

**Bachelor Project**



**Czech  
Technical  
University  
in Prague**

**F3**

**Faculty of Electrical Engineering  
Department of Cybernetics**

## **Drone detection using neural networks from combined RGB camera and LiDAR data**

**Adam Škuta**

**Supervisor: Matouš Vrba  
Supervisor–specialist: Martin Saska  
Field of study: Mathematical Engineering  
Subfield: Mathematical Modelling  
February 2017**



## Acknowledgements

Děkuji ČVUT, že mi je tak dobrou *alma mater*.

## Declaration

Prohlašuji, že jsem předloženou práci vypracoval samostatně, a že jsem uvedl veškerou použitou literaturu.

V Praze, 10. February 2017

## Abstract

**Keywords:** word, key

**Supervisor:** Matouš Vrba  
Ústav X,  
Uliční 5,  
Praha 99

## Abstrakt

**Klíčová slova:** slovo, klíč

**Překlad názvu:** Moje bakalářka se  
strašně, ale hrozně dlouhým předlouhým  
názvem — Cesta do tajů kdovíčeho

# Contents

<b>1 Introduction</b>	<b>1</b>
<b>2 Convolutional neural networks</b>	<b>3</b>
2.1 Sparse to dense . . . . .	3
2.2 YOLOv3 . . . . .	5
<b>3 Sensors</b>	<b>7</b>
3.1 Coordinate systems . . . . .	7
3.2 Camera Model . . . . .	8
<b>4 Dataset</b>	<b>9</b>
4.1 Unreal Engine . . . . .	9
4.2 AirSim . . . . .	10

2.1 Example of different network architectures available .....	3
--	---

# Chapter 1

## Introduction

The goal of this thesis is to prove whether a usage of LiDAR data coupled with images from RGB is useful for the localization of UAVs in contrast to the usage of image data alone. The LiDAR and RGB camera will be mounted on top of the scanner UAV. All the measurements will be taken inside a virtual environment, with a realistic UAV and sensor simulation. The dataset will then be preprocessed and used as the input to a Convolutional neural network for the object detection. The preprocessing will be as following:

- Coordinate transformation for the non-matching coordinate systems
- Projection of 3d points into a 2d image
- Using a Convolutional neural network to generate more points in a sparse LiDAR pointcloud

Infographic of  
the process





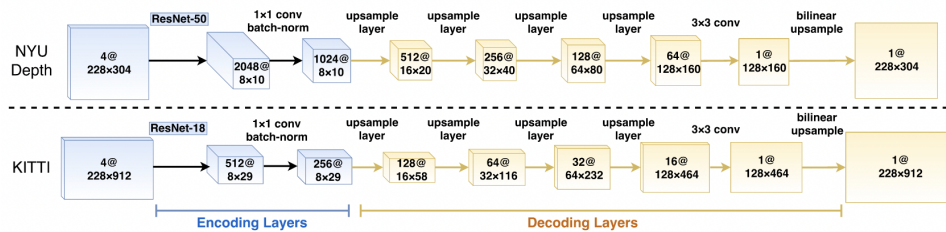
## Chapter 2

### Convolutional neural networks

A convolutional neural network will be used in two problems in the thesis. First one is used on a sparse LiDAR pointcloud generating more points and therefore making it more dense. Second one is used for object, in this case drone, detection, using RGB and RGBD data as the input.

#### 2.1 Sparse to dense

Sparse to dense is a Convolutional neural network written in PyTorch. It predicts the depth measurements from sparse depth dataset. The size of the network is modifiable and can be chosen as training parameters.



**Figure 2.1:** Example of different network architectures available

For the encoding and therefore input layers a ResNet-50 or ResNet-18 can be chosen, depending on the size of the input image for memory constraints. The decoding layers consist of 4 upsampling layers and a deconvolutional layer with either stride 2 or 3 or upprojection layer or upconvolutional layer as

a choice for training. The default loss function is least absolute deviations also known as  $L_1$  error:

$$L_1 = \sum_{i=1}^n |y_{true} - y_{predicted}| \quad (2.1)$$

where:

- $n$  is batch size
- $y_{true}$  are real depth values
- $y_{predicted}$  are predicted depth values

The input to the network are RGBD images and the output is a depth map with the same dimensions as input. The depth input  $D$  is sampled from the ground truth depth map  $D^*$  with the following formula:

$$D(i, j) = \begin{cases} D^*(i, j), & \text{with probability } p \\ 0, & \text{otherwise} \end{cases} \quad (2.2)$$

where:

- $i, j$  are coordinates of the input image
- $p = \frac{m}{n}$ , where  $m$  number of depth samples to be chosen at the start of training and  $n$  is the total amount of available depth samples

During training several input data augmentations take place. These augmentations include:

- Scaling the input image by a random number  $s \in [1, 1.5]$
- Rotating the input image by a random degree  $r \in [-5, 5]$
- Scaling the brightness, contrast and saturation of the RGB component of the image by a random number  $k \in [0.6, 1.4]$
- Normalizing the RGB component of the image
- Flipping the image horizontally with a 50% chance

## 2.2 YOLOv3

## Chapter 3

### Sensors

#### 3.1 Coordinate systems

In order to correctly label the data for training, a position of the second drone in relation to the camera mounted on the first one is required. The API call in AirSim returns a position in relation to its starting point. Therefore a transformation from the starting point of the second drone to the camera mounted on the first one is required. We can write this transformation as follows:

$$\mathbf{T} = \mathbf{T}_{d1}^c \mathbf{T}_{s1}^{d1} \mathbf{T}_{s2}^{s1} \quad (3.1)$$

where:

- $\mathbf{T}_{s2}^{s1}$  is transformation from the starting point of the second drone to the starting point of the first drone
- $\mathbf{T}_{s1}^{d1}$  is transformation from the starting point of the first drone to the body of the first drone
- $\mathbf{T}_{d1}^c$  is transformation from the body of the first drone to the cameras coordinate system

Transformation matrix  $\mathbf{T}$  can generally be described as follows:

$$\mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{p} \\ \mathbf{0}^T & 1 \end{bmatrix} \quad (3.2)$$

insert picture showing the positions of the drones

where:

- $\mathbf{R}$  is a 3x3 rotation matrix
- $\mathbf{p}$  is a 3x1 translation column vector
- $\mathbf{0}^T$  is a 1x3 row vector of zeros

## ■ 3.2 Camera Model

Picture of pin-hole camera

For the creation of the bounding boxes used for training a transformation from coordinate system of the camera to the pixel values of the image needs to be defined. For this task a pinhole camera model is used. The transformation is then defined as follows:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} f \frac{x}{z} + c_x \\ f \frac{y}{z} + c_y \end{bmatrix} \quad (3.3)$$

where:

- $x'$  and  $y'$  are pixel coordinate values on the image
- $x, y, z$  are coordinate values of a point to be transformed
- $f$  is focal length of the camera
- $c_x, c_y$  are offsets on the image plane

## Chapter 4

### Dataset

The dataset for this work can be generated in two ways. The first is real-life drone shots mixed with point clouds from LiDAR mounted on top of a drone. The second is generating a dataset using a realistic virtual environment where a drone, camera and LiDAR are being emulated very close to their real-life counterparts. An advantage to this approach is that a great variety of environments can be chosen a lot of them often inaccessible otherwise (power plant, airport, snowy mountains out of season etc.). Therefore this approach will be chosen for the task.

#### 4.1 Unreal Engine

Unreal Engine is a software tool used for creating realistic 3d environments, most often used as a video game engine. It is written in C++ and open-source supporting a variety of pre-built environments and assets. For this work three different environments will be used for the creation of the dataset:

citation  
<https://www.unrealengine.com/en-US/features>

- exact name. Snow

- exact name. Park

Pictures of the environments

exact name. City centre

Together

exact number of pictures

pictures and labels were generated using two drones. One drone was equipped with RGB camera and LiDAR sensor and was responsible for taking the pictures and pointclouds from LiDAR. The second one was used as a model for drone detection.

## 4.2 AirSim

airsim zdroj

Open-source plugin for Unreal Engine called AirSim was used for the generation of the dataset. It simulates realistic flight motions of drones as well as seven types of sensors, including RGB camera and LiDAR used for this task. AirSim supports both a C++ API as well as Python API, latter which was used for controlling the motion and capturing the dataset. Location of the second drone was generated through API call, which produces a location of the drone in global coordinate system of the map, which is later transformed to the local coordinates of the first drone carrying the LiDAR and RGB sensors.. The capturing drone traveled on each map on a 3d cube grid:

Transformacna matica?

Grafika kocky po ktorej lietal dron