pairs of .c and .h files were created, ready to be included in our application code.

### 3.8 Neural Network model design in Edge Impulse Studio

Designing a Neural Network with Edge Impulse is a much less involved process than the one we described above, because many steps of image preprocessing are automated. To start with NN design, we first had to upload our image data to the Edge Impulse Studio project. This can be done either by connecting an S3 bucket<sup>4</sup> with data with the Edge Impulse account, and transferring data to a specific project or by using the Edge Impulse command-line tools to upload image data from a computer directly to a project. We chose the S3 bucket approach, as once the data was uploaded it was trivial to transfer it to different projects.

After the data were uploaded, the rest of the NN design was done through the Edge

---

<sup>4</sup>Simple Storage Service or S3 is another service provided by Amazon, used for storing a large amount of data in the cloud.

Impulse web interface. In Figure 3.9 we can see the so-called Impulse Design tab, where we design a Neural Network by choosing different blocks. With input blocks, we tell what kind of data are we inputting, either image or time-series data, and with processing blocks we decide how are we going to extract features. Input and preprocessing blocks always return same output given the same input, while learning blocks are trainable and learn from previous experiences. For the learning block, we can choose to use a Neural Network provided by Keras, an anomaly detection algorithm, or a pre-trained model for transfer learning.

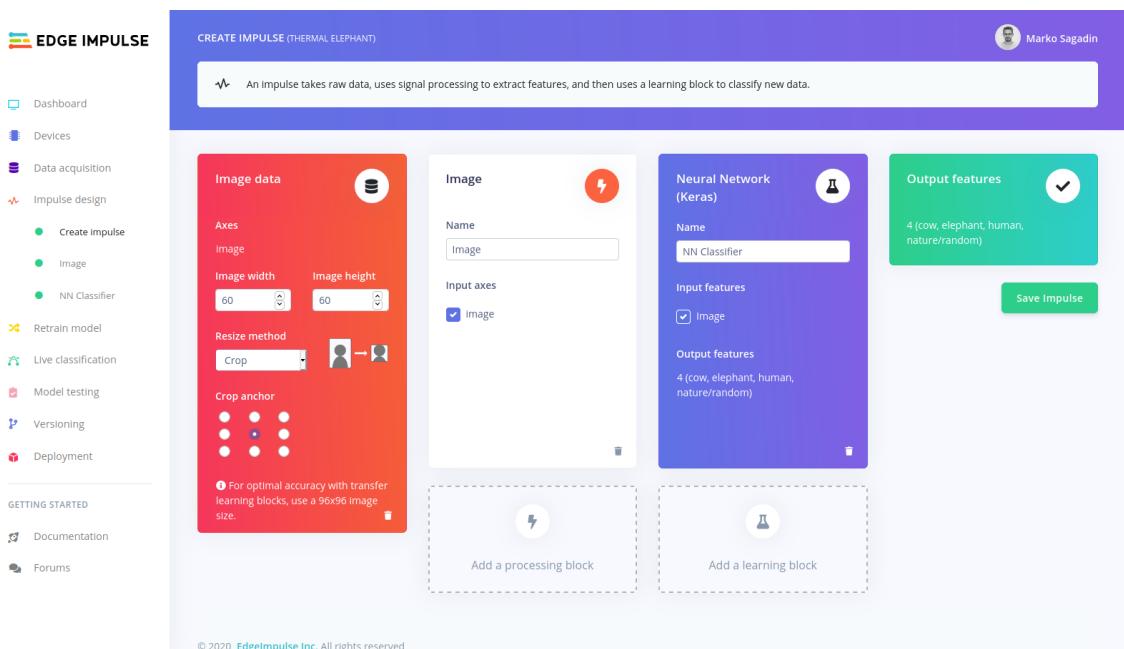


Figure 3.9: Creating a Neural Network in Edge Impulse Studio.

Since we were training with image data, we selected an image input block. As Edge Impulse did not support images of different image ratios at the time, we had to crop our images to 60 x 60 pixels. For the processing block, we selected the image processing block as this was the only possible choice. For the learning block we were using both Keras's Neural Network block and Transfer Learning block.

Both of the blocks are configurable. We could either define our network with different blocks representing layers, or switch to the text editor with Keras Sequential API code, where we could make our adjustments. Settings such as learning rate, number of epochs<sup>5</sup>, and confidence rating are also available, regardless of the option we chose.

Training of Neural Networks inside Edge Impulse Studio is done in a cloud, so as users we do not have to worry about setting up a development environment. We only have to start it and wait for it to finish. After training was done, Edge Impulse showed how well the model was performing on the validation data, how much flash and RAM would it need, and approximately how long on-device inference would take, based on the frequency and the processor of the microcontroller.

An example of the output is presented in Figure 3.10, where inferencing time is estimated for a Cortex-M4 microcontroller, running at 80 MHz. Edge Impulse also automatically converts trained model to an optimised full-integer model.

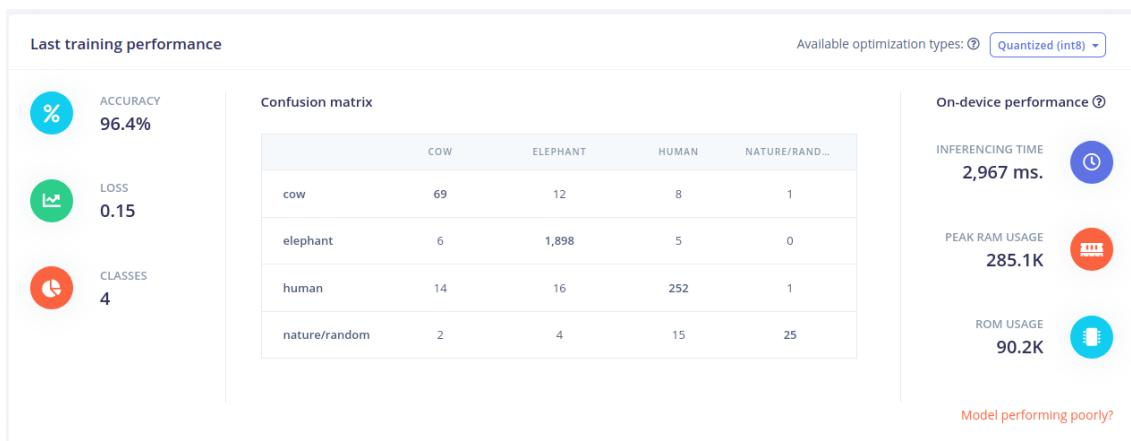


Figure 3.10: Training performance report.

The final step was deploying the trained model to the microcontroller. This step is fairly simple, as Edge Impulse provides a few example projects on their GitHub for the different platforms that it supports. As we wanted to compare the performance of the models on STM32f767ZI, we chose the Mbed platform. We copied the example Mbed project from GitHub, and in the Edge Impulse Studio we selected to generate an inferencing library with our model for the Mbed platform. We extracted the library, which consisted of C++ files, into an example project and compiled it. An example project just runs the inference continuously on one image and outputs results over the serial port. A performance comparison between this example project and our implementation is done in Section 5.2.3.

---

<sup>5</sup>Number of epochs tells us how many times the whole training set passed through model, during training process.

## 4 Design and implementation of the early warning system

General structure and tasks of an early warning system were already described in chapter 1.2. As it was mentioned before, an early warning system consists of two different components:

1. Several small embedded devices, which are deployed in the field. They capture images with a thermal camera, process them and send results over a wireless network.
2. One gateway, which is receiving results and relaying them to an application server over the Internet connection.

In this chapter, we focus on the structure and design of the deployed embedded system, both from hardware and firmware point of perspective. We also describe the construction of an application server, how received data is processed, stored and presented.

The general block diagram of an embedded system with a thermal camera is presented in Figure 4.1

The embedded system consists of two different microcontrollers with two distinct tasks, a thermal camera, PIR sensor, wireless communication module, power switch and battery.

Powerful, high-performance microcontroller and thermal camera are usually turned off, to conserve battery life. A less capable, but low-power microcontroller spends most of the time in low-power mode, waiting for a wakeup trigger from the PIR sensor. The PIR sensor points in the same direction as the thermal camera and

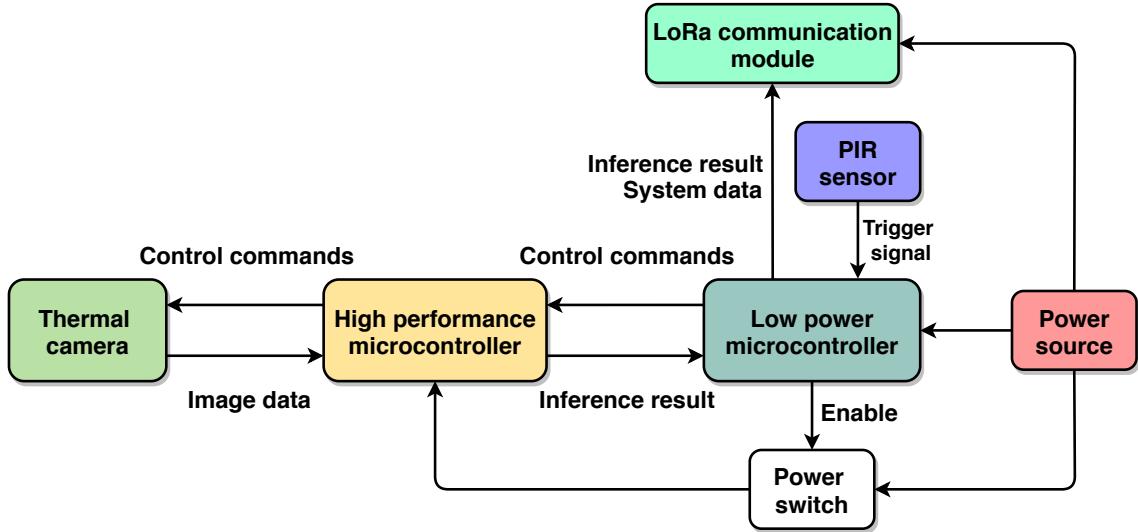


Figure 4.1: General block diagram of an embedded system

detects any IR radiation of a passing object.

If an object passes PIR's field of vision, PIR sensor produces trigger signal, which consequently wakes up a low power microcontroller. The microcontroller then enables the power supply to the high-performance microcontroller and thermal camera, and sends a command request for image capture and processing.

The thermal camera only communicates with a high-performance microcontroller, which configures it and requests image data. That data is then input into a Neural Network algorithm and probability results are then returned to a low power microcontroller. Low power microcontroller then packs the data and sends it over the radio through a wireless communication module. The power source to a high-performance microcontroller and thermal camera is then turned off to conserve power. Diagram of the described procedure can also be seen in Figure 4.2.



Figure 4.2: Diagram describing behavior of embedded early warning system

## 4.1 Hardware

In this section, we present concrete components that we used to implement the embedded part of the early warning system. The hardware version of the embedded system diagram is presented in Figure 4.3. The system consists of various development and evaluation boards.

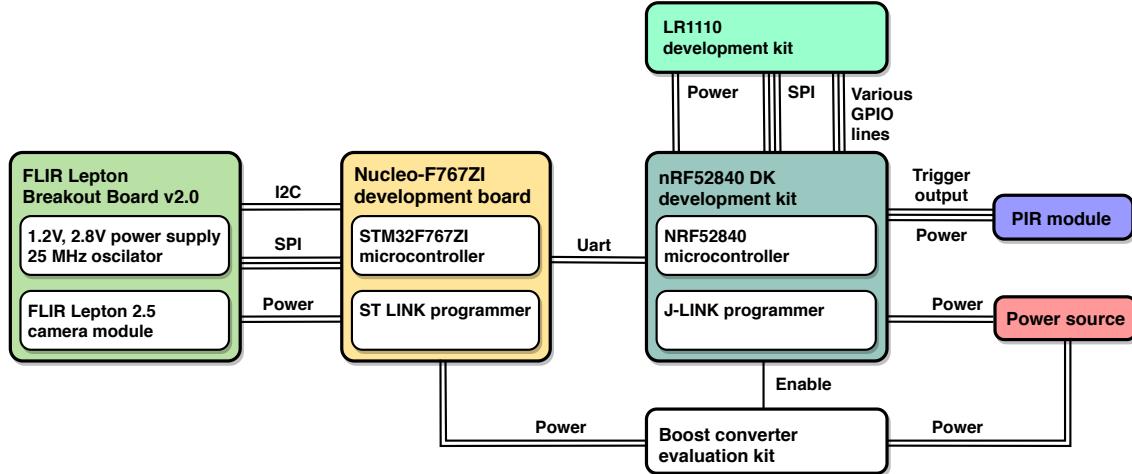


Figure 4.3: Hardware diagram of embedded early warning system

### 4.1.1 Nucleo-F767ZI

Nucleo-F767ZI (seen on Figure 4.4) is a development board made by STMicroelectronics. Board features STM32F767ZI microcontroller with Cortex-M7 core, which has 2 MB of flash, 512 kB of SRAM and can operate at clock speed of 216 MHz. It also features memory caches and a flash accelerator, which provide an extra boost in performance. It is convenient to program it, as it includes onboard ST-LINK programmer circuit.

We chose this microcontroller simply because it is one of the more powerful general purpose microcontrollers on the market. As we knew that Neural Networks are computationally expensive to compute and that models can be quite large in terms of memory, we selected it knowing that we can always scale down, if we have to.

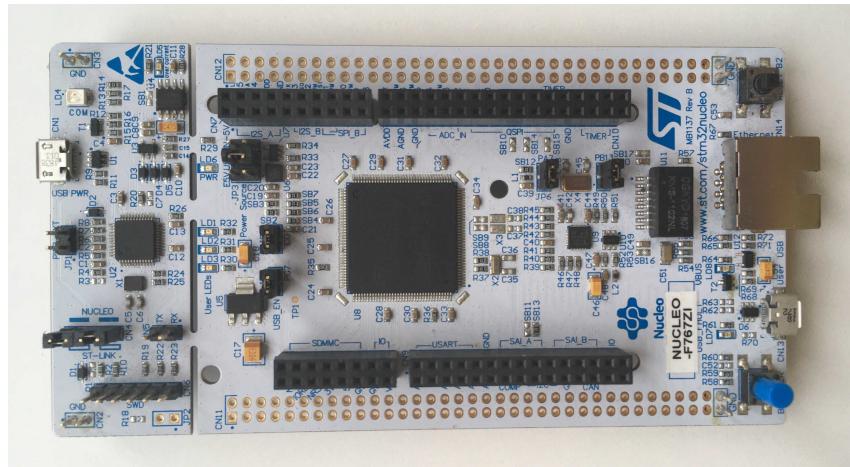


Figure 4.4: Nucleo-F767ZI development board

#### 4.1.2 nRF52840 DK

For the part of the system, which had to contain a low-power microcontroller and would control the communication module and power control for the Nucleo-F767ZI board, we decided to use the nRF52840 development kit. The development kit, made by Nordic Semiconductor, can be seen in Figure 4.5

The main logic on the board is provided by an nRF52840 microcontroller with Cortex-M4 core, which has 1 MB of flash, 256 kB of RAM and a Bluetooth 5 support. nRF52840 has a consumption of  $0.5 \mu\text{A}$  in sleep mode, which makes it ideal for our purpose.



Figure 4.5: nRF52840DK development board

#### 4.1.3 LR1110 development kit

For the role of the LoRa transceiver module, we decided to use Semtech's development kit which uses the LR1110 chip. LR1110 is a multi-functional solution as it contains LoRa transceiver, GNSS and WiFi geoposition scanning modules. Development kit seen in Figure 4.6 contains an LR1110 chip, three different antennas and its respective tuning networks. It comes in a convenient Arduino shield form factor, which means that we can directly attach it to nRF52840 DK without any jumper wires.

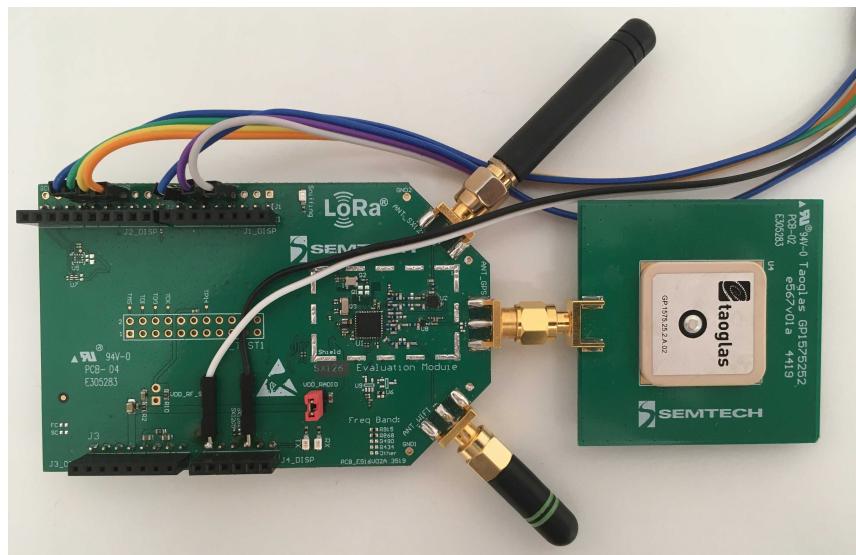


Figure 4.6: MAX17225ENT+T boost converter breakout board

#### 4.1.4 Boost converter evaluation kit

Power control of the Nucleo-F767ZI board and the FLIR camera is provided by the MAX17225ENT+T boost converter chip. Breakout board containing the chip is shown in Figure 4.7. Operating the boost converter chip is simple, its enable line can be directly connected to a microcontroller pin, driving it high will enables output and driving it low disables it. The output voltage is controlled by an external resistor.

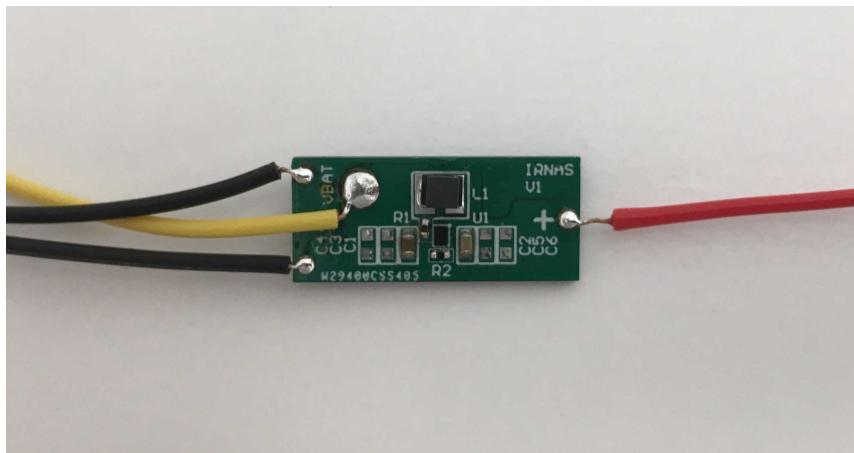


Figure 4.7: MAX17225ENT+T boost converter breakout board

#### 4.1.5 FLIR Lepton 2.5 camera module and Lepton breakout board

Section 2.6 described what kinds of thermal cameras exist and how do they work, and Section 2.6.1 described why FLIR Lepton 2.5 was chosen. However, not much was said about what sort of support circuitry FLIR camera needs and how do we make images with it.

FLIR Lepton camera is powered from two different sources, 1.2 V and 2.8 V, and requires a reference clock of 25 MHz. All of this is provided by the Lepton breakout board, which can be seen in Figure 4.8. The front side of the breakout board contains a FLIR module socket and the backside contains two voltage regulators and an oscillator. Breakout board can be powered from 3.3 to 5 V and also conveniently breaks out all communication pins in form of header pins.

FLIR Lepton module itself contains five different subsystems that work together and can be configured:

- AGC – Automatic Gain Control, affects image contrast and quality
- SYS – System Information
- VID – Video Processing Control
- OEM – Camera configuration for OEM customers
- RAD – Radiometry

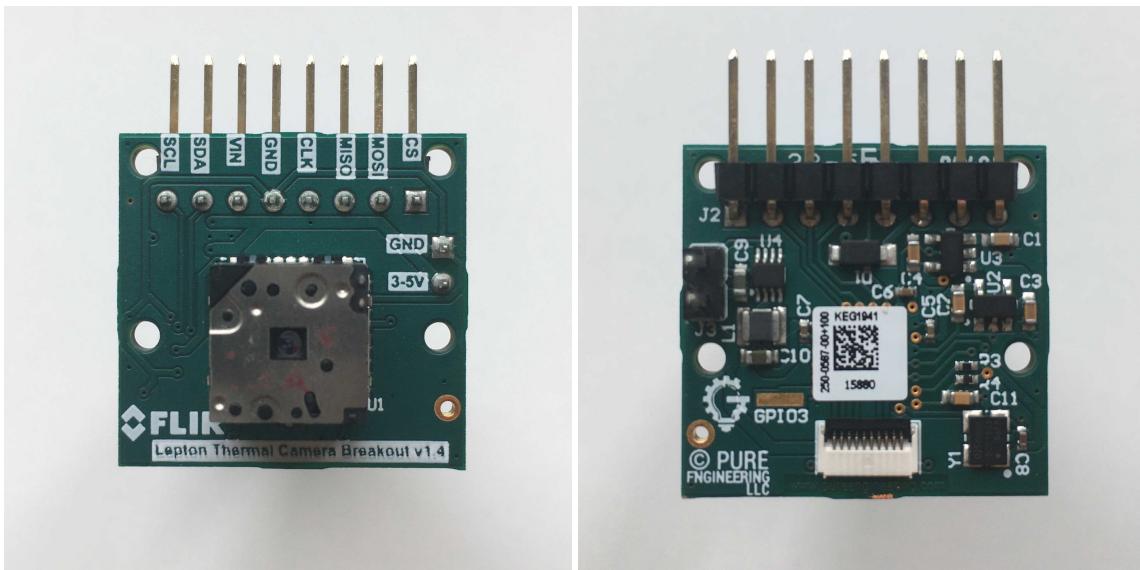


Figure 4.8: Front and back side of FLIR Lepton breakout board with thermal camera module inserted.

Task of AGC subsystem is to convert a dynamic range of an IR sensor into a compact range that is more suitable for storing and displaying images. In the case of FLIR Lepton, this is a 14-bit to 8-bit conversion. For our purposes, the AGC subsystem was turned on, as the input to our image classification model were 8-bit values.

The microcontroller communicates with FLIR camera over two interfaces: two-wire interface (TWI) is used for control of the FLIR camera and Lepton's VoSPI protocol is used for image transfer.

TWI is a variation of an I2C protocol, instead of 8 bits all transfers are 16 bits. The internal structure of the Lepton's control block can be seen in Figure 4.9. Whenever we are communicating with the FLIR camera we have to specify which subsystem are we addressing, what type of action we want to do (get, set or run), length of data and data itself.

Lepton's VoSPI protocol is used only to stream image data from the camera module to the microcontroller, which means that the MOSI line is not used. Each image fits into one VoSPI frame and each frame consists of 60 VoSPI packets. One VoSPI packet contains 2 bytes of an ID field, 2 bytes of a CRC field and 160 bytes of data<sup>1</sup>, which represents one image line.

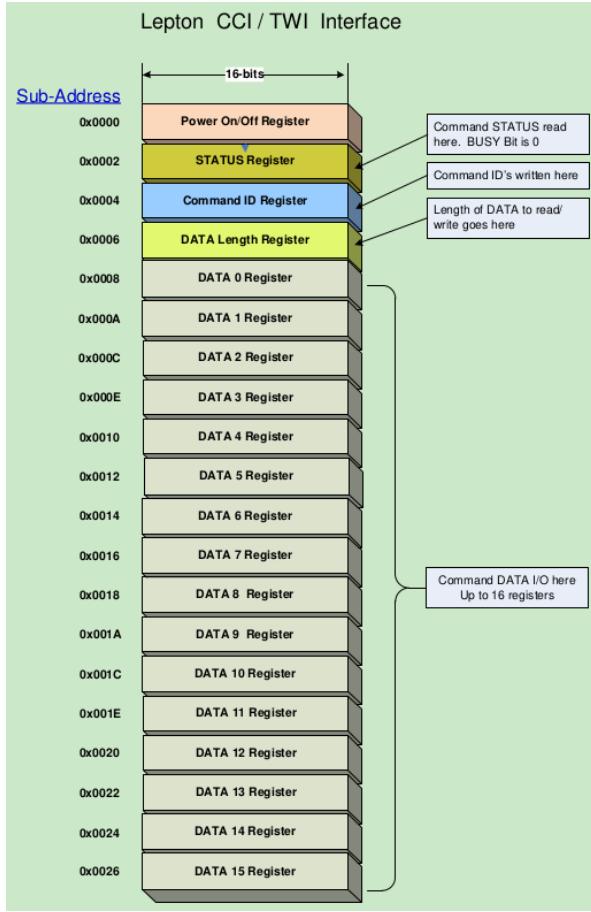


Figure 4.9: Command and control interface of FLIR Lepton camera.

The refresh rate of VoSPI frames is 27 Hz, however, only every third frame is unique from the last one. It is the job of the microcontroller to control the SPI clock speed and process each frame fast enough so that the each unique frame is not discarded.

#### 4.1.6 PIR Sensor

We used a cheap, generic PIR sensor, that can be seen in Figure 4.10. It has two potentiometers, which are used to adjust sensor's sensitivity and detection delay. PIR sensor runs on 3.3 V, which enables us to directly power it from the nRF52840 DK.

---

<sup>1</sup>Because images pixel values fit into the 14-bit range by default, it means that one-pixel value needs two bytes of data (two most significant bytes are zero). That means that each image line (80 pixels) is stored in 160 bytes. If AGC conversion is turned on, each pixel is then mapped into an 8-bit range, however, the size of one line in the VoSPI packet remains 160 bytes, 8 most significant bits are simply zeros.

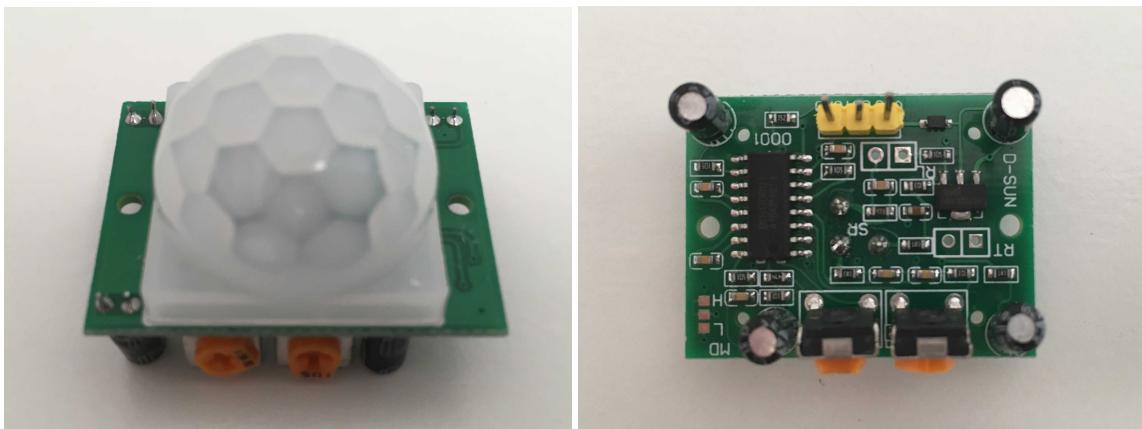


Figure 4.10: Front and back side of FLIR Lepton breakout board with thermal camera module inserted.

## 4.2 Firmware

### 4.2.1 Tools and development environment

For our firmware development we did not choose any of the integrated development environments, provided by different vendors. Instead we used terminal text editor Vim for writing and editing the code.

As we were programming two different microcontrollers, we were using different tools for each one.

#### 4.2.1.1 Development environment for STM32f767ZI

For building our firmware programs we used GNU Make, build automation that builds software according to user written *Makefiles*. To compile code we used the Arm embedded version of GNU GCC. To program binaries into our microcontroller we used OpenOCD.

For hardware abstraction library we used libopencm3, which is an open-source low-level library that supports many of Arm's Cortex-M processors cores, which can be found in a variety of microcontroller families such as ST's STM32, Toshiba's TX03, Atmel's SAM3U, NXP's LPC1000, Silabs's EFM32 and others. Libopencm3 provided us with linker files, startup routines, thinly wrapped peripheral drivers and a starting

template makefile, which served as a starting point for our project.

As libopencm3 does not provide `printf` functionality out of the box we used an excellent library by GitHub user mpaland [35].

#### 4.2.1.2 Development environment for nRF52840

To develop firmware for nRF52840 we decided to use The Zephyr OS, which is a small kernel, designed for IoT embedded systems. Besides usual RTOS functionalities such as tasks, mutexes, semaphores it also provides a common driver API for supported microcontrollers.

#### 4.2.2 Architecture design

STM32F767ZI firmware was designed to be very efficient and lean, only truly necessary parts of firmware were implemented.

As seen in Figure 4.11 we split the firmware into two hardware and application modules.

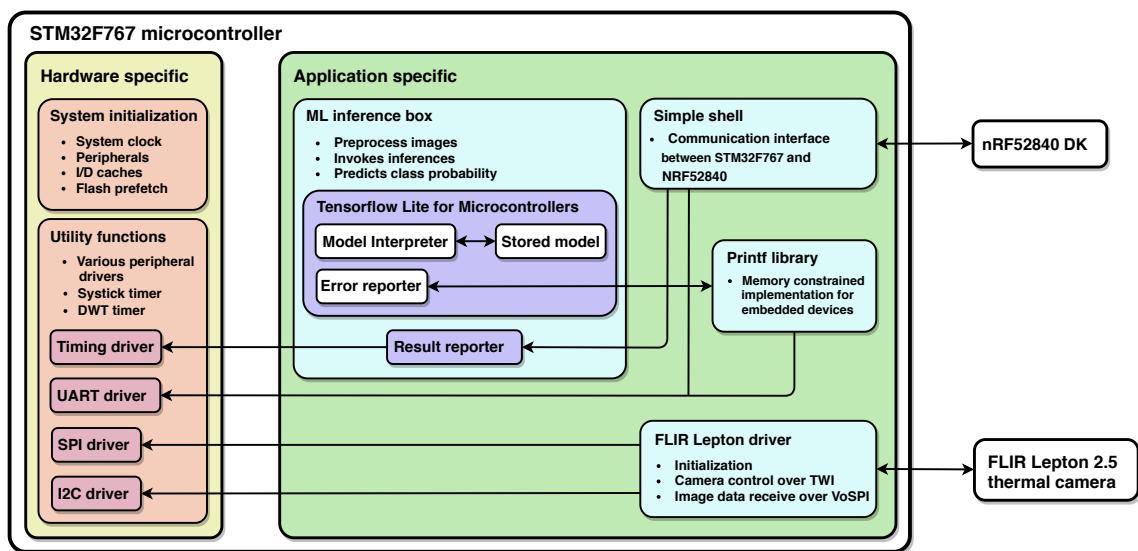


Figure 4.11: Architecture diagram of the firmware that is running on a STM32F767ZI microcontroller.

Hardware specific module is mostly using libopencm3 API to set the system clock and initialize peripherals. Small function wrappers had to be written to make use of

various peripheral drivers easier.

FLIR Lepton driver was written from scratch, as many libraries provided either by the camera manufacturer or open source communities were too complex and implemented many features that we did not require.

Thanks to TFLite Micro API, the ML inference module could be written as a simple black box. Image data goes in, predictions come out.

The architecture diagram for nRF52840 can be seen on Figure 4.12. For the nRF52840 microcontroller we did not had to write any peripheral drivers, as they were provided by Zephyr itself.

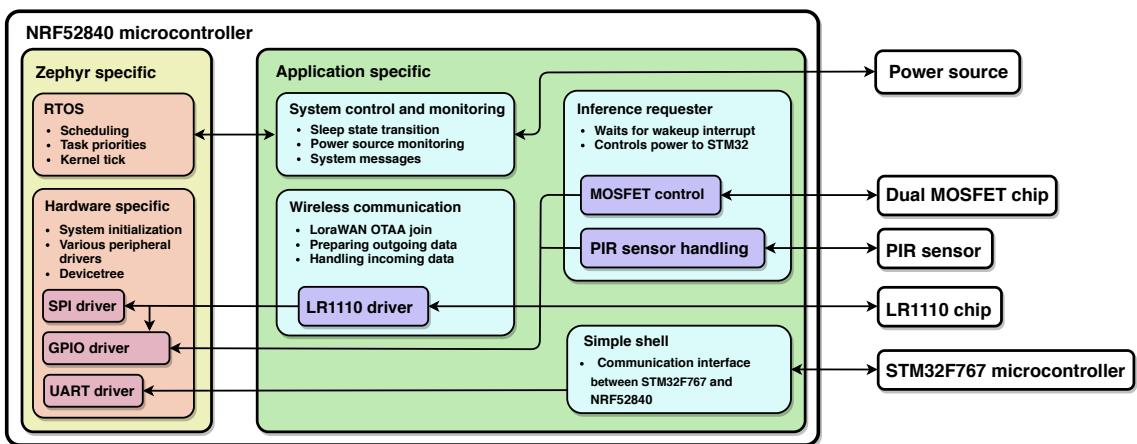


Figure 4.12: Architecture diagram of the firmware that is running on a nRF52840 microcontroller.

Utmost importance was achieving low power consumption. That meant that nRF52840 had to spend most of its time in low-power mode, only waking up for regular system checks and PIR trigger signals. When PIR trigger signal would be received, the interrupt would wake up the nRF52840, which would then enable a boost converter, thus enabling power to STM32F67 and FLIR Lepton camera.

For the communication interface, we decided to implement a simple UART shell module. nRF52840 microcontroller would act as a host and send commands to STM32F767ZI. STM32F767ZI would execute commands and send back results. Such setup enabled us to test each microcontroller separately.

Once we knew that UART communication worked correctly, we could issue commands

directly from the computer’s serial port, which enabled us to develop and test firmware for STM32F767ZI separately from nRF52840 firmware.

We also wrote a communication module, which takes care of controlling the LR1110 chip, joining the LoRaWAN network, preparing outgoing messages and sending them over the LoRaWAN network.

#### 4.2.3 MicroML and build system

A large part of this thesis was concerned with porting TFLite Micro to libopencm3, our platform of choice. To understand how this could be done, we first had to analyse how the code is built in TFLite Micro.

To compile source files and build binaries TFLite Micro uses GNU Make. The main makefile includes several platform specific makefiles, which dictate how firmware is built, and several bash scripts which download various dependencies. By providing command-line arguments users decide which example needs to be compiled and for which platform. The build system makes some assumptions about locations of the platform specific files, which in case of example projects are scattered over the whole TensorFlow GitHub repository.

We learned a useful principle while analysing the build process. Each time we would build an example for a new platform, Make would first compile all TensorFlow files, create a static library out of them, compile specific example source files and then link against the library in linking stage. If we wanted to build firmware for a different example, but for the same platform, Make would only have to compile source files of that example and link them with previously made library. As compilation of required TensorFlow files takes quite some time, this was an efficient option.

After analysing the TFLite Micro’s build system we created a list of requirements that we wanted to fulfil on our platform.

1. We wanted to keep project-specific code, libopencm3 code and TFLite Micro code separated.
2. We wanted a system, where it would be easy to change a microcontroller

specific part of a building process.

3. We wanted to reuse static library principle that we saw in TFLite Micro build process.

Covering different platforms and use cases made main TFLite Micro makefile quite complex and hard to understand. This meant that it would be hard to reuse it while porting to a new platform and we needed a different approach or reuse something else.

To solve our problem we started developing a small project that we named MicroML<sup>2</sup>. MicroML enables users to develop ML applications on microcontrollers supported by libopencm3. Project's directory structure can be seen in Figure 4.13

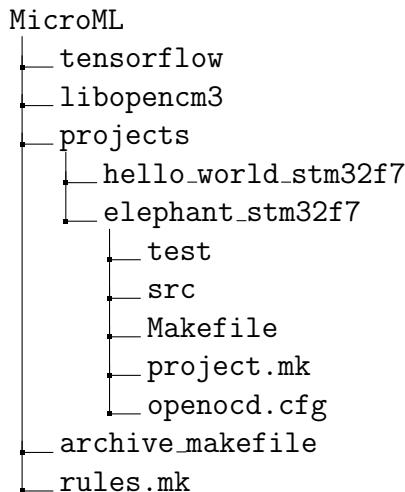


Figure 4.13: Directory structure of MicroML project.

Folders `tensorflow` and `libopencm` are directly cloned from their respective sources as Git submodules, which means that they are fixed at specific commits, usually at major release points. In folder `projects` users place all their specific projects. Besides source files, each project has to contain three specific files:

- **project.mk** - It contains information which files need to be compiled inside the project folder. In it is defined for which microcontroller the code needs to be compiled and what kind of optimisation flags should be used.

---

<sup>2</sup>Project is open-source and publicly available on GitHub [36].

- **openocd.cfg** - Configuration file that tells OpenOCD which programmer interface needs to be used to flash a microcontroller and the location of the binary file that needs to be flashed.
- **Makefile** - Project's makefile that gathers source files inside the project folder. It makes it possible to call `make` directly from projects directory, which eases the development process. It does not specify any building rules, those are specified in included `rules.mk` file in the root directory of the project.

Some initial commands need to be executed when the project is cloned from the GitHub for the first time. Figure 4.14 represents the complete build process.

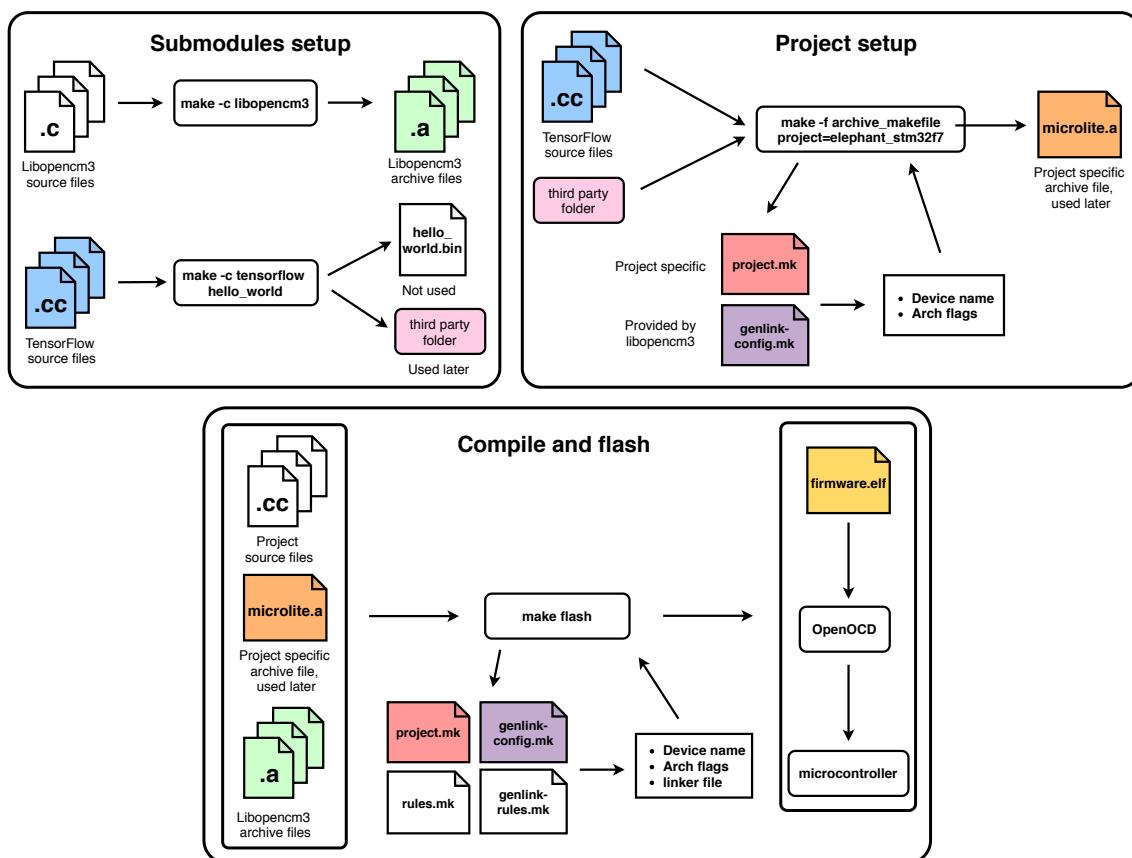


Figure 4.14: Build system of MicroML project.

In *submodules setup* stage we first compile both of the submodules, this step requires two makefiles that are already provided by each submodule. Compiling libopencm3 creates a group of archive files (static libraries), which contain all platform specific code. Compiling a TensorFlow Hello World example does not produce any archive

files that we would need, however, it does execute several scripts which download several different third party files. TFLite Micro library depends on these files, which means that MicroML does as well. *Submodules setup* stage only has to be executed once.

Whenever we start with a new project that will use ML algorithms, we need to go through *project setup* stage. From main directory, we call `make` command with `archive_makefile` and define `PROJECT` variable with the name of our project. `Archive_makefile` looks into `project.mk` and extracts `DEVICE` variable. Libopencm3's `genlink-config.py` script then determines which microcontroller specific compilation flags<sup>3</sup>are needed. All needed TensorFlow source files and third party files are then compiled and a project specific `microlite.a` archive file is created in our project's folder.

*Compile and flash* stage is then continuously executed during the development process. By calling `make flash` directly in our project folder we compile all project files, `microlite.a` and libopencm3 archive files that were created early. Libopencm3 helper scripts (`genlink-config.mk` and `genlink-rules.mk`) provide us with the microcontroller specific flags and a linker script. After compilation a `firmware.elf` is created, Make then automatically calls OpenOCD, which flashes a microcontroller.

As flashing a big binary to a microcontroller takes a long time, we also created a similar setup for testing inference directly on the host machine. That way we could test ML specific routines fast and quickly remove any mistakes found on the way.

#### 4.2.4 Running inference on a microcontroller

TFLite Micro API is fairly simple to use and general enough that it can be copied from project to project without many modifications. Figure 4.15 shows a simplified inference code example, copied from our project.

---

<sup>3</sup>For example, to compile firmware for STM32F767ZI we need flags `-mcpu=cortex-m7`, `-mthumb`, `-mfloat-abi=hard` and `-mfpu=fpv5-sp-d16`.They tell GCC compiler that we are compiling for a cortex-m7 processor, that we want to use thumb instruction set and that we want to use hardware floating-point unit with single precision.

```

1 // An area of memory to use for input, output,
2 // and intermediate arrays.
3 const int kTensorArenaSize = 200 * 1024;
4 static uint8_t tensor_arena[kTensorArenaSize];
5
6 int main()
7 {
8     // Debug print setup
9     tflite::MicroErrorReporter micro_error_reporter;
10    tflite::ErrorReporter *error_reporter = &micro_error_reporter;
11
12    // Map the model into a usable data structure
13    const tflite::Model* model = tflite::GetModel(full_quant_tflite);
14
15    // Pull in needed operations
16    static tflite::MicroMutableOpResolver<8> micro_op_resolver;
17    micro_op_resolver.AddConv2D();
18    micro_op_resolver.AddMaxPool2D();
19    micro_op_resolver.AddReshape();
20    micro_op_resolver.AddFullyConnected();
21    micro_op_resolver.AddSoftmax();
22    micro_op_resolver.AddDequantize();
23    micro_op_resolver.AddMul();
24    micro_op_resolver.AddAdd();
25
26    // Build an interpreter to run the model with.
27    static tflite::MicroInterpreter interpreter(model,
28                                                micro_op_resolver,
29                                                tensor_arena,
30                                                kTensorArenaSize,
31                                                error_reporter);
32
33    // Allocate memory from the tensor_arena
34    interpreter->AllocateTensors();
35
36    // Get information about the memory area
37    // to use for the model's input.
38    TfLiteTensor* input = interpreter->input(0);
39    TfLiteTensor* output = interpreter->output(0);
40
41    // Load data from image array
42    for (int i = 0; i < input->bytes; ++i) {
43        input->data.int8[i] = image_array[i];
44    }
45
46    // Run the model on this input and time it
47    uint32_t start = dwt_read_cycle_counter();
48    interpreter->Invoke();
49    uint32_t end = dwt_read_cycle_counter();
50
51    // Print probabilités and time elapsed
52    print_result(error_reporter, output, dwt_to_ms(end-start));
53}

```

Figure 4.15: Example of TensorFlow Lite inference code in C++.

As a first step, we need to define the size of `tensor_arena` array, which holds memory of input, output, and intermediate arrays. Exact size of `tensor_arena` is determined by trial and error: we set it to some big value and then decrease it in steps, until the code does not work any more.

In lines 9 and 10 we create an instance of `ErrorReporter` object. This object serves as a thin wrapper around platform specific `printf` implementation. If some part of TensorFlow code crashes, `ErrorReporter` notifies us what went wrong.

In the line 13 we pull in our ML model in hex dump format that we created with `xxd`. `Full_quant_model` is defined in a different file, not seen in this example.

In lines 16 to 24 we create an operation resolver. One way to do it is to specify each required operation specifically (which is done in the example) or simply pull in all operations. Latter approach is not recommended, as it results in large binary size. To find out exactly which operations were required we used online tool Netron [37], which showed us a deconstructed view of a trained model.

In lines 27 and 33 we create an `MicroInterpreter` instance and allocate memory to it that we specified with `tensor_arena` earlier. Lines 37 and 38 assign input and output of the interpreter to the new `TfLiteTensor` variables. This step enables us to do two things. Firstly, variables `input` and `output` now point to information about data format: We can find out how many dimensions are needed, what is the size of those dimensions and what is expected type of variable (`uint8_t`, `int8_t`, `float...`). In tests that we were running on the laptop, we tested exactly for these values to confirm that the model works as expected. Secondly, we now have a way to directly feed data into `input`, this is done in for loop on line 41. One of `TfLiteTensor` members is a union variable `data` which contains variables of all possible types. This type of structure enables us to load input with any kind of data, in our case `int8`.

In line 47 we finally invoke interpreter and run inference on input data. The whole expression is surrounded by the timing functions, which are used to keep track of time spent computing inference.

We finally call `print_results`, written by us, where we pass `error_reporter` for printing, `output` for extracting computed probabilities and elapsed time.

After the initial setup, we can load data, call `invoke`, and print results as many times we want.

### 4.3 Server-side components and software

In this section, we describe possible server-side construction of various frameworks which enable us to receive LoRaWAN messages, parse them, store them in a database and visualise them. We did not implement this specific setup as it was not required for testing purposes, however, at IRNAS we use this setup constantly for our IoT products and implementation of a such system would be trivial.

The system that we use consists of different tools, each one with a distinct task. These tools are The Things Network (TTN), Node-RED, InfluxDB and Grafana. The flow of information and tasks for each tool is presented in Figure 4.16.

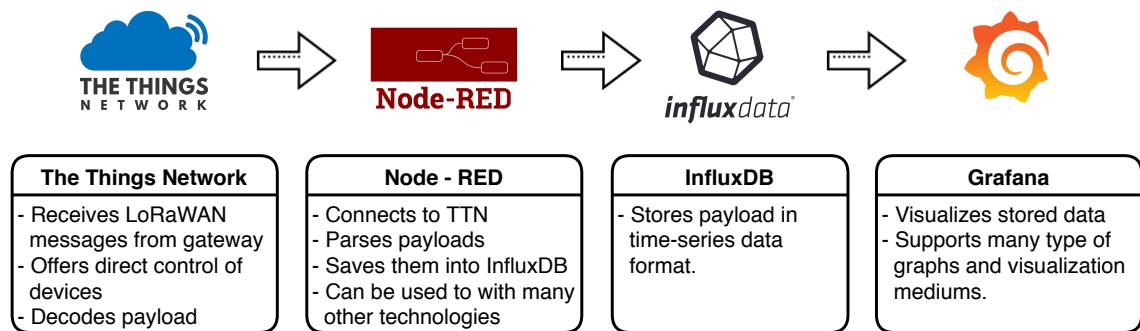


Figure 4.16: Server side flow of information. Icons source: [11]

TTN is responsible for routing packets that are captured by a gateway to the application server. Since it is open-source and free, anyone can register their gateway device into the network and thus helping to extend it. TTN is web-based, so we can see payload messages directly in the browser. Since data is usually encoded in binary format, we can provide a decoder-script written in JavaScript and TTN will automatically pass each message, thus decoding it.

Node-RED functions as a glue logic that parses packets and shapes them into a

format that is required by InfluxDB. Node-RED provides a browser-based flow editor, where actual programming can be done graphically. Logic is programmed by choosing different blocks called *nodes* and connecting them. This is convenient, as Node-RED provides different nodes for communicating with different technologies, such as MQTT, HTTP requests, emails, Twitter accounts and others. In our use case, we needed to use nodes seen in Figure 4.17. Node *Elephant Gateway* is connected to a specific application on TTN, which is used for the collection of packets from our devices in the field. Any packet that will appear in that TTN application will also appear in Node-RED. Node *Parse packet* extracts the information contained in each packet and stores it in a specific format, which is finally sent to the node *Elephant Database*.



Figure 4.17: Node-RED flow

*Elephant Database* is connected to InfluxDB, which acts as a time-series database. Any packet that is saved in it is automatically timestamped.

Data is then visualized in Grafana. Grafana is an open-source analytics and monitoring solution. Users define which database is set as a source and Grafana provides graphical controls which are at some point converted into SQL-like language, understandable to InfluxDB. Grafana provides different types of visualizations, such as graphs, gauges, heat maps, alert lists and others. In our use case, we could display information about various devices in the field, such as battery voltage, number of wakeup triggers, results of each inference and others.

Example of Grafana graph can be seen in Figure 4.18.

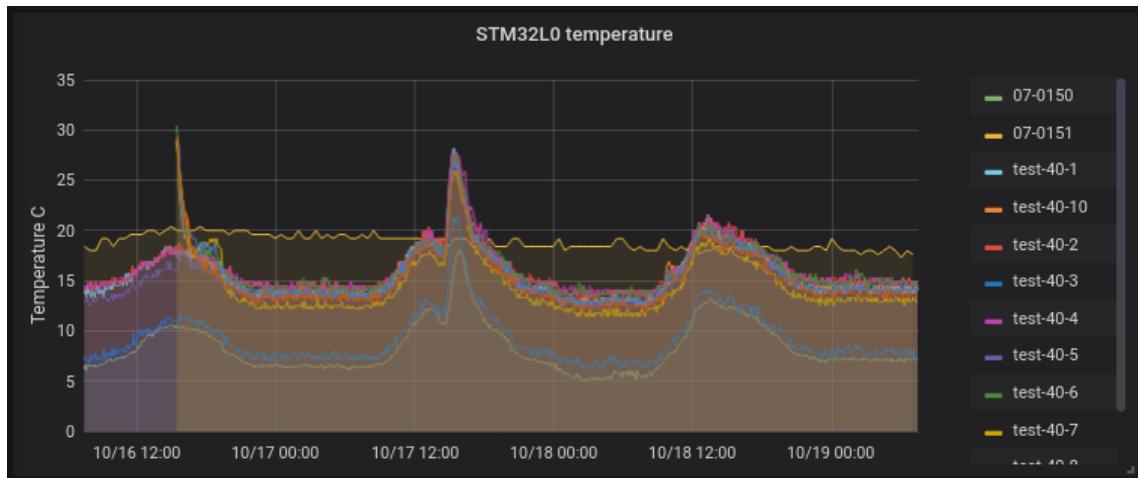


Figure 4.18: Example of Grafana graph.

One important quality of Node-RED, InfluxDB and Grafana is that they can run directly on an embedded Linux system, such as Raspberry Pi or BeagleBone, which greatly lowers the cost of hardware that is needed.

## 5 Measurements and results

### 5.1 Comparison of models

As mentioned in Section 3.6, we used Keras Tuner model to find hyperparameters that would yield the highest accuracy. Instead of hard-coding hyperparameters when building a model with Keras API, we defined a search space of possible values with `HyperParameter` class and used that as a hyperparameter.

We passed the created model to a `RandomSearch` class, with few other parameters such as batch size, number of epochs and maximum number of trials. As we started the hyperparameter search, Keras Tuner started picking a randomly set of hyperparameters, which were used to train a model. This process was repeated for a trial number of times. Used hyperparameters and achieved accuracy on the validation set for each trained model were logged in a text file for later use.

After training several different models we picked a few and compared them. Comparison of models trained in Edge Impulse Studio was also done.

#### 5.1.1 Hyperparameter search space and results analysis

General structure of CNN model was already described in Section 3.6 and in Figure 3.8. We decided to search for the following hyperparameters:

- Number of filters in all three convolutional layers (can be different for each layer)
- Size of filters in all three convolutional layers (same for all layers)
- Size of the dense layer

- Dropout rate
- Learning rate

Possible values of hyperparameters (also known as hyperparameter search space) are specified in Table 5.1.

Table 5.1: First hyperparameter search space

<b>Hyperparameter</b>	<b>Set of values</b>
<code>FilterNum1</code>	From 16 to 80, with a step of 8
<code>FilterNum2</code>	From 16 to 80, with a step of 8
<code>FilterNum3</code>	From 16 to 80, with a step of 8
<code>FilterSize</code>	3 x 3 or 3 x 4
<code>DenseSize</code>	From 16 to 96, with a step of 8
<code>DropoutRate</code>	From 0.2 to 0.5, with a step of 0.05
<code>LearningRate</code>	0.0001 or 0.0003
<b>Random search variable</b>	<b>value</b>
<code>EPOCHS</code>	25
<code>BATCH_SIZE</code>	100
<code>MAX_TRIALS</code>	300

Search space of `FilterNumX`, `DenseSize` and `DropoutRate` hyperparameters were chosen based on initial training tests conducted on a thermal image dataset and other various models that were trained on similar data. Value of `FilterSize` is usually 3 x 3, however, most of example ML projects that we could find on the Internet were training on image data of the same dimensions. We wanted to test how would a filter with the same ratio of dimensions as image data (3 x 4 and 60 x 80 respectively) perform. Hyperparameter `learning_rate` was chosen heuristically, we saw that 10 times higher values, such as 0.001 or 0.003, would leave model's accuracy stuck at suboptimal optima, from where it could not be improved anymore.

We also had to set 3 variables that directly affected how long will random search last. From initial tests we saw that models usually reached maximum possible accuracy around 20<sup>th</sup> epoch, to give some headroom we set the number of epochs to 25. We kept batch size relatively small, at 100, which meant that weights would get updated

regularly. Hyperparameter `MAX_TRIALS` had the biggest impact on the training time, we set it to 300.

The training lasted for about 12 hours. After it was done we compiled a list of all 300 trained models and their different hyperparameter values, number of parameters and achieved accuracies Part of it can be seen in Table 5.2.

Table 5.2: Partial results of first random search of hyperparameters

Model ID	FilterNum1	FilterNum2	FilterNum3	DenseSize	DropoutRate	FilterSize	LearningRate	Number of parameters	Accuracy[%]
0a	72	80	64	72	0.4	3x4	0.0003	1,514,400	98.35
1a	32	40	72	56	0.35	3x4	0.0001	1,260,332	98.31
2a	40	48	32	64	0.35	3x4	0.0001	656,797	98.31
3a	56	16	48	72	0.4	3x4	0.0001	1,057,924	98.28
4a	80	64	40	96	0.45	3x4	0.0003	1,245,788	98.28
96a	16	32	72	80	0.25	3x4	0.0001	1,762,508	98.00
97a	72	56	40	56	0.45	3x4	0.0003	748,580	98.00
98a	32	24	24	48	0.35	3x3	0.0001	358,308	98.00
99a	48	16	40	40	0.45	3x3	0.0003	493,412	98.00
100a	24	72	64	40	0.45	3x3	0.0003	844,684	98.00
191a	64	56	16	52	0.4	3x3	0.0001	386,996	97.76
192a	48	40	24	24	0.4	3x4	0.0001	208,172	97.73
193a	56	64	72	24	0.25	3x4	0.0003	617,692	97.73
194a	48	72	48	32	0.25	3x4	0.0003	544,652	97.73
295a	48	32	64	16	0.5	3x4	0.0001	351,012	95.87
296a	40	24	56	24	0.5	3x4	0.0001	431,572	95.77
297a	56	16	80	16	0.2	3x4	0.0001	411,020	95.63
298a	24	16	48	24	0.5	3x4	0.0001	359,924	94.46
299a	40	48	56	16	0.35	3x3	0.0003	310,860	82.86

After analyzing results we came to several conclusions:

1. We saw that almost all trained models, except for the last one, achieved accuracy above 90 %. This proved that the general architecture of the model was appropriate for the problem.
2. We could not see any visible correlation between a specific choice of a certain

hyperparameter and accuracy. This showed that selection of hyperparameters is a non-heuristic task, at least for our particular problem.

3. Filter of size  $3 \times 4$  did not perform significantly better compared to one with size  $3 \times 3$ .
4. There is a weak correlation between the number of parameters (model's complexity) and accuracy, however, models with small size and great accuracy exist, model *2a* is an example of that.
5. First 200 models cover an accuracy range of 0.62 %. However inside of this range model number of parameters varies hugely, for example, model *192a* has more than 8 times fewer parameters than model *96a*, although the difference in accuracy (0.27 %) is negligible.

It is apparent from results that large models are not necessary to achieve high accuracy on our training data, so we decided to run the random search of hyperparameters again.

This time we lowered the maximum and the minimum numbers of filters and size of the dense layer. We decreased all steps from 8 to 2, thus increasing the number of possible configurations. We decided to lower the bottom boundary of `DropoutRate` from 0.2 to 0.0, which means that some models will not be using dropout layer at all. We expected that training without dropout layer would produce suboptimal results, however, we wanted to test it. Redefined search space for second random search can be seen in Table 5.3 We increased the number of `MAX_TRIALS` from 300 to 500, as we were expecting that more models will end up underfitting and also because there would be more possible options because of smaller step size. Partial table of results of random hyperparameter search can be seen in Table 5.4.

Table 5.3: Second hyperparameter search space

Hyperparameter	Set of values
FilterNum1	From 4 to 48, with a step of 2
FilterNum2	From 4 to 48, with a step of 2
FilterNum3	From 4 to 48, with a step of 2
FilterSize	3 x 3 or 3 x 4
DenseSize	From 4 to 48, with a step of 2
DropoutRate	From 0.0 to 0.5, with a step of 0.05
LearningRate	0.0001 or 0.0003
Random search variable	value
EPOCHS	25
BATCH_SIZE	100
MAX_TRIALS	500

Table 5.4: Partial results of second random search of hyperparameters

Model ID	FilterNum1	FilterNum2	FilterNum3	DenseSize	DropoutRate	FilterSize	LearningRate	Number of parameters	Accuracy[%]
0b	40	20	20	48	0.25	3x4	0.0001	304,216	98.14
1b	44	10	28	42	0.2	3x4	0.0003	362,264	98.14
2b	18	38	26	38	0.1	3x4	0.0003	316,956	98.11
95b	20	16	34	40	0.3	3x3	0.0003	416,230	97.62
96b	46	42	28	32	0.4	3x3	0.0003	297,466	97.62
97b	30	26	30	34	0.2	3x3	0.0001	320,570	97.59
195b	28	16	40	24	0.1	3x3	0.0001	298,252	97.31
196b	44	30	32	20	0.3	3x4	0.0003	220,098	97.31
197b	46	40	10	40	0.1	3x3	0.0001	140,874	97.31
295b	20	8	34	26	0.3	3x3	0.0003	269,464	96.90
296b	18	16	10	20	0.3	3x4	0.0003	65,740	96.87
297b	8	22	28	16	0.1	3x3	0.0001	141,742	96.87
395b	10	20	12	30	0.0	3x3	0.0001	112,246	96.87
396b	24	24	46	18	0.2	3x3	0.0003	263,924	96.14
397b	6	18	12	24	0.4	3x4	0.0001	90,520	96.11
497b	42	30	22	6	0.4	3x3	0.0003	57,386	82.86
498b	4	4	20	12	0.4	3x3	0.0003	72,992	82.86
499b	32	36	36	4	0.15	3x3	0.0001	65,648	82.86

Some observations:

1. We can see that the accuracy of the best model *0b* compared to the best model *0a* from the previous search is only 0.21 % lower, although it has about 5 times fewer parameters.
2. Although that it might seem that `FilterSize` of 3 x 4 yields best results, we did not saw a strong tendency towards 3 x 3 or 3 x 4 filter size after manually analyzing best 30 models.
3. We can see that the worst three models have the same accuracy of 82.86 %, same as the worst-performing model from the first random search. There are 82.86 % images of elephants in the validation class, which means that the model probably assigned all validation images to elephant class and was satisfied with achieved accuracy.
4. We can see that the model *296b* has a quite low number of parameters, only 65,740 when compared to its neighbours.

### 5.1.2 Comparison of selected, re-trained models

Two random searches gave us a large number of different models to choose from. In every other ML application where the execution time would not be a constraint, we could simply take the best performing model and be done with it. In our case, we had to make a trade-off between model's accuracy and execution speed.

For comparison and later on device performance testing we decided to pick and retrain<sup>16</sup> models: *0a*, *2a*, *0b*, *172b*, *338b* and *460b*, their properties are listed in Table 5.5.

Chosen models vary greatly in the number of parameters. Models *0a*, *2a*, *0b* have high number of parameters but their accuracy is high. Models *172b*, *338b* and *460b* were chosen because of their small size and reasonably good accuracy.

---

<sup>16</sup>Retraining was required as Keras Tuner module only saved hyperparameter settings during search and not each trained model. As the weights are initially randomized, the accuracy of retrained models is going to be similar but not exact when compared to the accuracy returned by random search.

Table 5.5: Properties of selected models

Model ID	FilterNum1	FilterNum2	FilterNum3	DenseSize	DropoutRate	FilterSize	LearningRate	Number of parameters	Accuracy[%]
0a	72	80	64	72	0.4	3x4	0.0003	1,514,400	98.35
2a	40	48	32	64	0.35	3x4	0.0001	656,797	98.31
0b	40	20	20	48	0.25	3x4	0.0001	304,216	98.14
172b	42	44	8	14	0.1	3x4	0.0001	60,672	97.38
338b	4	18	6	10	0.05	3x4	0.0003	20,290	96.63
460b	6	28	4	8	0.1	3x4	0.0003	13,114	93.60

As we are dealing with an imbalanced dataset, where 82.86 % of our validation data consists of elephant images, accuracy is not the best metric to use when comparing models. Simply classifying all images into elephant class would yield an accuracy of 82.86 %, which sounds high, although it would not actually do any classification.

When analysing the performance of a model on an imbalanced dataset it is more appropriate to use precision and recall metrics<sup>2</sup>. They can give us a better idea of how well the model is performing on data of specific classes. Calculated metrics can be seen in Table 5.6, we abbreviated precision to P and recall to R for clarity.

Table 5.6: Precision and recall metrics of trained models

Model ID	0a	2a	0b	172b	338b	460b
<b>Metrics</b>						
accuracy[%]	98.18	98.04	98.04	96.80	96.28	93.4
Number of parameters	1,515K	657K	304K	61K	20K	13K
P of elephant class[%]	99.22	99.46	99.25	99.29	98.80	97.80
P of human class[%]	96.92	95.38	95.38	92.00	91.69	80.31
P of cow class[%]	90.99	93.69	90.09	84.68	75.68	69.37
P of nature/random class[%]	77.42	64.52	79.03	46.77	59.68	40.32
R of elephant class[%]	99.29	98.80	98.84	97.87	98.43	97.39
R of human class[%]	93.20	94.51	95.09	91.44	89.22	85.57
R of cow class[%]	94.39	92.04	96.15	89.52	84.00	81.91
R of nature/random class[%]	87.27	97.56	84.48	93.55	67.27	28.09

<sup>2</sup>Precision tells us what percentage of data points in a specific predicted class fall into that class. Recall tells us what percentage of data points inside a certain class were actually predicted correctly [14].

As we can see all six models are generally classifying elephants correctly, both precision and recall of elephant class are high, above 97 %, which is important. Precision and recall values of other classes are generally lower, especially for nature/random. We can see that top three models *0a*, *2a* and *0b* are quite similar in terms of precision and recall, which means that we can easily prefer *0b*, without sacrificing accuracy. Models *172b* and *338b* perform a bit worse when compared to top three models, however, they have a low number of parameters which should translate to lower inference time. Last model, *460b*, performs the worst and it should generally not be used.

Another way to compare models performance is to look at a confusion matrix. Figure 5.1 shows comparison between confusion matrices of *0a* model on the left and *460b* model on the right. In the case of *0a*, 19 elephant images were not classified correctly, and 17 images were wrongly classified as elephants. This is not ideal, however is much better compared to performance of *460b*, where 53 elephants were wrongly classified and 63 of images classified as elephants were not actually elephants.

		Model 0a				Model 460b				
		elephant	human	cows	nature/ random	elephant	human	cows	nature/ random	
True Label	elephant	2388	12	3	4	elephant	2354	22	5	26
	human	6	315	3	1	human	32	261	8	24
	cows	2	6	101	2	cows	9	11	77	14
	nature/ random	9	5	0	48	nature/ random	22	11	4	25
		elephant	human	cows	nature/ random	elephant	human	cows	nature/ random	
		Predicted Label				Predicted Label				

Figure 5.1: Confusion matrices of *0a* model (left) and *460b* model (right).

### 5.1.3 Comparison of Edge Impulse models

We wanted to take our 6 models and compare them against 6 Edge Impulse models that were created by using the same hyperparameters. However, at the time of writing Edge Impulse supported only model training on images of the same dimensions. Images with different dimensions could either be cropped or scaled to fit 1:1 ratio. Using the same hyperparameters in Edge Impulse Studio, that were used in our models, would always create a model with a smaller number of parameters. The smaller image creates a smaller network when compared to a bigger image, given that the rest of architecture does not change. That meant that we could not make a direct comparison between our models and models trained in Edge Impulse Studio. We also could not perform a random search of hyperparameters in Edge Impulse Studio, as this feature was not fully supported at the time of writing this thesis.

We decided to train a few differently sized models, using the same general CNN architecture as before, but with some minor changes in hyperparameter values. We also trained a few models with Transfer Learning technique. Edge Impulse offers scaled-down versions of pre-trained MobileNetV2<sup>3</sup>NN architecture, which we used.

Tables 5.7 and 5.8 show properties of Edge Impulse models using CNN architecture and Transfer Learning technique, respectively. Table 5.9 shows calculated precision and recall values of Edge Impulse models using both approaches.

We used only two different versions of MobileNetV2, 0.35, and 0.1, as we saw an accuracy drop in the reduction of width multiplier hyperparameter. In all cases the pre-trained model was followed by a one or two dense layers, with dropout layers in between.

---

<sup>3</sup>MobileNetV2 is a efficient, lightweight NN architecture, designed for image recognition tasks, suitable for mobile applications [14]. MobileNetV2 contains width multiplier hyperparameter, which scales up or down the total number of parameters, thus providing a trade-off between accuracy and computation complexity. Edge Impulse offers three different width multiplier options: 0.35, 0.1 and 0.05.

Table 5.7: Properties of Edge Impulse models using CNN architecture.

Model ID	FilterNum1	FilterNum2	FilterNum3	DenseSize	DropoutRate	FilterSize	LearningRate	Number of parameters	Accuracy[%]
0ei	72	80	64	72	0.4	3x4	0.0003	1,168,804	97.7
1ei	40	48	32	64	0.35	3x4	0.0001	503,196	97.5
2ei	40	20	20	48	0.25	3x4	0.0001	231,204	97.3
3ei	42	44	8	14	0.1	3x4	0.0003	52,272	96.6

Table 5.8: Properties of Edge Impulse models using Transfer Learning technique

Model ID	Width Multiplier	DenseSize1	DenseSize2	DropoutRate	LearningRate	Number of parameters	Accuracy[%]
0tl	0.35	16	N/A	0.1	0.0005	430,676	98.5
1tl	0.35	16	16	0.1	0.0005	430,948	98.4
2tl	0.35	32	32	0.1	0.0005	452,484	98.7
3tl	0.1	32	32	0.1	0.0005	135,732	95.7

Table 5.9: Precision and recall metrics of trained Edge Impulse models

Model ID	0e	1e	2e	3e	0tl	1tl	2tl	3tl
Metrics								
Accuracy[%]	97.7	97.5	97.3	96.6	98.5	98.4	98.7	95.7
Number of parameters	1,169K	503K	231K	52K	430K	431K	452K	136K
P of elephant class[%]	99.69	99.53	99.42	99.27	99.27	99.42	99.48	98.65
P of human class[%]	95.05	95.05	91.87	91.52	97.17	95.41	96.47	85.16
P of cow class[%]	82.22	78.89	82.22	79.57	92.22	91.01	92.22	75.56
P of nature/random class[%]	63.04	65.22	73.91	50.0	86.96	89.13	91.3	78.85
R of elephant class[%]	99.86	98.91	98.44	98.65	99.48	99.27	99.37	98.03
R of human class[%]	93.4	92.44	94.2	88.4	94.5	95.74	95.79	90.26
R of cow class[%]	90.24	86.59	84.09	79.57	92.22	89.01	94.32	67.33
R of nature/random class[%]	87.88	88.24	94.44	95.83	95.24	97.62	95.45	91.11

Some observations:

- Models using CNN architecture did not out perform our models in terms of accuracy. Models *2a* and *0b* both had accuracy of 98.04 %, while none of Edge Impulse models with CNN architecture did not pass 98 %.
- Most of the models trained with Transfer Learning technique outperformed our models.
- Model *2tl* performed exceptionally well, reaching an accuracy of 98.7 % while having a relatively small number of parameters.
- We saw that by decreasing width multiplier we did not benefit much in accuracy as much as we lost in model size. Even increasing the sizes of dense layers did not solve the problem.

## 5.2 On-device performance testing

Performance testing of all models was done on an STM32F767ZI microcontroller, running at 216 MHz. Testing of our models was done by using MicroML framework that we wrote, which directly called TFLite Micro API. Testing of Edge Impulse models was done on Mbed OS, as this platform is supported by Edge Impulse and they already provide an example for it. We could not test model *0a* on device as TFLite converter failed to produce a compilable model.

To time the execution of our code we used Data Watchpoint Trigger (DWT) which contains a counter that is directly incremented by the system clock. DWT does not use interrupts, therefore it does not introduce the overhead of calling interrupt routines as systick timer does.

Edge Impulse provides examples for testing out of the box, so not much work is needed to get the first-order approximation of performance. For profiling, the code execution Mbed API was used, which uses timer interrupts to track elapsed time. Figure 5.2 shows models ranked from the fastest to the slowest.

We can see that all models perform inference in less than 1 second, which was a

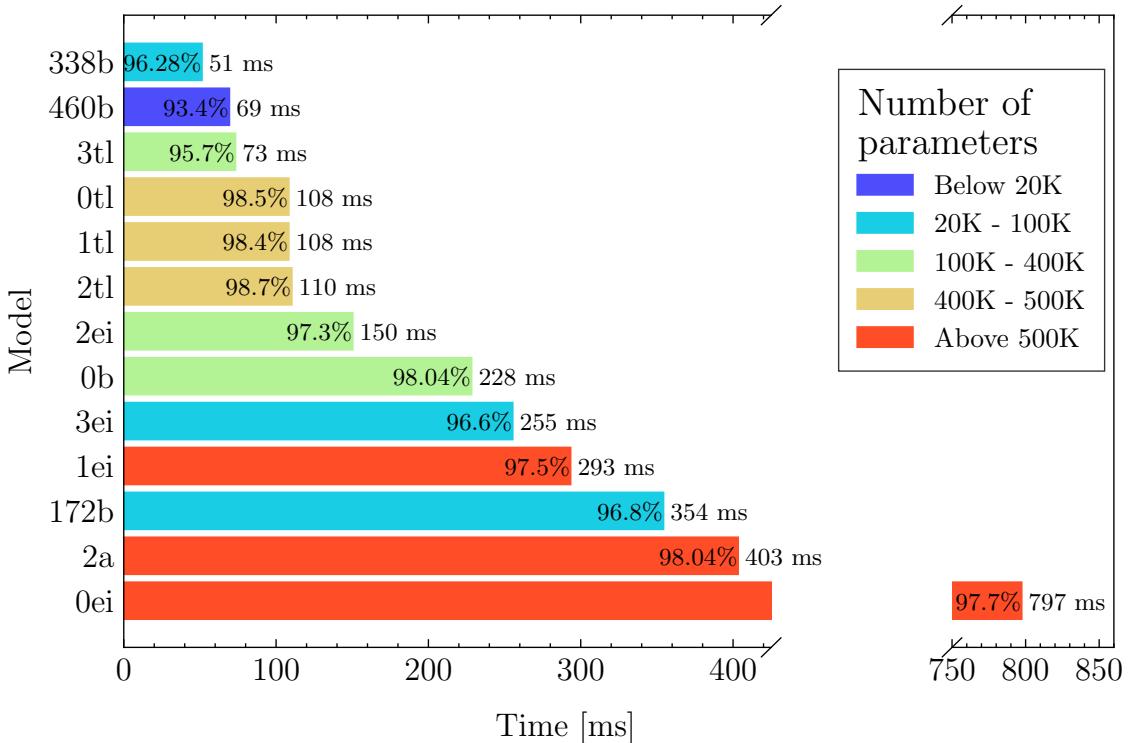


Figure 5.2: Comparison of time of inference of different models.

constraint that we set earlier in Section 3.1. Best time wise performing model was *338b* with inference time of 51 ms, there are also many models that perform inference under 300 ms.

We also discovered some unexpected trend in results. We assumed that inference time is proportional to the number of parameters if the general structure of the model remains the same. As can be seen from the Figure 5.2 there are few exceptions to this rule, model *338b* executed inference faster than *460b* although it has more parameters (20K versus 13K). Model *172b* was slower than *0b*, although it has five times less parameters. This behaviour is not exclusive only to our models, but it can be seen in Edge Impulse models as well, for example, models *2ei* and *1ei*.

Edge Impulse models trained with Transfer Learning technique *0tl*, *1tl*, *2tl* and *3tl* should not be compared to other models in this sense, as the architecture of MobileNetV2 contains additional different operations.

We can only speculate about the reason for this behaviour, since it is present both in our models and Edge Impulse models, we can assume this to be a TFLite bug.

### 5.2.1 Comparison of code sizes

We also wanted to inspect Flash and RAM sizes of binaries, that we compiled for on-device testing. For this task we used `arm-none-eabi-size` command line tool which returns sizes of `text`, `data`, `bss` sections in bytes, example of output can be seen in Figure 5.3. To compute used Flash we need simply add bytes from `text` and `data` sections, and to compute used RAM we add together `data` and `bss` sections<sup>4</sup>.

```
1 $ arm-none-eabi-size firmware.elf
2      text      data      bss      dec      hex filename
3 149124       388    47064   196576  2ffe0 firmware.elf
```

Figure 5.3: Example output of `arm-none-eabi-size` command.

Code sizes for all models are presented in Figure 5.4, models are ordered the same way as they have been in Figure 5.2.

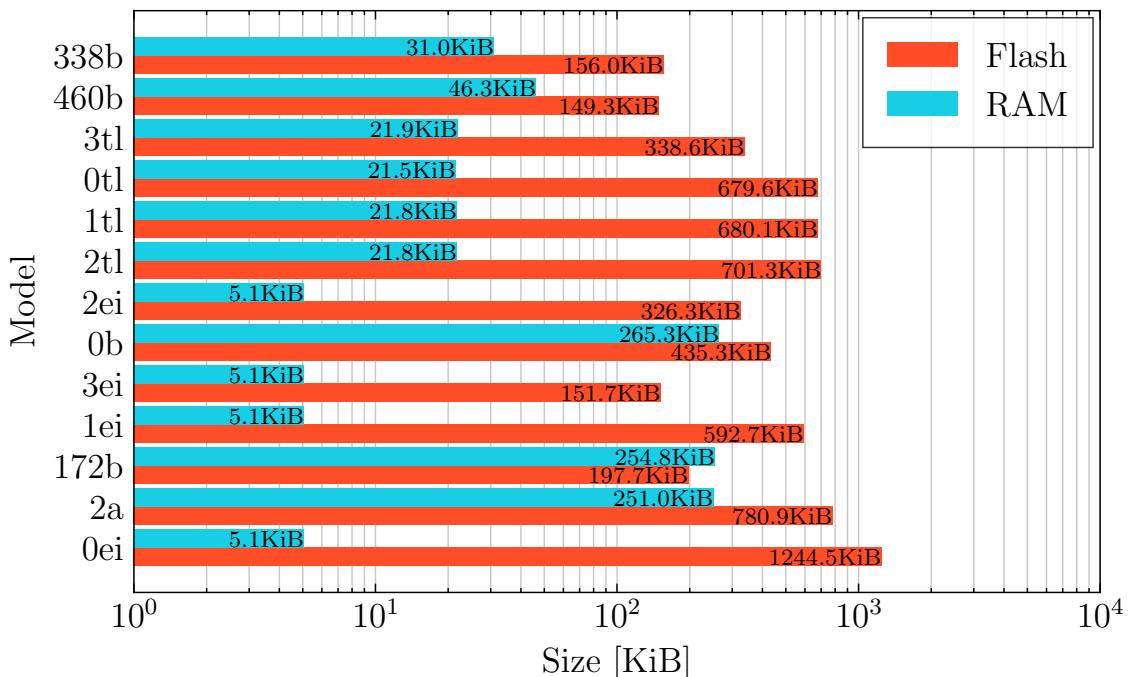


Figure 5.4: Comparison of Flash and RAM size of compiled example models.

We can see that all of our models generally use more RAM then edge Impulse models. This is due to how the inference is executed. TFLite Micro uses a generic interpreter

---

<sup>4</sup>Data section which contains initialized static variables is first placed into Flash memory and is copied to RAM before program enters main function. That is why we have to account for additional `data` section in Flash memory.

approach, where the model is loaded at runtime. Edge Impulse uses a compiled approach, which they named The EON<sup>TM</sup> Compiler [38]. The EON<sup>TM</sup> Compiler still uses TFLite Micro, however, it does not use its interpreter, but directly calls operation kernels. This means that linker knows exactly which operations are used and more data can be moved into Flash, thus eliminating unneeded code space [38].

### 5.2.2 Comparison of different optimisation options

To be able to run the ML inference at maximum possible efficiency under MicroML some extra amount of work and research was required. Figure 5.5 shows reductions in inference time of *0b* model while using different optimisation methods.

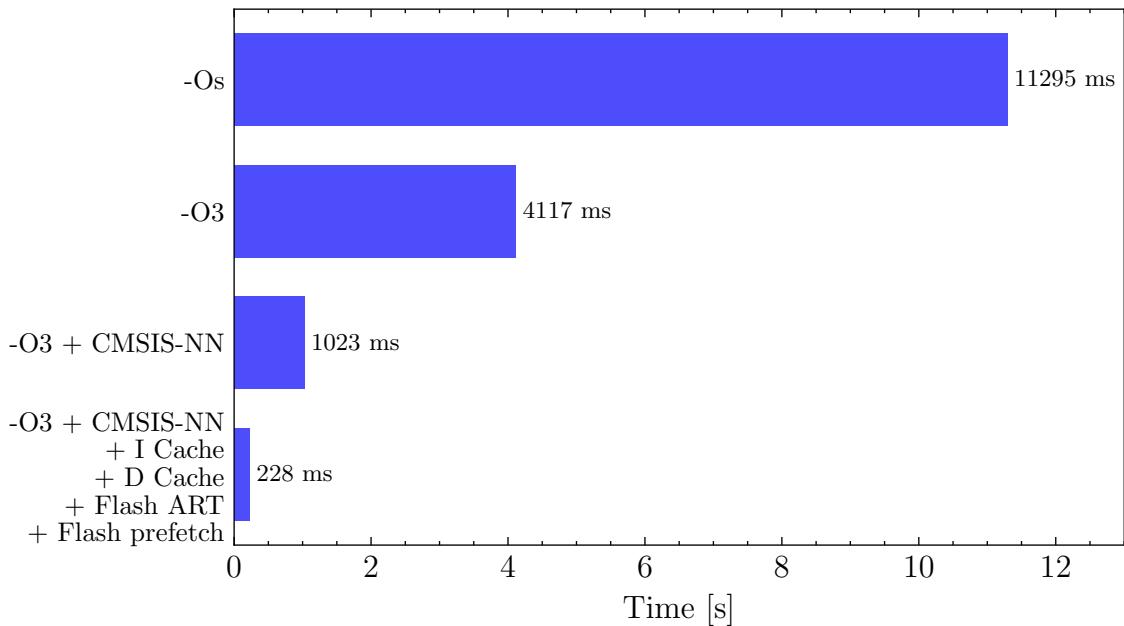


Figure 5.5: Inference time of *0b* model using different optimisations.

We started with no optimisations at all, while using only **-Os** compiler flag. **-Os** flag generally optimises for minimal size, it enables all **-O2** optimisations, except those that increase size. This optimisation level is often used, however, we found out that inference time of more than 11 seconds was too long.

Changing optimisation level to **-O3** drastically decreased the time of inference, down to 4117 ms. **-O3** turns on all **-O2** optimisations plus additional ones and completely

disregards any code size optimisations.

Changing compiler optimisations flags could not lower time of inference any further, so other approaches were needed. While reading through TensorFlow documentation we saw that it supports CMSIS-NN library for ARM microcontrollers. CMSIS-NN is a collection of efficient Neural Network kernels, that intends to maximise performance and lower code size of NN models implemented on ARM microcontrollers. TensorFlow provides wrappers for some of these kernels, such as convolutional, fully connected, pooling layers and others. Not much work was needed to use these highly efficient kernels, we only needed to specify in our Makefile that we want to compile CMSIS-NN kernels and not compile generic TensorFlow kernels. Time of inference dropped for about 3 seconds, down to 1023 ms.

As we saw that similarly sized Edge Impulse models were running much faster on Mbed platform compared to our MicroML code, while using the same microcontroller, we knew that there was another step left. The final performance increase was reached by using features only fully found in Cortex-M7 microcontrollers and partly in Cortex M3/4 microcontrollers. To achieve it we had to enable I and D caches, flash prefetch and flash ART.

ART stands for Adaptive real-time memory accelerator, which encompasses I/D caches and flash prefetch buffer. I and D caches are small, efficient portions of memory, which are located in the CPU of the microcontroller. They hold instructions and data respectively, if those are requested by next microcontroller instruction, they can be read much faster compared to reading them from flash memory. By enabling flash prefetch microcontroller reads additional sequential instructions into prefetch buffer, thus enabling execution without any wait states (if instruction flow is sequential). Elimination of wait states greatly improved the time of inference, it was decreased to 228 ms.

### 5.2.3 Scoring trained models

Choosing the best model for on-field deployment is a hard task due to many different metrics: precision and recall values, time of inference, and code size. To make this

job easier we devised a scoring system: each metric is going to be normalized and multiplied with some weight value. All products would then be summed up and the result would represent the final score. The Sum of the weights is equal to 100, which means that the possible maximum score is also 100. We decided to allocate 50 weight points to all precision and recall values, 30 points to the time of inference value, 5 points for Flash size and 15 points for RAM size. As we care more about precision and recall values of elephant, human and cow classes we gave them 7 weight points each, while nature/random class received only 4. We value Flash size less than RAM, as most of the microcontrollers have much less RAM than they do Flash, thus we gave 15 weight points to RAM size and only 5 to Flash size.

Since the time of inference, Flash and RAM sizes are properties, which should give a larger score, the smaller they are, we mapped them into a range between 0 and 1. The smallest value inside the set would be assigned 1, the biggest 0, the values in between are mapped linearly.

Scoring is described mathematically in 5.1, while the final results can be seen in Figure 5.6

$$\begin{aligned}
Score[i] &= 7K(P_{elephant}, i) + 7K(P_{human}, i) + 7K(P_{human}, i) + 4K(P_{ntr/rnd}, i) \\
&\quad + 7K(R_{elephant}, i) + 7K(R_{human}, i) + 7K(R_{human}, i) + 4K(R_{ntr/rnd}, i) \\
&\quad + 30 F(ToI, i) + 5 F(Flash, i) + 15 F(RAM, i) \\
K(X, i) &= \frac{(X[i] - MIN(X))}{MAX(X) - MIN(X)} \\
F(X, i) &= \frac{(X[i] - MAX(X))}{MIN(X) - MAX(X)}
\end{aligned} \tag{5.1}$$

Where:

$Score$  - Vector of calculated scores

$i$  -  $i^{th}$  model

$P_j$  - Vector of precision values of  $j^{th}$  class

$R_j$  - Vector of recall values of  $j^{th}$  class

$ToI$  - Vector of Time of Inference values

$Flash$  - Vector of Flash sizes

$RAM$  - Vector of RAM sizes

$K(X, i)$  - Normalizing function with vector X and element index i as arguments

$F(X, i)$  - Normalizing, inverting, function with vector X and element index i as arguments

$MAX(X)$  - Function that finds maximum element in vector X

$MIN(X)$  - Function that finds minimum element in vector X

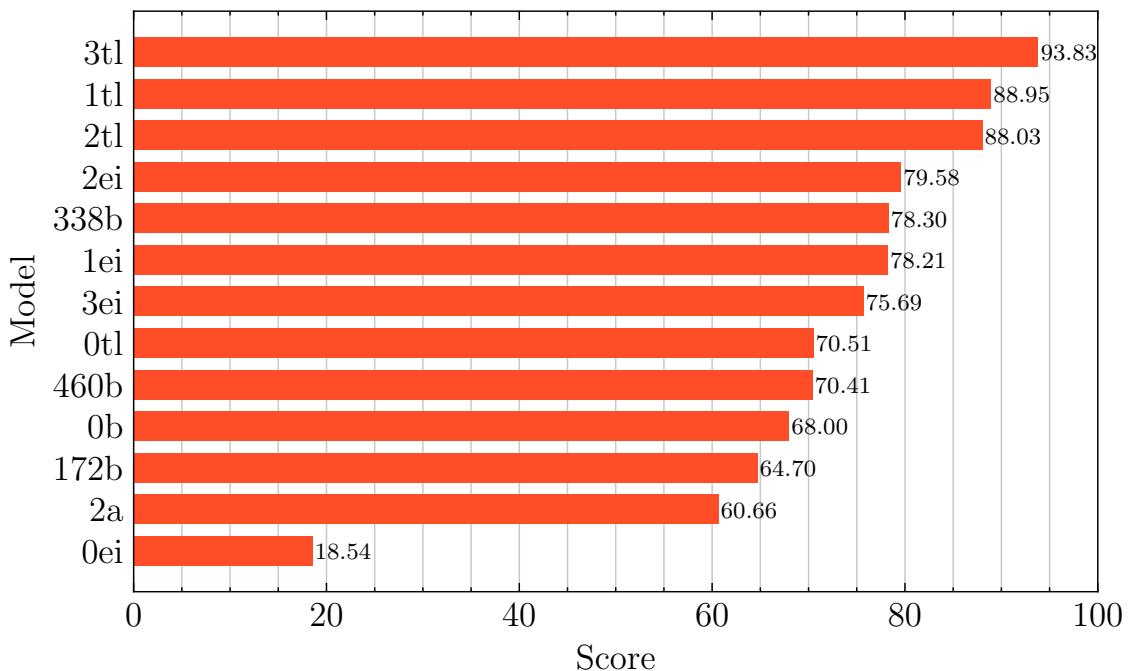


Figure 5.6: Inference time of  $0b$  model using different optimisations.

Final verdict: model  $3tl$  received the highest score, models  $1tl$ ,  $2tl$  follow. This should not be a surprise, models trained with Transfer Learning achieve high accuracies and execute inference in about 100 ms or less while compiled approach for computing Neural Networks keeps used Flash and RAM sizes to a minimum.

### 5.3 Power profiling of an embedded early warning system

## Bibliography

- [1] Nyhus, P. J. Human–Wildlife Conflict and Coexistence. *Annual Review of Environment and Resources*, 41, (2016), 11, pages 143–171.
- [2] SARPO, WWF. Human Wildlife Conflict Manual. Available on:  
[https://wwf.panda.org/our\\_work/wildlife/human\\_wildlife\\_conflict/hwc\\_news/?84540/Human-Wildlife-Conflict-Manual](https://wwf.panda.org/our_work/wildlife/human_wildlife_conflict/hwc_news/?84540/Human-Wildlife-Conflict-Manual), [12.06.2020].
- [3] WILDLABS, WWF. Human Wildlife Conflict Tech Challenge: Polar Bear Case. Available on:  
<https://www.wildlabs.net/hwc-tech-challenge-polar-bear-case>, [14.06.2020].
- [4] WILDLABS, WWF. Human Wildlife Conflict Tech Challenge: Tiger Case. Available on:  
<https://www.wildlabs.net/hwc-tech-challenge-tiger-case>, [14.06.2020].
- [5] WILDLABS, WWF. Human Wildlife Conflict Tech Challenge: Asian Elephant Case. Available on:  
<https://www.wildlabs.net/hwc-tech-challenge-elephant-case>, [14.06.2020].
- [6] Ganesh, S. Human-elephant conflict kills 1,713 people, 373 pachyderms in 3 years. Available on: <https://www.thehindu.com/news/national/human-elephant-conflict-kills-1713-people-373-pachyderms-in-3-years/article26225515.ece>, [14.06.2020].

- [7] Guha N., In the Heart of the Conflict: Understanding the Human Elephant Dynamics in Udaguri. Available on: <https://www.econe.in/post/in-the-heart-of-the-conflict-understanding-the-human-elephant-dynamics-in-udaguri>, [08.09.2020].
- [8] Save our species, Human wildlife conflict - global challenge: local solutions. Available on: <https://www.saveourspecies.org/news/human-wildlife-conflict-global-challenge-local-solutions>, [08.09.2020].
- [9] The Week, Arresting image of human elephant conflict wins photo prize. Available on: <https://www.theweek.co.uk/89566/arresting-image-of-human-elephant-conflict-wins-photo-prize>, [08.09.2020].
- [10] WILDLABS, WWF. Human Wildlife Conflict Tech Challenge. Available on: <https://www.wildlabs.net/hwc-tech-challenge>, [14.06.2020].
- [11] Icons8 - Icons used in various figures. Available on: <https://icons8.com/>, [21.9.2020].
- [12] WILDLABS, WWF. HWC Tech Challenge Winners Announced. Available on: <https://www.wildlabs.net/resources/news/hwc-tech-challenge-winners-announced>, [20.06.2020].
- [13] Dangerfield A. Progress report – January 2019 – Thermal imaging for human-wildlife conflict. Available on: <https://blog.arribada.org/2019/01/10/progress-report-january-2019-thermal-imaging-for-human-wildlife-conflict>, [20.06.2020].
- [14] Geron, A. *Hands-on machine learning with Scikit-Learn and TensorFlow: concepts, tools, and techniques to build intelligent systems, 2nd edition*. O'Reilly Media, Sebastopol, CA, 2019.
- [15] Burkov, A. *The Hundred-Page Machine Learning Book*. Andriy Burkov, 2019.

- [16] Li F., Karpathy A., “Cs231n: Convolutional neural net- works for visual recognition.” Stanford University course. Available on: <http://cs231n.stanford.edu/>, [25.06.2020].
- [17] Zhang, Y., Suda, N., Lai, L., and Chandra, V. Hello edge: Keyword spotting on microcontrollers. *ArXiv*, abs/1711.07128, (2017), 2.
- [18] Louis, M. S., Azad, Z., Delshadtehrani, L., Gupta, S., Warden, P., Reddi, V. J., and Joshi, A. Towards deep learning using tensorflow lite on risc-v. *Third Workshop on Computer Architecture Research with RISC-V (CARRV)*, 1, (2019), 6.
- [19] Warden P., Why the future of machine learning is tiny. Available on: <https://petewarden.com/2018/06/11/why-the-future-of-machine-learning-is-tiny/>, [06.07.2020].
- [20] Situnayake D., Make deep learning models run fast on embedded hardware. Available on: <https://www.edgeimpulse.com/blog/make-deep-learning-models-run-fast-on-embedded-hardware/>, [08.07.2020].
- [21] Dive into deep learning, Convolutional Neural Networks. Available on: [http://d2l.ai/chapter\\_convolutional-neural-networks/conv-layer.html](http://d2l.ai/chapter_convolutional-neural-networks/conv-layer.html), [17.9.2020].
- [22] TensorFlow, GitHub repository. Available on: <https://github.com/tensorflow/tensorflow>, [21.9.2020].
- [23] Rouse M., internet of things (IoT). Available on: <https://internetofthingsagenda.techtarget.com/definition/Internet-of-Things-IoT>, [27.9.2020].
- [24] Ubidots, LoRaWAN vs NB-IoT: A Comparison Between IoT Trend-Setters. Available on: <https://ubidots.com/blog/lorawan-vs-nb-iot/>, [28.9.2020].
- [25] Polymorph, IoT connectivity comparison (GSM vs LoRa vs Sigfox vs NB-Iot). Available on: <https://www.polymorph.co.za/>

- iot-connectivity-comparison-gsm-vs-lora-vs-sigfox-vs-nb-iot/, [28.9.2020].
- [26] Bäumker, E., Garcia, A., and Woias, P. Minimizing power consumption of lora (R) and lorawan for low-power wireless sensor nodes. *Journal of Physics: Conference Series*, 1407, (2019), 11, page 012092.
- [27] Knight M., A GitHub repository containing open-source implementation of the LoRa CSS PHY. Available on:  
<https://github.com/BastilleResearch/gr-lora>, [29.9.2020].
- [28] Wong G.W., LoRa Rolls Into Philly. Available on:  
<https://www.electronicdesign.com/technologies/embedded-revolution/article/21805205/lora-rolls-into-philly>, [29.9.2020].
- [29] Vollmer, M. and Möllmann, K. P. *Infrared Thermal Imaging: Fundamentals, Research and Applications*. Wiley-VCH, Boston, Massachusetts, 2018.
- [30] Dangerfield A., HWC Tech Challenge Update: Comparing thermopile and microbolometer thermal sensors. Available on:  
<https://www.wildlabs.net/resources/case-studies/hwc-tech-challenge-update-comparing-thermopile-and-microbolometer-thermal>, [18.07.2020].
- [31] Bhan, R., Saxena, R., Jalwania, C., and Lomash, S. Uncooled infrared microbolometer arrays and their characterisation techniques. *Defence Science Journal*, 59, (2009), 11, page 580.
- [32] MoviTherm, What is NETD in a Thermal Camera? Available on:  
<https://movitherm.com/knowledgebase/netd-thermal-camera/>, [18.07.2020].
- [33] Dangerfield A., Progress report – February 2020 – Thermal imaging for human-wildlife conflict. Available on:  
<https://blog.arribada.org/2020/02/17/>

- [progress-report-feburart-2020-thermal-imaging-for-human-wildlife-conflict/](https://github.com/progress-report-feburart-2020-thermal-imaging-for-human-wildlife-conflict/), [02.10.2020].
- [34] GroupGets - LeptonModule, GitHub repository. Available on: <https://github.com/groupgets/LeptonModule>, [21.9.2020].
- [35] Mpaland, A printf / sprintf Implementation for Embedded Systems, GitHub. Available on: <https://github.com/mpaland/printf>, [27.10.2020].
- [36] Sagadin M., MicroML, Quick-start machine learning projects on microcontrollers with help of TensorFlow Lite for Microcontrollers and libopencm3, GitHub repository. Available on: <https://github.com/MarkoSagadin/MicroML>, [27.10.2020].
- [37] Roeder, L., Netron, Visualizer for neural network, deep learning, and machine learning models. Available on: <https://netron.app/>, [30.10.2020].
- [38] Jongboom J., Introducing EON: neural networks in up to 55less ROM. Available on: <https://www.edgeimpulse.com/blog/introducing-eon>, [20.11.2020].