

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

Title

Author:
Your name

Supervisor:
Name of supervisor

Submitted in partial fulfillment of the requirements for the MSc degree in Type of Degree of
Imperial College London

September 2015

Abstract

Your abstract.

Acknowledgments

Comment this out if not needed.

Contents

1	Introduction	1
2	Background	2
2.1	Background, Computer Vision	2
2.1.1	Deep Learning Preliminaries	2
2.1.2	Convolutional Networks	2
2.1.3	Image Detection	4
2.1.4	Image Segmentation	7
2.1.5	Instance Segmentation	8
2.1.6	Recurrent Neural Networks:	10
2.1.7	Bayesian Optimisation Layer Dimensions:	10
2.2	Background, SLAM	11
2.2.1	Sensors	12
2.2.2	Camera to World Frame	13
2.2.3	Projecting to World Frame	14
2.2.4	World to Camera Frame	16
2.2.5	Bicycle States	16
2.2.6	The SLAM Problem	17
2.3	Background Object Tracking	21
2.3.1	Kernelised Correlation Filter	21
3	Contribution	24
3.1	Depth Measurement Outlier Removal	24
3.2	Frame-to-Frame Tracking	26
3.2.1	Kernelised Correlation Filter	27
3.2.2	BRISK - Selector	27
3.2.3	Class Matching	28
3.2.4	Kalman Filter ROI Matching	28
3.2.5	ROI Lives	30
3.2.6	Combining Matching Methods	31
3.3	Movement Prediction	35
3.3.1	Objects Prediction	35
3.3.2	Bike Model	36
3.4	System Evaluation	38
3.4.1	Mask-RCNN	38
3.4.2	Okvis with Dynamic Objects	38
3.4.3	Depth Camera	38

4 Ethics	39
4.0.1 Modification for Military Purposes	39
4.0.2 Other Ethics Issues	39
4.1 Ethics Checklist	39
5 Experimental Results	42
6 Conclusion	43

Chapter 1

Introduction

There were c. 18,500 cycling accidents and 3,500 deaths on British roads in 2016 [1]. Of these accidents 80% occur during the day and driver/rider error accounts for 71% of the total collisions [2]. The most common place for a cycling incident is at a road junction with 75% of collisions occurring there [3]. With cycling becoming a more accepted form of transport, the need to protect and modify bikes with a comprehensive safety system is important.

This report will examine state of the art technologies that are used for autonomous vehicles to investigate their applicability to a cycling environment. The three main topics involve image detection/segmentation, SLAM algorithms and collision/path prediction.

Chapter 2

Background

2.1 Background, Computer Vision

2.1.1 Deep Learning Preliminaries

Computer vision and its applications has been revolutionised recently by Deep Learning frameworks. These have led to super-human accuracy in image recognition, for example Inception v3 [4], release in 2016, reached a 3.46% error rate on the ImageNet [5] dataset. This is known as an image classification task; the photo contains one object in the dataset and that object generally occupies the majority of the frame. Similar networks are used with additional features for image detection and segmentation which are integral to detecting moving objects such as pedestrians and vehicles in an urban environment.

The current state-of-the-art methods in computer vision use a form of neural network. These networks are characterised by a series of non-linear nodes that are organised in layers. At each layer the network learns progressively abstract and higher-level features of the image to be detected until the final (usually dense layer) which determines the predicted class. Therefore the 'depth' of a neural network often refers to the number of layers and hence the inherent complexity of features that a network can learn. This is, however, only a basic representation of the problem and some clever and shallower architectures are able to outperform deeper networks. It is important to note that the depth and width, more specifically the number of parameters in the model generally reduces speed of prediction and increases the memory required to store/load the model into RAM. This is of little consequence in most applications, however for the limited computational power that can be carried on a bicycle it will be of concern in this report.

2.1.2 Convolutional Networks

A CNN generally contains convolutional, pooling, fully connected and normalisation layers. This form of network became particularly popular in computer vision after Geoffrey Hinton's 2012 paper [6], with an ImageNet classification error of c.16%, significantly lower than the state-of-the-art at the time which was 26%. Convolutional networks use the assumption that each input is an image and therefore certain properties can be encoded into the architecture of the network. This allows a significant reduction in the number of tunable network parameters.

A regular neural network, as mentioned earlier Section 2.1.1, contains a number of fully

connected layers i.e. each neuron in a layer is connected to every neuron in the previous and next layers respectively. However convolutional neural networks are arranged in three dimensions: width, height, depth (Figure 2.1). In this architecture the neurons in a layer are only connected to some of the layer before, unlike in the fully connected case, this significantly reduces the number of parameters [7].

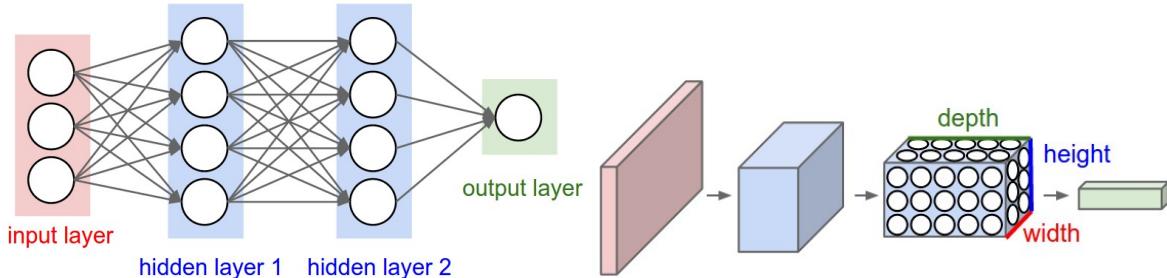


Figure 2.1: Regular Neural Network (Left) [7] , CNN (right) [7].

Convolutional Layers

A convolutional layer learns a set of filters, each filter is compact in size, say $(7 \times 7 \times 3)$. This filter slides (convolves) over the image, then a dot product is taken between the filter values and the inputs to that filter (Figure 2.2). The filter slides across with a respective stride, which is a hyperparameter, a stride of one equates to moving one pixel across at a time. A convolutional layer is therefore able to pick up features such as a corner, or higher level features such as the texture of fur. This process creates a 2D activation map containing the response of the filter at each position on the input volume. A number of filters are passed over and a number of activation maps are created, relating to the layer depth which is a tuneable hyperparameter.

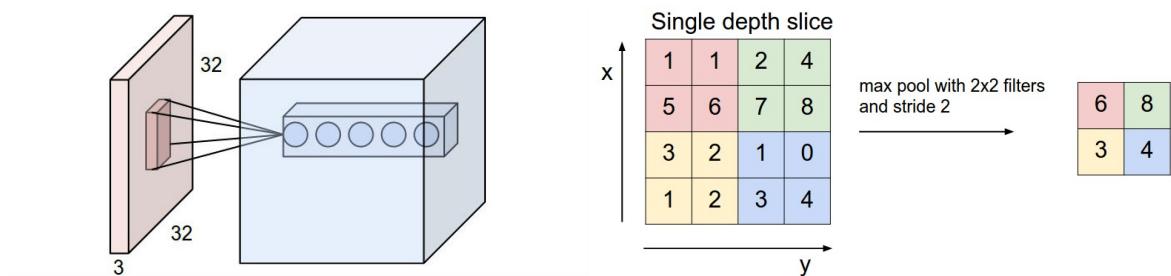


Figure 2.2: CNN Filter (left) [8], Max Pooling [9]

Pooling Layers

Pooling layers are used to reduce the number of parameters in a network. This is a much debated topic, Geoffrey Hinton's paper [10] discusses the problems of convolutional networks in his 2017 paper and especially focuses on the disadvantages of pooling layers. It is however a current, although fading, method of parameter reduction. For the purposes of optimising

network architecture for smaller hardware such as is portable on a bike this is still a recognised and useful method.

A pooling layer such as that depicted in (Figure 2.2) passes a kernel over the input (much like in a convolutional layer) and at each spacial location it calculates a metric. For example, max pooling takes the largest value contained within the kernel grid at each consecutive spacial location. This process reduces the input size with respect to the kernel size and the implemented stride.

2.1.3 Image Detection

Image detection is an area of research that uses the CNN networks described in Section 2.1.2. This is, however, a more difficult task as there can be multiple objects/classes of interest within one image. A naive approach to this problem would be to slide a window across an image and classify at each location. When the network classifies, say a cat at a particular window location it is likely that a cat is located in that area.

Although the aforementioned approach is valid it would be very slow, not only would it be needed to slide the window over the whole image, the size of the window would also need to be varied. The network cannot know the area of the photo that the object occupies and therefore what size to make the window. There are a number of differing approaches to this problem R-CNN [11] and other later iterations vs YOLO (you only look once) [12]. The R-CNN approach is very similar to the vanilla, naive solution, discussed above (commonly known as exhaustive search) and is used by later techniques for instance segmentation.

R-CNN

R-CNN, proposed by Ross Girshick et al. [11] in 2014, incorporates a region proposal method. This process produces 2000 proposals which are then passed through a large CNN, and each region is classified using class-specific linear SVMs. This process equates to a 53.7% mean average precision (mAP) compared to 35.1% by state-of-the-art methods at the time. In addition to this performance gain the system doesn't rely on far more complicated ensemble methods.

The region proposal method used is called selective search [13], this uses a fast segmentation algorithm by Felzenszwalb [14] followed by a greedy algorithm to iteratively group regions together using similarity metrics. The four metrics proposed in the paper are s_{colour} , $s_{texture}$, s_{size} , s_{fill} the aim is to use a number of metrics to pick up the differences between regions during the growth period and also limit one region's growth at the expense of others (hence the size metric).

The results of the region proposal (Figure 2.3) are a series of windows that are likely to contain an object. These are then simply warped to 227×227 pixels to fit the input dimension's of the CNN network that follows, containing five convolutional layers and two fully connected layers. The output of the CNN is then fed as a 4096 dimensional feature vector to a set of trained SVMs, one for each class, that are used to score the likelihood for that class. This approach is slow at 13s/image on a GPU and 53s/image on a CPU, the separation of a CNN and SVM is peculiar and adds to this timescale. A later iteration of this method Fast-RCNN is discussed in Section 2.1.3.

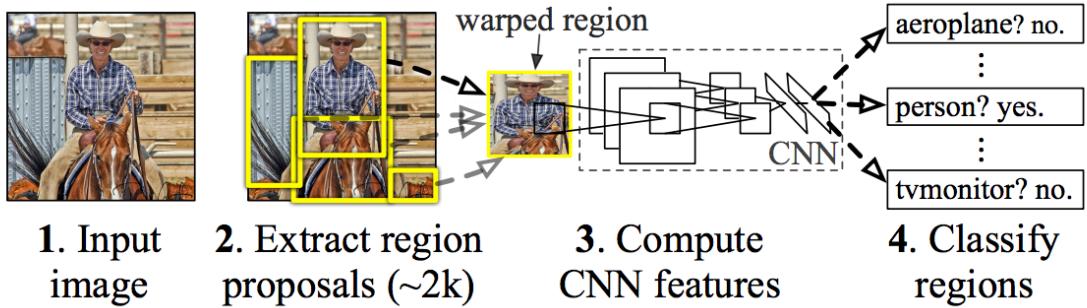


Figure 2.3: R-CNN Region Proposal & Architecture [11]

Fast-RCNN

R-CNN performs a CNN forward pass for each region proposal and does not share computation between regions [15]. Spacial pyramid pooling networks (SPPnet) [16] are designed to mitigate this problem by creating a feature map (which is always the same shape) from an input image and then classifying each object proposal from a vector extracted from the shared feature map. This paper still uses a multiple-stage classification pipeline but is a marked improvement on R-CNN with regards to speed.

Fast-RCNN uses the SPPnet concept but also incorporates a single-stage pipeline. The process inputs an image and multiple RoIs into a fully convolutional network (Figure 2.4). Each ROI is then pooled into a fixed-size feature map as in SPPnet and then mapped to a feature vector by fully connected layers. As can be seen in Figure 2.4 the network then branches into 2 separate outputs, one 4-point regression output of the bounding box location and a K-size classification output, where K is the number of object classes. The model therefore must contain a multi-task loss function that accounts for both bounding box and softmax loss (Equation 2.1).

$$L(p, u, t^u, v) = L_{cls}(p, u) + \lambda[u \geq 1]L_{loc}(t^u, v) \quad (2.1)$$

in which $L_{cls}(p, u) = -\log(p_u)$ is the log-loss for true class u , and L_{loc} is the bounding box regression loss defined as:

$$L_{loc} = \sum_{i \in \{x, y, w, h\}} smooth_{L_1}(t_i^u - v_i), \quad (2.2)$$

in which

$$smooth_{L_1} = \begin{cases} 0.5x^2 & \text{if } x < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases} \quad (2.3)$$

is a L_1 loss. This architecture allows for fast training with the forward pass of the network and therefore processes images $146\times$ faster than R-CNN due to sharing computation between layers and the improved pipeline design.

Figure 2.4 depicts the network architecture, this incorporates an ROI pooling layer which maps the convolutional feature map from an ROI (created by selective search) to a fixed extent feature map of a set $H \times W$, where H and W are layer hyperparameters. Each ROI is a rectangular window defined by a tuple of four elements (r, c, h, w) , where (r, c) specify

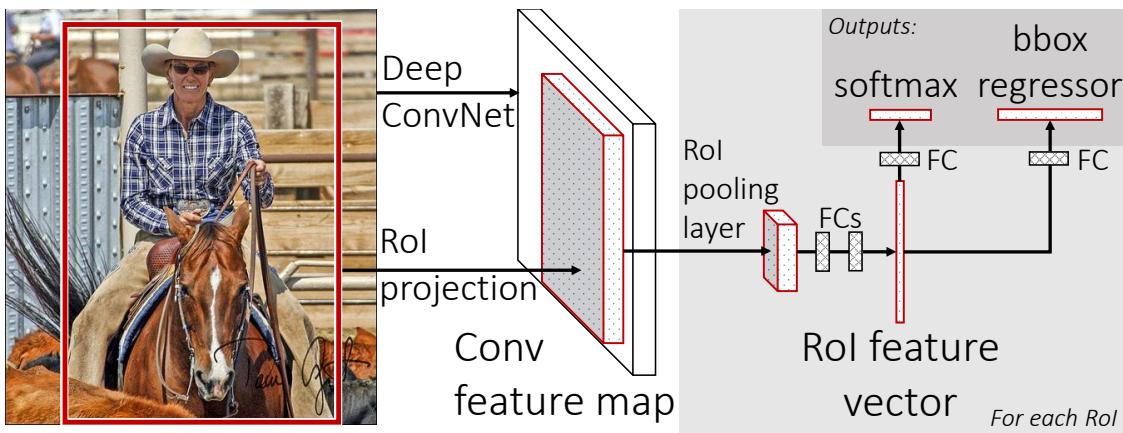


Figure 2.4: Fast R-CNN Architecture [15]

the top left corner and (h, w) the width. ROI max pooling then divides the ROI window into a grid and then max-pools each grid into the corresponding output cell. The size of the grid is proportionate to the ROI and so the output shape is constant regardless of ROI input size.

Faster R-CNN

Although the previous method [15] removed the need for a multi-stage classification pipeline (CNN and SVM), it still relies on the selective search [13] method for region proposal. This part is now the computational bottleneck that slows down the process. Faster R-CNN proposed by Ren et al. [17] seeks to alleviate this issue by combining region proposal into the CNN. The paper suggests the use of a region proposal network (RPN) to produce the ROIs used later in the network.

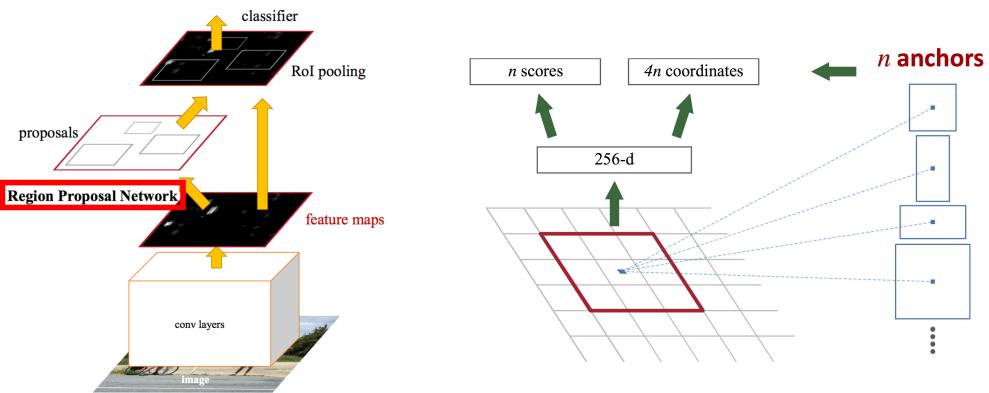


Figure 2.5: Fast R-CNN Architecture [15]

The region proposal network (RPN) outputs a set of rectangular object proposals, each of these proposals has an objectiveness score i.e. how likely the network thinks the proposal contains an object. This part of the network comes after a number of convolutional layers (in this paper 5 or 13 such layers). The region proposals are then created off the feature map from these layers by sliding a smaller network of size $n \times n$ over the convolutional feature map (Figure 2.5). At each sliding window location a number of region proposals are predicted simultaneously. These consist of 3 aspect ratios and 3 scales, therefore 9 total region

predictions or anchors at each location. The network is then followed by two sibling 1×1 convolutional layers, one for *reg* (region proposal) and one for *cls* (probability of object or not). The *reg* layer outputs $4k$ results, which represents the encoding coordinates of the k boxes and the *cls* later scores the probability that there is an object or no object in each proposal, therefore *cls* has $2k$ outputs. This results in a convolutional layer output of depth 9.

The loss for the network is calculated from two metrics, the output of the *reg* layer is a binary, i.e. positive example, or not. A positive example is defined by how much the box overlaps with the ground truth box of the object and so an IoU (Intersection-over-union) overlap > 0.7 is considered a positive example and an IoU < 0.3 is considered a negative label. For the first positive case, if no positive examples are found, the boxes with highest IoU are labelled positive. With this it follows that:

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*) \quad (2.4)$$

Where i is the anchor index, p_i is the predicted probability of anchor i containing an object. The ground truth label p_i is denoted as a 1 for a positive anchor and 0 for a negative therefore L_{reg} is only activated for a positive example, t_i and t_i^* represent the 4 coordinates of the bounding boxes, and L_{cls} is the classification loss. The two terms are normalised (N_{cls} and N_{reg}) and balanced by a weighting parameter λ . The final term to discuss, the bounding box regression is calculated, simply, as follows:

$$\begin{aligned} t_x &= (x - x_a)/w_a, & t_y &= (y - y_a)/h_a, \\ t_w &= \log(w/w_a), & t_h &= \log(h/h_a), \\ t_x^* &= (x^* - x_a)/w_a, & t_y^* &= (y^* - y_a)/h_a \\ t_w^* &= \log(w^*/w_a), & t_h^* &= \log(h^*/h_a) \end{aligned} \quad (2.5)$$

where x , y , w , and h are the box centre coordinates plus width and height. Predicted box, and anchored box are denoted by x , x_a , x^* (likewise for y , w , h).

Finally after the regions have been proposed, c.6000 per image, some regions overlap with others, and so, to reduce redundancy non-maximum suppression is used on their *cls* score with an IoU threshold of 0.7 to leave around 2000 proposals per image. Non-maximum suppression is a technique that sets all neighbouring areas to zero around a local maximum and hence only keeps the maximum value and suppresses the others.

2.1.4 Image Segmentation

Unlike the image detection problem described in Section 2.1.3, this task aims to segment an image on a pixel-by-pixel basis into respective classes (Figure 2.6). A conventional classification network takes an image and outputs a prediction of which class the image is i.e. a vector of probabilities or a ‘score’. This process uses fully-connected layers to produce the final output vector which throws away important spacial information used in a segmentation task. Therefore Long et al. [18] propose converting these final layers to convolutional layers, maintaining the spacial information and creating a heatmap output (Figure 2.6).

The heatmaps produced from the network then need to be upsampled from their respective dimensions to that of the output image. This process is essentially a reverse or deconvolution. Upsampling, therefore, with factor f is convolution with a fractional input stride of $1/f$.



Figure 2.6: Image Segmentation [19], Convolutional Heat Map [18]

Essentially each filter is placed on the output image and then multiplied by the input pixel. Hence the image increases in size depending on stride and kernel size as one pixel value is interpolated over a grid. The final problem therefore is how to choose the filter values, this can be done in a number of ways but the main accepted approach is to simply learn them as part of the tuning process.

An important concept is that of patchwise training which feeds the network a number of patches from an input image (small patches surrounding objects of interest) instead of an entire input image. This both helps to balance the classes, ensures the input has enough variance, and is a correct representation of the input set. However, this paper argues that this can be done from a fully convolutional training regime with incorporating a DropConnect-like mask [20] between the output and the loss. A DropConnect mask is similar to the proposed regularisation technique of Dropout [21] except that it randomly sets the network weights to 0 rather than the activation functions. This, on the output layer (connected to the network loss) is found to have the same effect as patch-wise sampling in Long et al.'s paper.

Although conventional proven classification architectures performed to state-of-the-art when modified for segmentation (56.0 mean IU), Long at al. [18] go on to design a bespoke network that achieves 62.7 mean IU. Mean pixel intersection over union (Mean IU) is a standard performance metric where the mean is taken over all classes, including the background. The paper describes a network architecture that uses skips [22]. Skips allow predictions from lower-level coarser layers to interact with finer, latter layer predictions. This lets the model make local predictions (finer) within the context of the global structure (coarser). This addition made for a marked increase in performance by simply up-sampling earlier layers and summing with latter layer outputs for a final image prediction.

2.1.5 Instance Segmentation

Instance segmentation is a problem that combines research in Sections 2.1.3 & 2.1.4, it aims to not only perform pixel-wise segmentation of classes but also determine instances of objects within an image. As discussed in Section 2.1.5 there is a simple method of combining these networks into one simple classification pipeline.

Mask R-CNN

Mask R-CNN was proposed by Facebook AI Research (FAIR) namely Kaiming He et al. in March 2017. It uses the Faster R-CNN network described in Section 2.1.3 which builds from

advances in Sections 2.1.3 & 2.1.3. The method works by simply adding a branch from the existing network for predicting segmentation masks as in Section 2.1.4 (Figure 2.7). The branch is a small FCN that is applied to each ROI for segmentation. There are a number of slight modifications to each approach to produce the final Mask R-CNN that are documented below.

RoIs are produced from the RPN as in Faster R-CNN [17] then an $m \times m$ mask is predicted

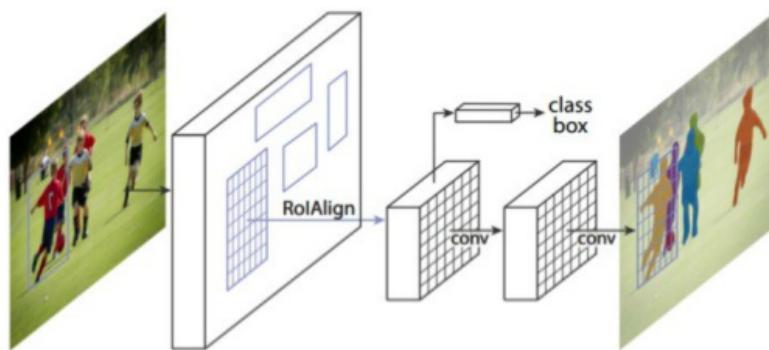


Figure 2.7: Mask R-CNN [23]

from each ROI without the use of a fully-connected (fc) layer and therefore keeps spacial information as in Section 2.1.4. It was found during experimentation that it is preferable to decouple the class classification and the mask production of the network. This stops multiple classes disrupting eachother during mask production, hence a mask is produced for each class. The class of each respective mask is then determined from the class output branch and that binary mask is selected.

The loss function for the network must now incorporate the three different branches; mask, class prediction, bounding-box prediction (Equation 2.6).

$$L = L_{cls} + L_{box} + L_{mask} \quad (2.6)$$

Where L_{cls} & L_{box} are identical to Section 2.1.3. The mask branch loss is determined by a per-pixel sigmoid between the ground-truth class (k) and the corresponding k th binary mask (excluding other classes) and is calculated with an average binary cross-entropy loss.

RoI pool as discussed in Section 2.1.3 is modified in this approach. Instead of subdividing an ROI into discrete regions e.g. a grid of 4×5 cells. An ROI does not always divide equally into these grid cells and so rounding of grid sizes is needed, this quantisation of sub-windows can cause misalignments of the ROI and extracted features. The new RoIAlign layer introduced in this paper mitigates this problem by not rounding grid-cells. Mask R-CNN uses bi-linear interpolation (linear interpolation with 2d) to compute the exact values of input features at four regularly samples locations in each ROI window.

2.1.6 Recurrent Neural Networks:

2.1.7 Bayesian Optimisation Layer Dimensions:

In order to get the best network architecture and hyper-parameters such as learning rate it is possible to naively try combinations and record which performs the best. However, when trying to optimise networks efficiently it is better to provide a range of possible working parameters and optimise them in an informed and, more importantly, automated manner. Bayesian optimisation allows such a method and can quickly check the parameter space for a good combination to minimise the network loss. Bayesian optimisation uses an acquisition function which evaluates a Gaussian process. It is cheaper to evaluate the proxy function that approximates the true function rather than run the neural network. Hence it is a very fast and effective method of optimisation.

Gaussian Processes: In a typical regression sample, the aim is to determine the parameters to a function (or family of functions) that best describe some data. A Gaussian process involves making the assumption that a function can be described simply by an infinite number of points along a function $f(x) = [f_1, f_2, \dots, f_\infty]$, a finite number of which are gaussian distributed. The aim is to then best approximate the relationship between these points via a covariance kernel $k(\cdot, \cdot)$ and a mean function $m(\cdot)$. A gaussian distribution is described by Equation 2.7 and a kernel function (Matern) by Equation 2.8.

$$p(x|\mu, \Sigma) = (2\pi)^{\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} e^{(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu))} \quad (2.7)$$

$$k_{Mat,3/2}(x_i, x_j) = \sigma_f^2 \left(1 + \frac{\sqrt{3}||x_i - x_j||}{l} \right) e^{(-\frac{\sqrt{3}||x_i - x_j||}{l})} \quad (2.8)$$

where l describes the smoothness of the data and σ_f the amplitude of the latent function. An example of a matern kernel fitted to some data is given in Figure

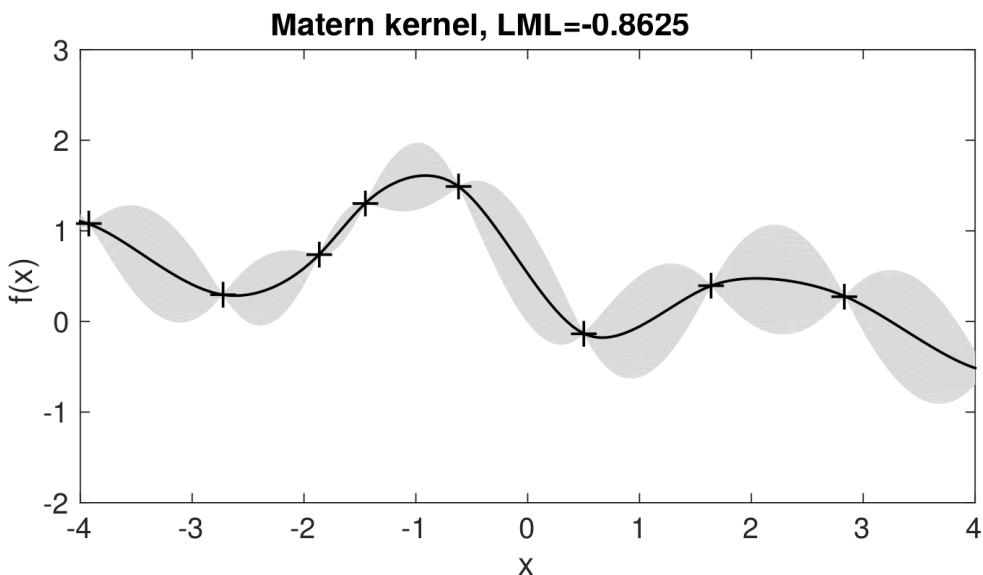


Figure 2.8: Gaussian Process: Matern Kernel [24]

Bayesian Optimisation: is a process that can maximise a set of hyper-parameters with a reduced evaluation of the true function. It uses an acquisition function that balances state exploration and exploitation i.e. how much to search local minima for the best solution vs exploring for global minima. The true function (NN model) is evaluated and data collected to initialise the process. Then the acquisition function determines where to evaluate next in the true function by querying the proxy function. An example acquisition function the lower confidence bound [25] is presented in Equation 2.9.

$$\alpha(\mathbf{x}) = -(\mu(\mathbf{x}) - (k)\sigma(\mathbf{x})) , k > 0 \quad (2.9)$$

The lower confidence bound balances exploration of global minima by the covariance function and weighting k , which increases at each iteration, with exploitation of local minima with the mean function. The Bayesian optimisation process roughly translates to:

Algorithm 1: Algorithm for Bayesian Optimisation

```

update_matched_roi (old_matched, new_matched, old_frame, new_frame)
Input :
  parameters                                     // A list of parameters to optimise
  parameterlimits                                // A list of ranges to optimise the parameters within
  function                                         // A pointer to the function to optimise
Output:
  parameters                                     // The optimal parameters

// Note that if there are no objects in old frame new_frame.id = zeros
D = [{x0, y0}]                                /* Initialise the dataset */
foreach step ∈ iterations do                  /* For every not-matched new ROI */
    update_gp(D[step − 1])      /* Update the Gaussian process with data from the last step */
    xstep = argmax( $\alpha(\mathbf{x})$ )    /* Evaluate the acquisition function for where to evaluate next
    */
    ystep = function(xstep)          /* Evaluate true function (NN model) */
    D.append({{xstep, ystep}})           /* Add data to dataset */
end
return(Xbest)                                /* Return optimal parameters */

```

2.2 Background, SLAM

In order to detect collisions between the bike and obstacles the software must have an understanding of the global environment. The detection described in Section 2.1 allows mapping, and tracking of dynamic objects in the frame, but distance, path prediction and eventually collision prediction rely strongly on an accurate understanding of 2.5-3D space.

2.2.1 Sensors

The sensors used in this project are an RGBD camera and an IMU. These are both available in the Intel Realsense Camera (ZR300, Figure 2.9). This camera has a USB cable that both provides the power and sends camera information. An extension cable allows for a laptop to be attached to the camera, from a backpack whilst cycling. Potential problems with this setup are due to juddering of the bike during operation; this could cause tricky footage to analyse and can data loss when the hard-drive is writing to disk. The latter issue can easily be solved by writing to an SSD while cycling.



Figure 2.9: Intel Realsense Camera (ZR300)

Camera Distortion

In order to fully understand the bike's location, the environment can be considered from two frames of reference, the world and camera frame. The camera, however, cannot be assumed to produce a perfect image and is hence prone to distortion - some types of camera distortion are shown in (Figure 2.10).

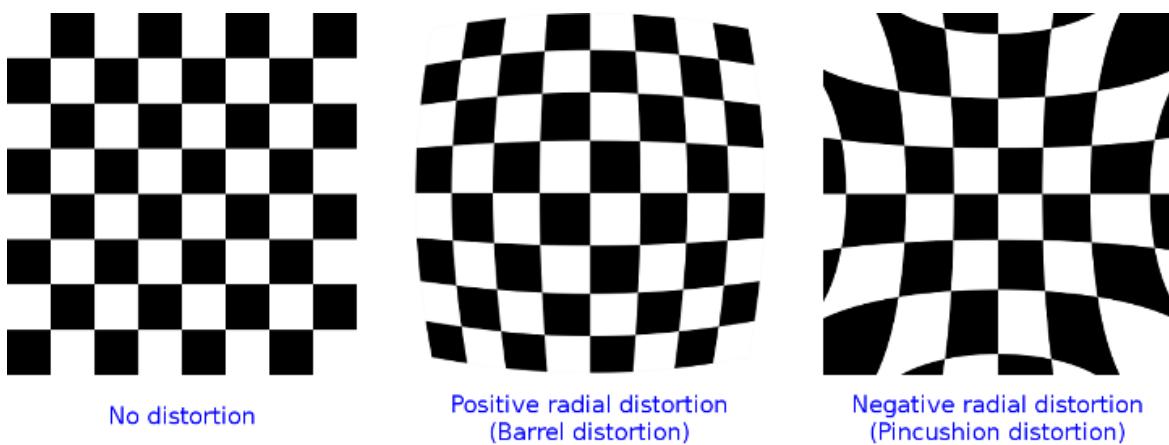


Figure 2.10: No Distortion (left)[26], Positive/negative radial distortion (middle/right)

A well-known distortion model is "radial-tangential" distortion. As the name implies there is both a radial and tangential component to the distortion (Equation 2.10)

$$r^2 = x_1'^2 + x_2'^2$$

$$\mathbf{x}'' = \mathbf{d}(\mathbf{x}') = \frac{1 + k_1 r^2 + k_2 r^4 + k_3 r^6}{1 + k_4 r^2 + k_5 r^4 + k_6 r^6} \begin{bmatrix} x_1' \\ x_2' \end{bmatrix} + \begin{bmatrix} 2p_1 x_1' x_2' + p_2(r^2 + 2x_1'^2) \\ p_1(r^2 + 2x_2'^2) + 2p_2 x_1' x_2' \end{bmatrix} \quad (2.10)$$

where x_1', x_2' is the 3d point mapped to the unit plane in z, $k_1 - k_6$ are radial distortion parameters, p_1, p_2 are tangential distortion parameters. The model parameters can be determined during camera calibration.

2.2.2 Camera to World Frame

Calculating the 3D Ray

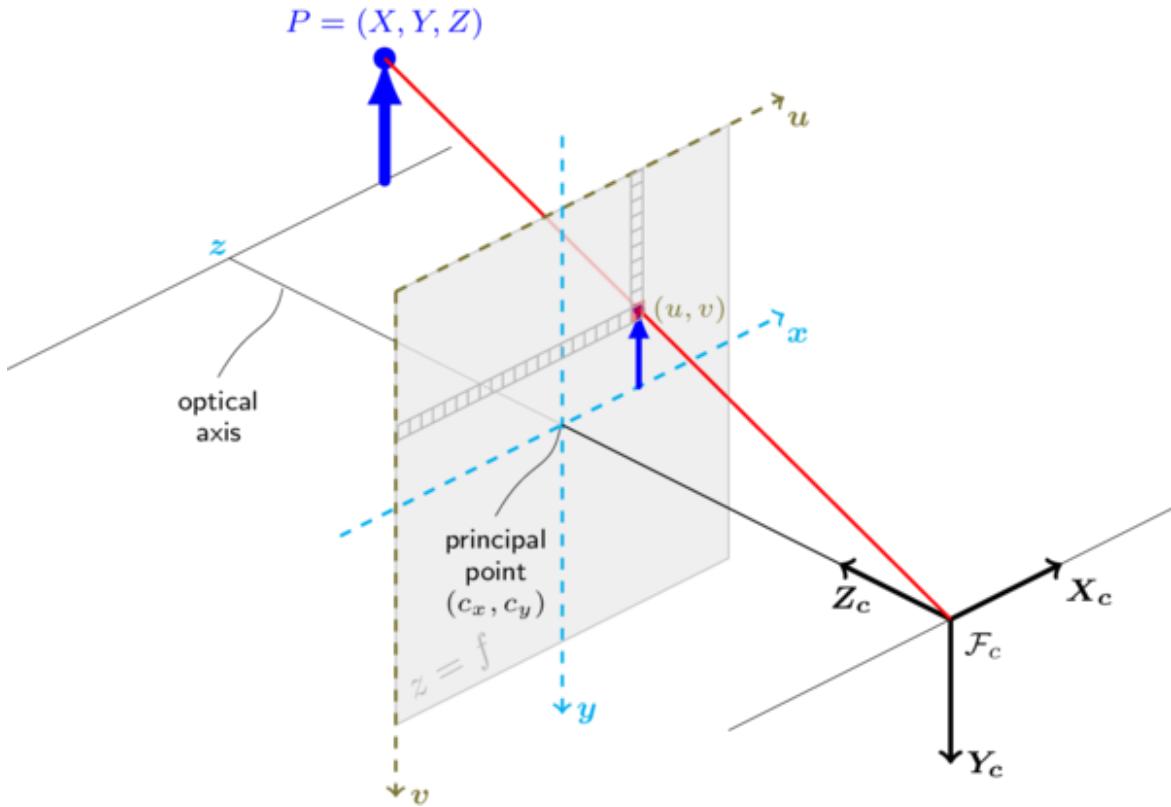


Figure 2.11: Pinhole Camera Model

Assuming an object or a point on the camera image, it is important to be able to transform such a point into the world frame. To track and gain a 3d appreciation of the environment. Hence for visual context (Figure 2.11) is a pinhole camera model. The only difference, in practice, in an intermediate un-distorting step to get the point's ray. This is due to the focal lens in the camera, which is not represented in the pinhole model. Hence to find a ray describing a 3d point the image must first be scaled, (Equation 2.11)

$$\mathbf{x}'' = \mathbf{k}^{-1}(\mathbf{u}) = \begin{bmatrix} \frac{1}{f_1} & 0 \\ 0 & \frac{1}{f_2} \end{bmatrix} \left(\mathbf{u} - \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} \right) \quad (2.11)$$

where f_1, f_2 are x/y focal lengths in pixels, and c_1, c_2 is the principal point (image centre) in pixels. The point can be projected to its undistorted location using the reverse of Equation 2.10.

$$\mathbf{x}' = \mathbf{d}^{-1}(\mathbf{x}'') \quad (2.12)$$

and finally the ray is calculated by remembering that the z component is of length 1.

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}' \\ 1 \end{bmatrix} \quad (2.13)$$

Calculating the 3D Location

Now that the ray has been determined we have a point on the 3D unit plane of $z = 1$. To find the 3D point in the camera frame, a depth measurement, obtained from the depth camera at the same (undistorted pixel location) is required. Hence the final 3D point is:

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}' \times \text{depth}(\mathbf{x}') \\ \text{depth}(\mathbf{x}') \end{bmatrix} \quad (2.14)$$

Assuming a static camera location, this would be all that is needed to map the space in-front of the camera. However, as the camera moves, the origin, and orientation of the camera frame changes with respect the world frame. Hence the software must not only determine the 3D location of a point in the camera frame but also translate that to a world frame representation.

2.2.3 Projecting to World Frame

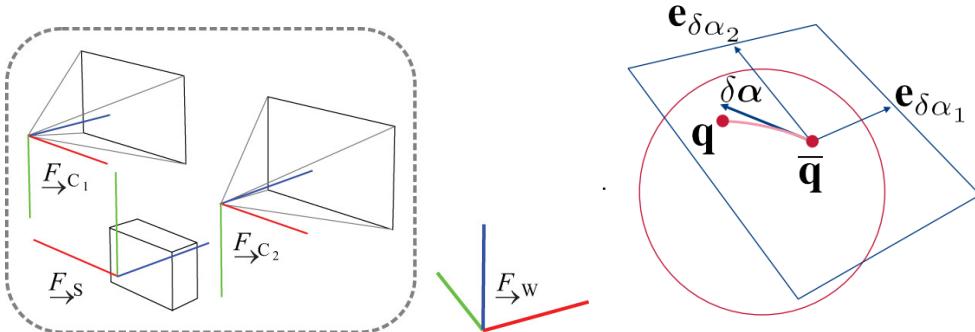


Figure 2.12: World vs Camera Frame (left)[27], Hamiltonian Quaternion (right)

To convert between world ($\overset{F}{\rightarrow}_W$), camera frame ($\overset{F}{\rightarrow}_{C_i}$), and IMU frame ($\overset{F}{\rightarrow}_S$), an appropriate method of describing the orientation, and location of an object in 3D space is needed. A Hamiltonian Quaternion is capable of correctly mapping the orientation of an object in 3D space, it does this by use of a 4D representation of one real part, and 3 imaginary;

$$\begin{aligned} q &= q_w + q_x i + q_y j + q_z k, \\ i^2 = j^2 = k^2 &= ijk = -1 \end{aligned} \quad (2.15)$$

With this extra representation of orientation it is possible to map the camera frame to the world frame. Given a position vector in the camera frame $p_c = [x, y, z, 1]$, it can be converted as follows:

$$\begin{aligned} T_{WC} &= T_{WS} \cdot T_{SC} \\ p_W &= T_{WC} \cdot p_C \end{aligned} \quad (2.16)$$

A transformation between frames is represented as \mathbf{T}_{BA} where B and A represent reference frames, so \mathbf{T}_{WC} converts between the world frame and camera frame. It is composed of a position and orientation transformation:

$$T_{WC} = \begin{bmatrix} C_{WC} & r_{WC} \\ 0_{1 \times 3} & 1 \end{bmatrix} \quad (2.17)$$

where $T_{WC} \in \mathbb{R}^{4 \times 4}$, the position translation matrix $r_{WC} \in \mathbb{R}^{3 \times 1}$, and the orientation matrix $C_{WC} \in \mathbb{R}^{3 \times 3}$. Note also that the inverse of this matrix can be easily calculated as C_{WC} is symmetric and therefore positive definite and it holds that $C_{CW} = C_{WC}^{-1} = C_{WC}^T$. Therefore:

$$T_{CW} = T_{WC}^{-1} = \begin{bmatrix} C_{WC}^T & -C_{WC}^T r_{WC} \\ 0_{1 \times 3} & 1 \end{bmatrix} \quad (2.18)$$

These transformation matrices must be determined from a SLAM algorithm such as OKVIS which is discussed in the following sections.

Camera to World Summary

Algorithm 2: Algorithm to transform camera point to world frame

cam2world (*cameraPoints*, *focalMatrix*, *imageCentre*, *T_{ws}*, *T_{sc}*, *cameraModel*)

Input :

<i>cameraPoints</i>	<i>// N points in the camera frame</i> $\in \mathbb{R}^{3 \times N}$
<i>focalMatrix</i> , <i>imageCentre</i>	<i>// Focal Length Matrix, Principal Points [x,y]</i>
<i>T_{ws}</i> , <i>T_{sc}</i>	<i>// IMU to World/Camera Matrix</i>
<i>cameraModel</i>	<i>// Camera Model</i> $\in \mathbb{R}^{4 \times 4}$

Output: *N* points in the world frame *worldPoints* $\in \mathbb{R}^{3,N}$

<i>x</i> = <i>scale_point</i> (<i>imageCentre</i> , <i>focalMatrix</i> , <i>cameraPoints</i>)	<i>/* Equation (2.11) */</i>
<i>x</i> = <i>undistort</i> (<i>cameraModel</i> , <i>x</i>)	<i>/* Equation (2.10) */</i>
<i>T_{wc}</i> = <i>matrix_multiply</i> (<i>T_{ws}</i> , <i>T_{sc}</i>)	
<i>x</i> = <i>add_depth</i> (<i>depth</i> , <i>x</i>)	<i>/* Equation (2.14), and Vstack [depth,1] to bottom of x */</i>
<i>worldPoints</i> = <i>transform_to_world</i> (<i>T_{wc}</i> , <i>x</i>)	<i>/* Equation (2.16) */</i>

2.2.4 World to Camera Frame

In order to correctly plot the Kalman filter and give other tracking algorithms the expected location on the picture the ability to transform from world to camera frame is needed. The process is similar to the camera to frame except the reverse, namely:

Algorithm 3: Algorithm to transform world point to camera frame

world2cam (*worldPoints*, *focalMatrix*, *imageCentre*, *T_{ws}*, *T_{sc}*, *cameraModel*)

Input :

<i>worldPoints</i>	<i>// N points in the camera frame</i> $\in R^{3 \times N}$
<i>focalMatrix</i> , <i>imageCentre</i>	<i>// Focal Length Matrix, Principal Points [x,y]</i>
<i>T_{ws}</i> , <i>T_{sc}</i>	<i>// IMU to World/Camera Matrix</i>
<i>cameraModel</i>	<i>// Camera Model</i> $\in R^{4 \times 4}$

Output: *N* points in the camera frame *cameraPoints* $\in \mathbb{R}^{3,N}$

<i>T_{wc}</i> = <i>matrix_multiply</i> (<i>T_{ws}</i> , <i>T_{sc}</i>)	
<i>T_{cw}</i> = <i>inverse</i> (<i>T_{wc}</i>)	/* Equation (2.18) */
<i>x</i> = <i>transform_to_camera</i> (<i>T_{cw}</i> , <i>worldPoints</i>)	/* Equation (2.16) */
<i>x</i> = <i>distort</i> (<i>cameraModel</i> , <i>x</i>)	/* Equation (2.10) */
<i>x</i> = <i>remove_depth</i> (<i>depth</i> , <i>x</i>)	/* Reverse Equation (2.14) */
<i>cameraPoints</i> = <i>scale_point</i> (<i>imageCentre</i> , <i>focalMatrix</i> , <i>x</i>)	/* Equation (2.11) */

2.2.5 Bicycle States

In order to find the bike location and orientation in each frame; a complete representation of the bike's state is needed. At each taken image (*k*) the bike can be assumed to have a set state \mathcal{X}_R^k . Contained within \mathcal{X}_R are a set of pose states, \mathcal{X}_T , and speed/bias states \mathcal{X}_{sb} , Equation 2.19. Landmarks are contained within the set \mathcal{X}_L .

$$\begin{aligned}\mathcal{X}_T &:= [\mathbf{wr}_S^T, \mathbf{q}_{WS}^T]^T \\ \mathcal{X}_{sb} &:= [\mathbf{s}v^T, \mathbf{b}_g^T, \mathbf{b}_a^T]^T \\ \mathcal{X}_R &:= [\mathbf{wr}_S^T, \mathbf{q}_{WS}^T, \mathbf{s}v^T, \mathbf{b}_g^T, \mathbf{b}_a^T]^T\end{aligned}\tag{2.19}$$

The elements \mathbf{wr}_S^T , \mathbf{q}_{WS}^T are the position and the quaternion body orientation, $\mathbf{s}v^T$ the velocity, \mathbf{b}_g & \mathbf{b}_a the accelerometer and gyroscope biases. Since $\delta\mathbf{q}$ describes the axis-angle perturbation about the estimate $\bar{\mathbf{q}}_{WS}$ it follows that $\mathbf{q}_{WS} = \delta\mathbf{q} \otimes \bar{\mathbf{q}}_{WS}$ "EXPAND THIS BIT)". Hence the error states combining with (Equation ??) are:

$$\begin{aligned}\delta\mathcal{X}_T &:= [\delta r^T, \delta\alpha^T]^T \\ \delta\mathcal{X}_{sb} &:= [\delta v^T, \delta b_g^T, \delta b_a^T]^T \\ \delta\mathcal{X}_R &:= [\delta r^T, \delta\alpha^T, \delta v^T, \delta b_g^T, \delta b_a^T]^T\end{aligned}\tag{2.20}$$

2.2.6 The SLAM Problem

If measurements were exact, mapping of an environment would be easy. This, however is not the case, there is often noise associated with each taken measurement and movement. For example distance measurements and IMU readings are subject both to a noise and a bias (calibration error) which can both be assumed to be Gaussian distributed.

$$\tilde{\mathbf{z}} = \mathbf{b}_c + s\mathbf{M}_z + \mathbf{b} + \mathbf{n} + \mathbf{o} \quad (2.21)$$

where \mathbf{z} is the correct measurement, \mathbf{b}_c is a long-term constant bias, s is a scaling factor, \mathbf{M} - misalignment, \mathbf{b} - time varying bias, n - noise, and o contains other un-modelled influences. This means that no measurement can be trusted exclusively for state-space information.

Historically, before the recent advances in SLAM, it was first assumed that errors have a compounding effect. The position of a bike, would become increasingly unknown as the environment is explored and the bike is displaced from the origin. This makes intuitive sense if say there is only IMU measurements available. At each step, more uncertainty is introduced from the noise and biases (shown in Equation 2.21) and the bike's position becomes more and more uncertain.

For some time it was also presumed to be the case even when external measurements of the environment are taken. This, however, is not the case as proved by Durrant-Whyte [28] since there is a high degree of correlation between estimates of different landmarks and locations in a map, these correlations also increase with successive observations (Figure 2.13). As an example of this (Figure 2.13 : Right) displays error in measurements of landmark observation, the mean error is common between all landmarks and so with successive landmark observations, the errors in location estimates become highly correlated. This means that although \mathbf{m}_i 's location may be quite uncertain but the relative location between two landmarks $\mathbf{m}_i - \mathbf{m}_j$ is almost certain.

In the case of the bicycle set-up, there is both an IMU, and a camera. These provide the basic inputs to the OKVIS software [27]. Okvis is able to determine the bike's location by optimising the total error, hence combining both the camera and IMU error.

Batch NonLinear Least Squares

Given an error function for the camera and IMU in the form:

$$J(\mathbf{x}) := \frac{1}{2} \sum_{n=1}^N \mathbf{e}_n(\mathbf{x})^T \mathbf{W}_n \mathbf{e}_n(\mathbf{x}) \quad (2.22)$$

where $\mathbf{e}_n(\cdot)$ is an error term that is weighted by \mathbf{W}_n , and would consist, in this case of the camera and IMU error. The only way to minimise this error is through observations which are also presumed to have an associated noise term.

$$\mathbf{z}_n = g_n(\mathbf{x}) + \mathbf{v}_n, \quad (2.23)$$

where z_n is the observation, $g_n(\mathbf{x})$ the observation model, and $v_n \sim \mathcal{N}(0, Q_n)$ and are independent. Hence as v_n is mean 0, the error can be simply formulated by:

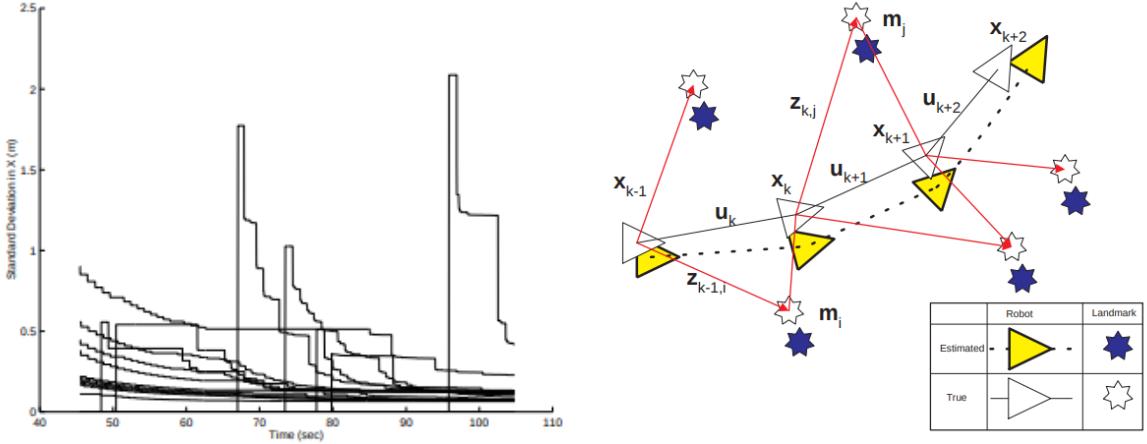


Figure 2.13: Left: Convergence in landmark uncertainty. Over time, standard deviations reduce monotonically to a lower bound [29]. Right: Error between estimated and true landmarks is common.

$$\mathbf{e}_n(\mathbf{x}) = z_n - g_n(\mathbf{x}), \quad \text{error} = \text{observation} - \text{predicted observation} \quad (2.24)$$

Since observation noise between the IMU and camera v_n are independent the covariance matrix is diagonal. Hence $W_n = E[\mathbf{e}_n \mathbf{e}_n^T]^{-1} = Q^{-1}$. The aim of least squares is to find the state vector \mathbf{x} that minimises the error function, hence the optimisation problem is posed as:

$$\mathbf{x} = \operatorname{argmin}(J(\mathbf{x})) \quad (2.25)$$

If the error term $\mathbf{e}_n(\mathbf{x})$ was linear w.r.t \mathbf{x} the solution would be as simple as setting the $\frac{\delta J(\mathbf{x})}{\delta \mathbf{x}}^T$ to zero and solving the respective simultaneous equations. However this is often not the case, and so the solution to the equation must be solved iteratively, with for example, gauss-newton formulation of the equations and a solver such as google-ceres.

Frame to Frame Matching

Keypoint Detection The Okvis [27] paper uses a corner detection method to determine keypoints and then matches these keypoints using a binary descriptor (BRISK [30]). In the BRISK paper a modification of the FAST feature extractor is used, however, in the OKVIS implementation a modified Harris Corner detection algorithm is utilised. A simplistic way to think of this is how one would go about determining a corner in an image. Imagine an edge, it is simply a location on an image where the brightness gradient changes rapidly in one direction (perpendicular to the edge). Hence a corner is a similar change but in two such directions. Imagine a location (x, y) and a small displacement from that location $(x + \Delta x, y + \Delta y)$. An area of rapid gradient change would have very different intensity values at (x, y) and surrounding points. Hence a corner will maximise the sum of squared differences between points.

$$f(x, y) = \sum_{x_k, y_k \in W} (I(x_k, y_k) - I(x_k + \Delta x, y_k + \Delta y))^2 \quad (2.26)$$

This can be further distilled by the Taylor expansion to:

$$f(x, y) = \sum_{(x,y) \in W} (I_x(x, y)\Delta x + I_y(x, y)\Delta y)^2 \quad (2.27)$$

or in matrix form:

$$\begin{aligned} M &= \begin{bmatrix} \sum_{(x,y) \in W} I_x^2 & \sum_{(x,y) \in W} I_x I_y \\ \sum_{(x,y) \in W} I_x I_y & \sum_{(x,y) \in W} I_y^2 \end{bmatrix} = R^T \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} R \\ f(x, y) &= (\Delta x \ \Delta y) M \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix} \end{aligned} \quad (2.28)$$

where M is the structure tensor (also commonly known as the Second Moment Matrix). The Harris corner response then is taken from the eigenvalues of M , intuitively since the gradient changes rapidly the second moment matrix eigenvalues are large in that direction Figure 2.14.

$$R = \det(M) - k\text{trace}(M)^2 = \lambda_1 \lambda_2 - k(\lambda_1 + \lambda_2)^2 \quad (2.29)$$

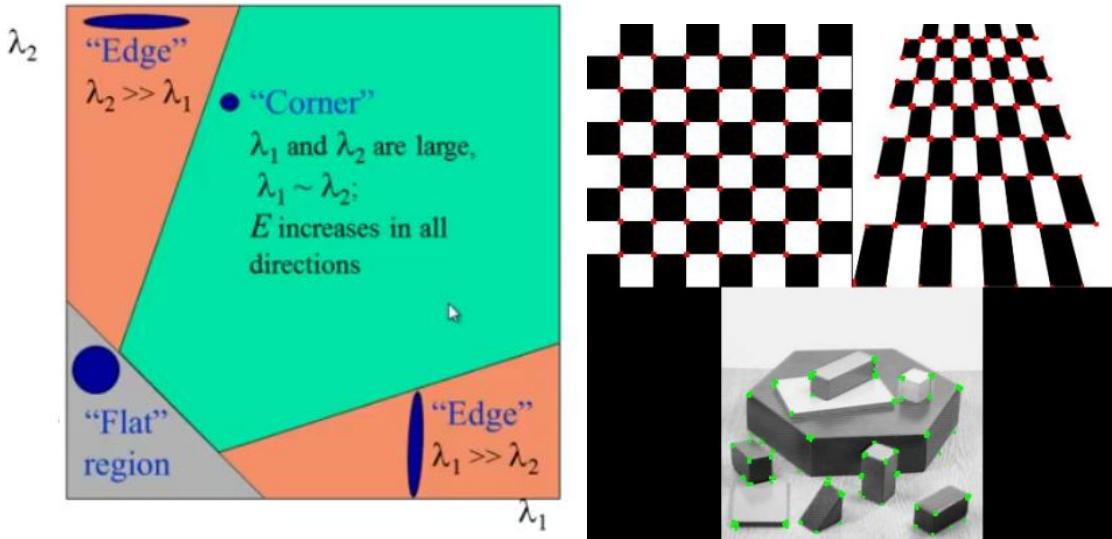


Figure 2.14: Left: Structure Tensor Eigenvalues, Right: Harris corner detection result [31]

Keypoint Matching Binary descriptors use a sampling pattern to create pairs of points (lines). The algorithm then compares the two points (p_1, p_2) and if $p_1 > p_2$ a 1 is recorded and otherwise 0. This creates a binary string, that along with orientation representation describe a keypoint. When matching keypoints between images a simple, and fast $\sum XOR(img1, img2)$ reveals keypoint similarity from which a threshold can be used to determine matches.

The BRISK sampling pattern (Figure 2.15: Right) uses a standard Gaussian smoothing at each sample point. The standard deviation of the sample is represented by the red circles in the figure. This ensures the sample is less prone to noise than other traditional binary detectors e.g. BRIEF.

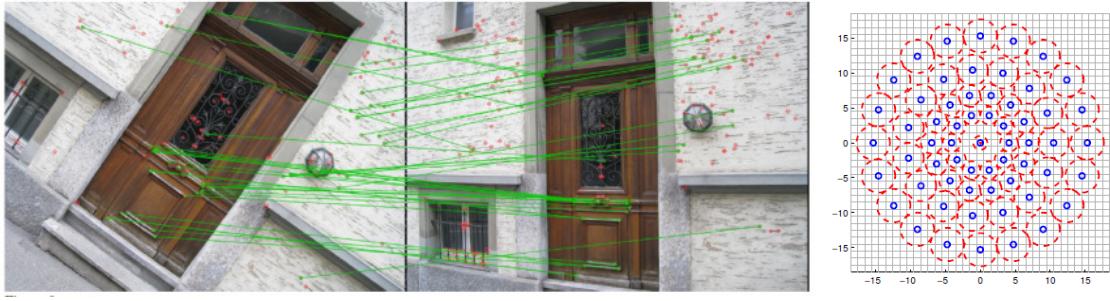


Figure 2.15: Left: Brisk keypoint matching example, Right: Brisk Sampling Pattern [30]

Camera Re-projection Error

As discussed in Section 2.2.4, the camera uses a lens to take a photograph. This, especially on cheaper cameras, can distort the image. The distortion is intrinsic to the camera and it needs to be calibrated to determine a distortion model $d(\cdot)$. This is then used to determine the re-projection error [32].

$$\mathbf{e}_r^{i,j,k} = \mathbf{z}^{i,j,k} - d_i(\mathbf{T}_{CiS}^k \mathbf{T}_{SW}^k w \mathbf{l}^j) \quad (2.30)$$

where $\mathbf{z}^{i,j,k}$ denotes the measurement image coordinates.

IMU Kinematics

Therefore the kinematics of the IMU can be modelled as:

$$\begin{aligned} {}_W\dot{\mathbf{r}}_S &= \mathbf{C}_{WS} {}_s\mathbf{v}, \\ \dot{\mathbf{q}}_{WS} &= \frac{1}{2}\Omega({}_S\tilde{\omega})\mathbf{q}_{WS}, \\ {}_S\dot{\mathbf{v}} &= ({}_S\tilde{\mathbf{a}} + \mathbf{w}_a - \mathbf{b}_a) + C_{SW} {}_W\mathbf{g} - ({}_S\omega) \times {}_S\mathbf{v}, \\ \dot{\mathbf{b}}_g &= \mathbf{w}_{bg}, \\ \dot{\mathbf{b}}_a &= -\frac{1}{\tau}\mathbf{b}_a + \mathbf{w}_{ba}, \end{aligned} \quad (2.31)$$

where:

$$\begin{aligned} \Omega({}_S\tilde{\omega}) &= \begin{bmatrix} -{}_S\omega \\ 0 \end{bmatrix}^\oplus \\ {}_S\omega &= {}_S\tilde{\omega} + \mathbf{w}_g - \mathbf{b}_g \end{aligned} \quad (2.32)$$

${}_W\dot{\mathbf{r}}_S$ represents the position derivative likewise ${}_S\dot{\mathbf{v}}$ for velocity, ${}_S\tilde{\mathbf{a}}$ accelerometer measurements, ${}_W\mathbf{g}$ earth's gravitational acceleration vector, and finally $\mathbf{w} := [w_g^T, w_a^T, w_{bg}^T, w_{ba}^T]^T$ are independent Gaussian white noise processes with zero mean.

IMU Kinematics Linearised

In order to optimise the cost function, it is important to linearise the IMU kinematics. The error dynamics take the form:

$$\delta\mathcal{X} \approx \mathbf{F}_c(\bar{x}_R)\delta\mathcal{X}_R + \mathbf{G}(\bar{x}_R)\mathbf{w} \quad (2.33)$$

$$\begin{bmatrix} \delta\dot{\mathbf{r}} \\ \delta\dot{\mathbf{a}} \\ \delta\dot{\mathbf{v}} \\ \delta\dot{\mathbf{b}}_g \\ \delta\dot{\mathbf{b}}_a \end{bmatrix} = \begin{bmatrix} 0_{3 \times 3} & [\bar{C}_{WS}\bar{v}]^\times & \bar{C}_{WS} & 0_{3 \times 3} & 0_{3 \times 3} \\ 0_{3 \times 3} & 0_{3 \times 3} & 0_{3 \times 3} & -\bar{C}_{WS} & 0_{3 \times 3} \\ 0_{3 \times 3} & -\bar{C}_{WS}[W\mathbf{g}]^\times & -[S\bar{v}]^\times & -[S\bar{v}]^\times & -I_3 \\ 0_{3 \times 3} & 0_{3 \times 3} & 0_{3 \times 3} & 0_{3 \times 3} & 0_{3 \times 3} \\ 0_{3 \times 3} & 0_{3 \times 3} & 0_{3 \times 3} & 0_{3 \times 3} & -\frac{1}{\tau}\mathbf{I}_3 \end{bmatrix} \begin{bmatrix} \delta\mathbf{r} \\ \delta\mathbf{a} \\ \delta\mathbf{v} \\ \delta\mathbf{b}_g \\ \delta\mathbf{b}_a \end{bmatrix} + \begin{bmatrix} 0_{3 \times 3} \\ \bar{C}_{WS}\mathbf{W}_g \\ \bar{C}_{WS}\mathbf{W}_a \\ \mathbf{W}_{bg} \\ \mathbf{W}_{ba} \end{bmatrix} \quad (2.34)$$

This formulation of the error dynamics is then used to derive the IMU error term (Equation 2.35), for further details see [27].

$$\mathbf{e}_s^k(\mathbf{x}_R^k, \mathbf{x}_R^{k+1}, \mathbf{z}_s^k) = \begin{bmatrix} {}^W\hat{\mathbf{r}}_S^{k+1} - {}^W\mathbf{r}_S^{k+1} \\ 2[\mathbf{q}_{WS}^{k+1} \otimes \mathbf{q}_{WS}^{k+1-1}] \\ \hat{\mathbf{x}}_{sb}^{k+1} - \mathbf{x}_{sb}^{k+1} \end{bmatrix} \quad (2.35)$$

2.3 Background Object Tracking

Object tracking has been recently advanced by means of discriminative learning methods, these methods try to distinguish between the object and the environment that it is in. The model must therefore be fed negative samples to understand the relevant environment.

2.3.1 Kernelised Correlation Filter

In the recent paper by Henriques et al. [33] the tools are developed to supply thousands of negative samples without iterating over them explicitly. The paper makes use of circulant matrices to provide negative samples. These are obtained by translating the data using a *cyclic shift operator*, which in the 1D signal case is:

$$P = \begin{bmatrix} 0 & 0 & \dots & 0 & 1 \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix} \quad (2.36)$$

Therefore the product of $P\mathbf{x} = [x_n, x_1, x_2, \dots, x_{n-1}]$ and any shift can be modelled using $P^u\mathbf{x}$. Because of its cyclic nature, the signal \mathbf{x} repeats when $u = n$ and so the full set of shifted signals is obtained with $P^u\mathbf{x}|u = 0, \dots, n-1$. An example of this is shown in the 1D case by Figure 2.16.

Therefore to compute the regression with shifted samples, the data matrix X takes the form:

$$X = C(\mathbf{x}) = \begin{bmatrix} x_1 & x_2 & \dots & x_{n-1} & x_n \\ x_n & x_1 & \dots & x_{n-2} & x_{n-1} \\ x_{n-1} & x_n & \dots & x_{n-3} & x_{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_2 & x_3 & \dots & x_n & x_1 \end{bmatrix} \quad (2.37)$$

Noting that all circulant matrices can be diagonalised by the Discrete Fourier Transform (DFT) X can be written as:

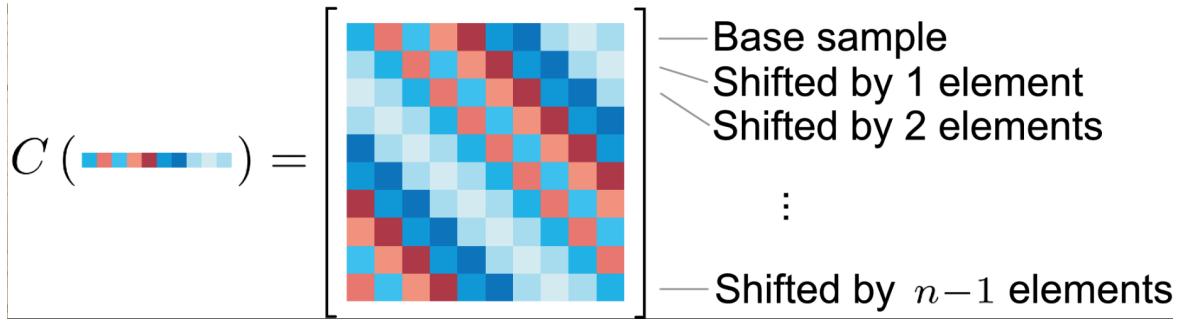


Figure 2.16: Illustration of circulant matrix.[33]

$$\begin{aligned} F(\mathbf{z}) &= \sqrt{n}F\mathbf{z} \\ \hat{\mathbf{x}} &= F(\mathbf{x}) \\ X &= F \text{diag}(\hat{\mathbf{x}}) F^H \end{aligned} \tag{2.38}$$

This is possible because the DFT is a linear operation.

Kernels

Using a kernel, the model can implicitly use a high-dimensional feature space. The kernel function used to do this (e.g. Gaussian), Equation 2.39 creates an $n \times n$ kernel matrix $K_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2\right) \tag{2.39}$$

and therefore the function becomes:

$$f(\mathbf{z}) = \mathbf{w}^T \mathbf{z} = \sum_{i=1}^n \alpha_i k(\mathbf{z}, \mathbf{x}_i) \tag{2.40}$$

with the solution of kernelised Ridge Regression given by:

$$\alpha = (K + \lambda I)^{-1} \mathbf{y} \tag{2.41}$$

where K is the kernel matrix and α is the vector of coefficients α_i . Given a kernel function that combines data through a commutative operation, it can be proven that given circulant data the kernel matrix K is also circulant and therefore the DFT trick can be used. This leads to a diagonalisation of Equation 2.41 to Equation 2.42.

$$\hat{\alpha} = \frac{\hat{\mathbf{y}}}{\hat{\mathbf{k}}^{\mathbf{x}\mathbf{x}} + \lambda} \tag{2.42}$$

where $\hat{\mathbf{k}}^{\mathbf{x}\mathbf{x}}$ is the first row of the kernel matrix $K = C(\hat{\mathbf{k}}^{\mathbf{x}\mathbf{x}})$, and $\hat{\mathbf{a}}$ denotes the DFT of a vector. This means that only an $n \times 1$ kernel matrix needs to be computed, due to the circulant kernel property and DFT. This makes the method far faster than conventional kernel methods that require an $n \times n$ sized kernel matrix.

Detection

The basic premise is that the α term in Equation 2.42 is now initialised to a learnt set of coefficients that correctly describe the object. Therefore, in order to detect the object of interest the regression function $f(\mathbf{z})$ must be evaluated at a number of candidate patches. These patches can themselves also be described by circulant shifts. Defining the kernel matrix K^z between cyclic shifts of base samples x and z , with elements $k(P^{i-1}z, P^{j-1}x)$. The cyclic property, as described in Section 2.3.1 allows only the first row of the kernel matrix to be relevant, and hence :

$$K^z = C(\mathbf{k}^{\mathbf{x}\mathbf{z}}) \quad (2.43)$$

where k^{xz} is the kernel correlation of x and z . Finally the regression equation can be computed for all candidate patches by:

$$f(z) = (K^z)^T \alpha \quad (2.44)$$

This can be further sped up by the correct choice of kernel function and a clever kernel trick to modify the computational cost of kernel computation to $O(n \log n)$.

Chapter 3

Contribution

3.1 Depth Measurement Outlier Removal

In order to get the depth measurement of an object, the mask output from Mask R-CNN is overlayed onto the depth map and a mean depth is taken. A number of points are not detected and therefore large quantities of the depth map is 0. Therefore, the median depth is selected only from coordinates at which $depth > 0$. Figure 3.1 shows some examples of a depth camera image taken in the Imperial Library.

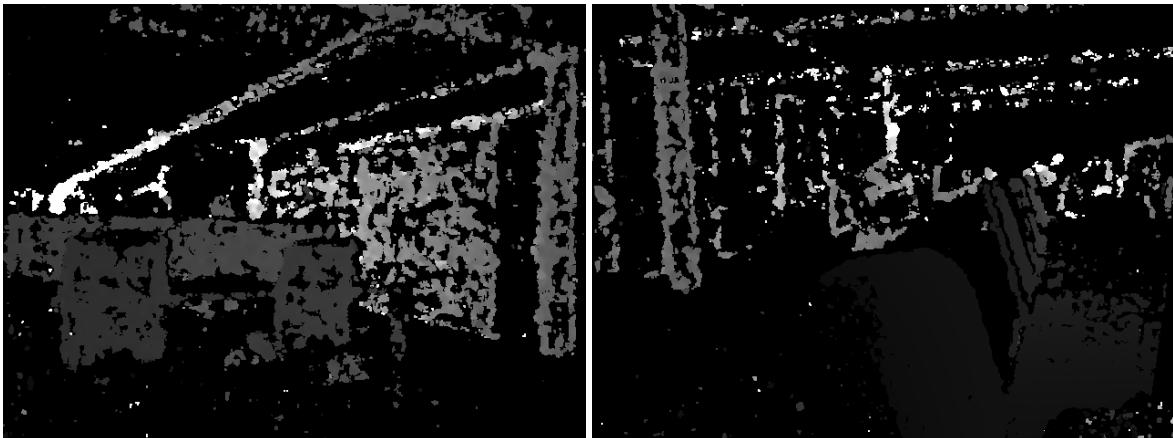


Figure 3.1: Depth Camera Output

As can be seen from Figure 3.1 the data is sparse and unreliable for objects at a distance. This is less so as the depth is scaled [0-255] for a picture and so some measurements seem less reliable than they are but it is still a challenge. Often when the depth measurements are incorrect an unfeasibly large value is returned $> 50m$. This is a real problem for tracking objects, not necessarily after a period of successful tracking as the value can be discarded as an outlier but on initialisation or after only a couple of data points. Not only is the depth incorrect but that problem is in both x and y when transformed to the world frame. An example of this phenomenon is displayed in Figure 3.2 where a bottle is tracked from frame to frame, which also demonstrates the ability of the tracking method to be non-cycling domain specific.

This can occur in a number of instances and far away objects are most prone as there is little

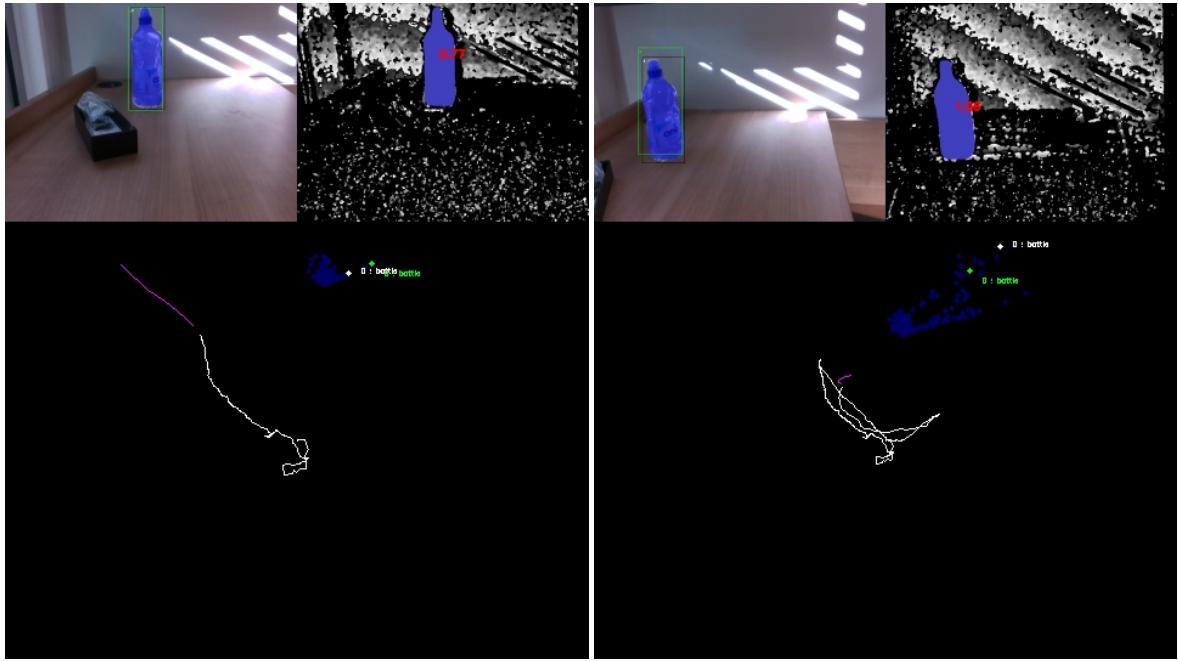


Figure 3.2: Left: Bottle with working depth, Right: Incorrect Depth measurement. In both cases the white dot is the bottle's location, the green dot is the Kalman filter prediction, blue dots are previous locations (the bottle did not move), and the white line is the camera location.

area covered on the depth image to measure the distance from. However, since the object is defined by a bounding box, prior knowledge of the classes can be used to detect such an outlier. Intuitively if two objects are the same size in an image and the first is closer than the second, then the latter object must be larger than the first. Hence using the camera to world model, if the depth is incorrectly large the height of the object in the world frame will be far larger than is known to be possible i.e. a 10m tall human.

Given the knowledge that human sizes roughly vary from a 4ft child to a 7ft adult, poor depth predictions convert people to outside of that range. If that is the case then it must be modified to an approximate value, say the depth of an averagely sized human 1.5m. This logic can be extended to other objects such as cars and bikes. Therefore, assuming a bounding box of $[x_{tl}, y_{tl}, x_{br}, y_{br}]$, where tl denotes top-left corner and br is the bottom right, the matrix P can be formed.

$$\mathbf{P} = \begin{bmatrix} x_{tl} & x_{tl} \\ y_{tl} & y_{br} \end{bmatrix} \quad (3.1)$$

Completing the initial steps of Cam2World (Section 2.2.2) until Equation 2.14, and adding the 1 required for the proceeding transform, the process yields:

$$\mathbf{K} = \begin{bmatrix} k_{(1,1)}d & k_{(1,2)}d \\ k_{(2,1)}d & k_{(2,2)}d \\ d & d \\ 1 & 1 \end{bmatrix} \quad (3.2)$$

Where $[k_{(1,1)}, k_{(2,1)}]$ are from $[x_{tl}, y_{tl}]$, $[k_{(1,2)}, k_{(2,2)}]$ from $[x_{tl}, y_{br}]$, and $d = \text{depth}(\mathbf{x}')$ from the depth image. For the next steps, the transformation matrix T_{WC} is used to calculate the world position.

$$\mathbf{T}_{WC} = \begin{bmatrix} \dots r_1 \dots \\ \dots r_2 \dots \\ \dots r_3 \dots \\ \dots r_4 \dots \end{bmatrix} \quad (3.3)$$

where r_i denotes row i of T_{WC} , and $r_4 = [0_{3 \times 1}, 1]$. Therefore replicating a simple matrix multiplication in Equation 3.4, the row r_3 is the only row from \mathbf{T}_{WC} that is needed for the object height.

$$\text{WorldPoints} = \begin{bmatrix} x_{w,tl} & x_{w,bl} \\ y_{w,tl} & y_{w,bl} \\ z_{w,tl} & z_{w,bl} \end{bmatrix} = \mathbf{T}_{WC} \mathbf{K} = \begin{bmatrix} \dots r_1 \dots \\ \dots r_2 \dots \\ \dots r_3 \dots \\ \dots r_4 \dots \end{bmatrix} \begin{bmatrix} k_{(1,1)}d & k_{(1,2)}d \\ k_{(2,1)}d & k_{(2,2)}d \\ d & d \\ 1 & 1 \end{bmatrix} \quad (3.4)$$

Hence the height of the object $z_{w,tl} - z_{w,bl}$ is obtained from:

$$\begin{aligned} z_{w,tl} - z_{w,bl} &= r_3 \mathbf{K}[:, 1] - r_3 \mathbf{K}[:, 2] \\ z_{w,tl} - z_{w,bl} &= r_3 (\mathbf{K}[:, 1] - \mathbf{K}[:, 2]) \end{aligned} \quad (3.5)$$

where $\mathbf{K}[:, 1]$ defines the first column of \mathbf{K} , and defining vector A as:

$$Ad = \mathbf{K}[:, 1] - \mathbf{K}[:, 2] = \begin{bmatrix} k_{(1,1)}d - k_{(1,2)}d \\ k_{(2,1)}d - k_{(2,2)}d \\ 0 \\ 0 \end{bmatrix} = d \begin{bmatrix} k_{(1,1)} - k_{(1,2)} \\ k_{(2,1)} - k_{(2,2)} \\ 0 \\ 0 \end{bmatrix} \quad (3.6)$$

then an approximation of the depth, given the average object height o_{av} is:

$$depth \approx \frac{o_{av}}{r_3 A} \quad (3.7)$$

Once this depth approximation is found the relevant models can get a depth update that is approximately correct unlike the actual camera depth measurement. This allows better modelling and tracking of object movement from frame-to-frame.

Depth stuff summary

hello

3.2 Frame-to-Frame Tracking

For each frame Mask RCNN outputs a set of bounding boxes, and masks. These are not within any particular order and so a method must be devised to track objects from one frame to the next. This is a hard challenge due to object occlusion, variation in shadow/lighting conditions and rapid in-frame movement. An obvious method is intersection over union (IOU) of the bounding boxes on a frame-by-frame basis. This is good in theory, however, it lacks the ability to compensate for rapid movement (of the object and the camera angle), occlusion and object overlapping. A set of methods have therefore been developed to overcome this problem, which combined provide a robust tracking system.

The method must both track objects in the frame, determine which objects are no-longer in-frame and find any new objects that enter the frame. Information about each object must

be stored as to determine the velocity and direction of travel/ if a collision is likely between the cyclist and the object in question. Figure 3.3 displays an example of a working system.

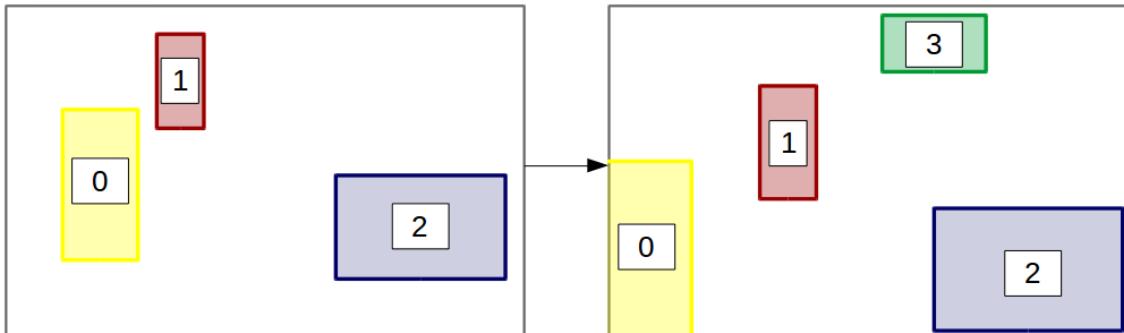


Figure 3.3: ROI Old to New Frame Example.

3.2.1 Kernelised Correlation Filter

A Kernelised correlation filter - one of the state-of-the-art methods used today for tracking objects from one frame to another is described in Section 2.3.1. A minor modification is required for use in a multi-object scenario which is described below:

Modifications: Since the bbox prediction is so fast for the KCF tracker; it can be used in conjunction with the other tracking methods mentioned in Sections 3.2.2, 3.2.3, and 3.2.4. The kcf tracking algorithm outputs a bounding box which is the predicted coordinates of the object in the next frame. This can be compared to Mask-RCNNs output through a simple matching process. A box in the new frame predicted by Mask RCNN $[x_{mask}, y_{mask}, width_{mask}, height_{mask}]$ can be scored by sampling the pdf from a multivariate Gaussian, in 4D, with mean $[x_{kcf}, y_{kcf}, width_{kcf}, height_{kcf}]$ and covariance matrix. The covariance matrix is an arbitrary choice but determines how far the system looks for matches from the predicted kcf location. The covariance matrix in the system is set to $20 \times I_{4 \times 4}$ which translates to a $\sigma^2 = 20$ pixels for $x, y, width, height$ respectively, and assumes independence between predicted variables. This score can then be used in Section 3.2.6 to be combined with the other tracking modules.

3.2.2 BRISK - Selector

BRISK as described earlier in Section-2.2.6 is a binary descriptor that can match keypoints across two photographs. This can be used also in the tracking section to match keypoints across ROIs, although not thoroughly accurate, it offers a method of matching ROIs regardless of location. This should help re-detect objects after periods of occlusion or overlapping with others that IOU would be unable to do. The process is simple to match the ROIs. First keypoints and descriptors are calculated for each ROI, these are then compared using the fast binary matching process of BRISK, then the Hamiltonian distance is calculated for each matching keypoint... TODO : EXPAND THIS SECTION

3.2.3 Class Matching

An obvious additional requirement is the matching of classes between frames, as a person's ROI should not match to the ROI of a car in the next frame. Hence for each ROI, the class is stored for access by the next frame. This can both reduce expensive computation time of other metrics but provides an easy ROI comparison tool. It can be noted though that Mask RCNN is not perfect and often mis-classifies objects. It is, however, reliable enough to be assumed as correct.

3.2.4 Kalman Filter ROI Matching

A Kalman filter, developed by Rudolf E. Kalman has two main parts, a predictive step and an update. The predictive step tries to find a-priori estimates of the state for the next time-step by projecting forward the current state and error covariance estimates. This estimate is then updated using the measurements (camera) of the object location. The iterative process continues improving the Kalman filter estimates at each stage.

This was first tested with just a camera frame state i.e. using the [x, y, width, height] positions and distances in pixels of each ROI on the input image. It was quickly discovered that the model was not sufficiently robust as it did not take into account the bike movement. Hence the filter was unable to track objects efficiently. Hence, once the Okvis framework was added, the filter was converted into tracking world frame coordinates which is far more successful.

Kalman Predictive Step

Assuming a linear system with Gaussian process, the state can be modelled as:

$$\hat{X}_k = F_k \hat{X}_{k-1} + w_{k-1} \quad (3.8)$$

where k is the time step, $\hat{X}_k \in \mathbb{R}^n$ is the state vector, F_k is the $n \times n$ state transition matrix from state k to $k+1$ and w_k is $\mathcal{N}(0, Q)$. The state vector is described as follows:

$$\hat{X}_k = [Cx, Cy, Cz, w, h, Cx_{vel}, Cy_{vel}, Cz_{vel}, w_{vel}, h_{vel}] \quad (3.9)$$

where (Cx, Cy) is the centre point of the bounding box of the object, Cz is the median distance of the masked ROI, w the width and h the height. Cx_{vel} describes the centre point velocity in the x direction, and $Cy_{vel}, Cz_{vel}, w_{vel}, h_{vel}$ the y/z direction, w the width and h the height. $[Cx, Cy, w, h]$ describe the measurable variables of the system, and there are no controllable parameters in this system. From basic equations of motion, a general measurable parameter γ , and corresponding velocity γ_{vel} follows the update rule:

$$\begin{aligned} \gamma_{k+1} &= \gamma_k + \gamma_{vel} \times \Delta t \\ \gamma_{vel,k+1} &= \gamma_{vel,k} \end{aligned} \quad (3.10)$$

where Δt is the time step between each frame. Therefore, the update matrix for the filter $F_k \in R^{10 \times 10}$ is as follows:

$$F_k = \begin{bmatrix} I_5 & \Delta t \times I_5 \\ 0_{5 \times 5} & I_5 \end{bmatrix} \quad (3.11)$$

The final term $w_k = \mathcal{N}(0, Q)$ describes the process noise in the system and is normally distributed with error covariance matrix Q . This is updated by the kalman filter at each step. Since the state describes 3D coordinates it is useful to determine how far the average walking human would walk (in m) from one frame to the next, given an average walk of 1.5m/s:

$$dist = \frac{1.5}{30} \approx 0.05 \quad (3.12)$$

Speeding this up a bit for runners, the x, y, z variance would not be above 0.1 so that is a good initialising parameter for the filter x, y, z variance states. The width and height is in pixel space and so would change a more rapidly, and so is scaled accordingly.

$$\begin{aligned} v &= [0.1, 0.1, 0.1, 7, 7, 0.01, 0.01, 0.01, 0.7, 0.7] \\ Q_k &= diag(v) \end{aligned} \quad (3.13)$$

where $Q_k \in \mathbb{R}^{10 \times 10}$ and is the diagonalised matrix of v .

Kalman Update Step

The measurement step can also predict the current state. With $Z_k \in \mathbb{R}^{10}$ as the system measurement vector:

$$Z_k = H_k X_k + v_k \quad (3.14)$$

where H_k is an 10×5 measurement matrix relating X_k to the measurement Z_k , and v_k is $\mathcal{N}(0, R)$. Since only the Cx, Cy, Cz, w, h is measurable H_k takes the form:

$$H_k = [I_5 \quad 0_{5 \times 5}] \quad (3.15)$$

and the measurement noise $v_k = \mathcal{N}(0, R_k)$ is set with covariance matrix. As discussed earlier the depth measurement is most prone to error, and the width/height are in pixel space. The Kalman measurement noise covariance therefore starts as:

$$\begin{aligned} r &= [0.2, 0.2, 0.2, 20, 20] \\ R_k &= diag(v) \end{aligned} \quad (3.16)$$

where $R_k \in \mathbb{R}^{5 \times 5}$ and is the diagonalised matrix of r .

Kalman Equations

The prior estimates for the next time step by the Kalman predictions step (no control vector):

$$\begin{cases} \hat{X}_k(-) = F_k \hat{X}_{k-1}(+) \\ P_k(-) = F_k P_{k-1}(+) F_k^T + Q_k \end{cases} \quad (3.17)$$

where $\hat{X}_k(-)$ is a priori state estimate and $P_k(-)$ is the priori estimate error at step k . It is essentially the previous state \times the transition matrix to the new state ($\hat{X}_k(-)$). Then a similar approach for the error covariance between the previous and current state. The correction step is then:

$$\begin{cases} \tilde{y}_k = Z_k - H_k \hat{X}_k(-) \\ S_k = R_k + H_k P_k(-) H_k^T \\ K_k = P_k(-) H_k^T S_k^{-1} \\ \hat{X}_k(+) = \hat{X}_k(-) + K_k \tilde{y}_k \\ P_k(+) = (I - K_k H_k) P_k(-) (I - K_k H_k)^T + K_k R_k K_k^T \\ \tilde{y}_{k|k} = Z_k - H_k \hat{X}_k(+) \end{cases} \quad (3.18)$$

where K_k is the $n \times m$ Kalman gain matrix, $\hat{X}_k(+)$ is the posterior state estimate according to the actual measurement Z_k and the predicted measurement $H_k \hat{X}_k(-)$ and $P_k(+)$ is the posterior state estimate error covariance.

Kalman Matching ROIS

From one frame to another it can be assumed that a number of bounding boxes (ROIS) are consistent in both frames. Hence a metric to match a box from a prior frame to the current frame is required. Using the Kalman filter, a prediction of where the ROI is at time k is easily obtained. But this does not give a likelihood measure of the boxes in the new frame. However, in the kalman update step \tilde{y}_k , the innovation residual, and S_k , the innovation covariance matrix are calculated. These can be used to determine the marginal likelihood of a newer ROI's $[Cx, Cy, w, h]$ given the kalman $P_k(-)$ and $\hat{X}_k(-)$. Skipping the formulation, the log-marginal likelihood is:

$$\begin{aligned} l &= \log p(\mathbf{z}) \\ l^{(-1)} &= 0 \\ l^{(k)} &= l^{(k-1)} - \frac{1}{2} (\tilde{y}_k^T S_k^{-1} \tilde{y}_k + \log |S_k| + d_y \log 2\pi) \end{aligned} \quad (3.19)$$

where d_y is the dimension of the measurement matrix. The part dependent on the new ROI is $\tilde{y}_k^T S_k^{-1} \tilde{y}_k$ as S_k , from Equation 3.18 does not involve \tilde{y}_k . Therefore the matching step can be performed by simply matching each ROI_{old} in the old frame to the ROI_{new} in the new frame that has the lowest scalar value from:

$$M = \tilde{y}_k^T S_k^{-1} \tilde{y}_k \quad (3.20)$$

with a threshold value to ensure no incorrect matches. The matching procedure must be the same for all methods and matches the maximum value. Therefore the *Match Metric* is as follows:

$$Match\ Metric = \begin{cases} 0 & \text{if } K_{threshold} - M \leq 0 \\ K_{threshold} - M & \text{otherwise} \end{cases} \quad (3.21)$$

3.2.5 ROI Lives

Each region of interest may be tracked from one frame to the next. However, if that object is occluded then it may not be spotted in the next frame. This, unless accounted for could lead to information loss regarding that object's state. Therefore a life system is in place whereby an object has a number of lives that are used when the object is not spotted in the next frame. These are decremented until 0 when the object is forgotten and no longer tracked.

3.2.6 Combining Matching Methods

Taking the example provided in Figure 3.3, there are three objects in the older frame and four in the new frame. BRISK (Section 3.2.2) is not reliable to match all objects from one frame to the next and sometimes cannot find matching keypoints between the same object in each frame (especially deformable objects such as people). Therefore, a system to combine methods is needed. For each ROI_{old} and each ROI_{new} the match metrics for each method (KCF, BRISK, Kalman) is calculated. This forms a matrix $matches \in \mathbb{R}^{new \times old \times methods}$ where old, new are the number of ROIs in the old and new frames, and $methods$ is the number of matching methods implemented (in this case two) but is designed to be easily extendible to more. The rules for the matching procedure are as follows:

After the matching process is complete, there are a number of extra steps that must take place before continuing to the next frame. For example, any oldROIs that are not detected in the new frame must have a life decremented and be added to the list of new frame ROIs with the relevant location and information.

Then the matched ROIs pairs from the old and new frame must exchange relevant information in preparation for the next frame:

Finally any new ROIs that were not in the old frame need to be assigned an individual ID to be recognised over the next frames.

Algorithm 4: Algorithm to match old and new frames**match**
 $(old_frame), new_frame, [matching_methods(oldROI, newROI)], matches, threshold$
Input :
 old_frame // ROIs as structures from the old frame

 new_frame // ROIs as structures from the new frame

 $[matching_methods(oldROI, newROI)]$ // List of pointers to functions for matching ROIs
// (KCF, Kalman Filter, Brisk)

 $matches$ // A matrix of zeros $\in \mathbb{R}^{(methods, oldROIs, newROIs)}$
 $threshold$ // The minimum match score considered to denote an inter-frame ROI match
Output:
 $old_matched$ // A list of matched oldROIs, $old_matched[i]$ matches to $new_matched[i]$
 $new_matched$ // A list of matched newROIs, $new_matched[i]$ matches to $old_matched[i]$
 $// Calculate score of match ($oldROI \rightarrow newROI$) for each method, $oldROI, newROI$$
foreach $oldROI \in old_frame$ **do** /* For every ROI in last frame */

foreach $newROI \in new_frame$ **do** /* For every ROI in current frame */

foreach $method \in matching_methods(oldROI, newROI)$ **do** /* For each method */

 $| \quad matches[method_i, oldROI_i, newROI_i] = matching_methods(oldROI, newROI)$
end
end
end
 $// Multiply along methods axis to produce one combined match matrix.$
foreach $method \in matching_methods(oldROI, newROI)$ **do** /* For every method */

if $method_i > 0$ **then** /* If method index is > 0 */

 $| \quad matches[0, :, :] = mult(matches[0, :, :], matches[method_i, :, :])$ /* Multiply along axis

 $*/$
end
 $matches = matches[0, :, :]$
 $// Match new and old ROIS.$
 $max = matches.max()$
while $max > threshold$ and $len(new_matched) < len(new_frame)$ **do**
 $| \quad old, new = matches.max()_{i,j}$ /* Find index of max element */

if $old \notin old_matched$ and $new \notin new_matched$ **then** /* If not already matched */

 $| \quad new_matched.append(new)$ /* Add new index to already matched list */

 $| \quad old_matched.append(old)$ /* Add old index to already matched list */

 $| \quad matches[new, old] = -1$
 $| \quad max = matches.max()$
end

Algorithm 5: Algorithm to append oldROIs that aren't detected in new frame

add_old_rois (*old_matched*, *old_frame*, *new_frame*)

Input :

<i>old_matched</i>	// A list of matched oldROIs, <i>old_matched</i> [<i>i</i>] matches to <i>new_matched</i> [<i>i</i>]
<i>old_frame</i>	// ROIs as structures from the old frame
<i>new_frame</i>	// ROIs as structures from the new frame

Output:

<i>new_frame</i>	// An updated list of newROI objects
------------------	--------------------------------------

// Update newROIs to include undetected ROIs from old frame.

```

foreach old  $\notin$  old_matched do          /* For every old roi not found in new frame */
    older = old_frame[old]           /* Get olderROI object */
    if older.lives > 1 then          /* If the ROI has a life to spare */
        older.lives -= 1             /* Decrement a life */
        update_location(older, TCW) /* Update camera/world location from Kalman prediction
        /*
        update_model(older, TCW)      /* Update pedestrian model from Kalman prediction */
        add_information(new_frame, older)  /* Append older ROI information to new frame
        /*
        /* (masks,lives,id,class,colour,briskKeypoints,briskDescriptors,KCFtracker) */

end

```

Algorithm 6: Algorithm to update newROIs with matched oldROIs

```

update_matched_rois (old_matched, new_matched, old_frame, new_frame)
Input :
old_matched                                // A list of matched oldROIs, old_matched[i] matches to new_matched[i]
new_matched                                // A list of matched newROIs, new_matched[i] matches to old_matched[i]
old_frame                                    // ROIs as structures from the old frame
new_frame                                    // ROIs as structures from the new frame
Output:
new_frame                                    // An updated list of newROI objects

// Add older ROI details to matched new ROI
foreach old,new ∈ old_matched,new_matched do /* For every matched old/new ROI pair */
    older = oldROIs[old]                      /* Get olderROI object */
    newer = newROIs[new]                      /* Get newerROI object */
    newState = newer.worldPoint                /* Get new location of ROI */
    older.kalman.correct(newState)             /* Update Kalman Filter with new location */
    older.model.add_points(newState)           /* Add new points to movement model */
    if older.lives < lifeThreshold then        /* If lives < threshold */
        | older.lives += 1
    update_object(newer,older)                  /* Update new ROI with older ROI information */
                                                /* (Kalman Filter,lives,colour,id,model) */
end

```

Algorithm 7: Algorithm to create new IDs for object in new frame

```

update_matched_rois (old_matched, new_matched, old_frame, new_frame)
Input :
old_matched                                // A list of matched oldROIs, old_matched[i] matches to new_matched[i]
new_matched                                // A list of matched newROIs, new_matched[i] matches to old_matched[i]
new_frame                                    // ROIs as structures from the new frame
Output:
new_frame                                    // An updated list of newROI objects

// Note that if there are no objects in old frame new_frame.id = zeros
maxId = new_frame.id.max() + 1           /* Get current maximum ID + 1 */
foreach new  $\notin$  new_matched do          /* For every not-matched new ROI */
    newer = new_frame[new]            /* Get newerROI object */
    newer.id = maxId                   /* Assign new ID */
    maxId += 1                          /* Increment the maximum ID */
end

```

3.3 Movement Prediction

Each object and the bike must be modelled from frame-to-frame and the subsequent path predicted for collision detection. The further in the future that can be accurately modelled, the better the warning system. The following section describes both the more complex bicycle model and the simple path prediction for pedestrians.

3.3.1 Objects Prediction

The objects are modelled by the Kalman filter from frame to frame and this provides a good estimate of both the direction and predicted position. Hence in a future position, say in 10 frames would simply be obtained from multiplying by a modified state-transition model.

$$F_{k+n_frames} = \begin{bmatrix} I_5 & \Delta t \times n_frames \times I_5 \\ 0_{5 \times 5} & I_5 \end{bmatrix} \quad (3.22)$$

where δt is the time difference from one frame to the next and n_frames is the number of frames in the future to predict the position at.

$$\hat{X}_{k+n_frames}(-) = F_{k+n_frames} \hat{X}_{k-1}(+) \quad (3.23)$$

The uncertainty of that position also increases with time, and can be taken directly from the error covariance prediction matrix. This however is updated iteratively and so it is easy enough just to repeat Equation 3.24 n_frames times.

$$P_{k+n_frames}(-) = F_{k+n_frames} P_{k-1}(+) F_{k+n_frames}^T + Q_k \quad (3.24)$$

This method then produces a predicted position and the uncertainty in that position which can be used to create spheres of possible object positions in future frames. These can then be simply checked for intersection with the bike model predictions.

3.3.2 Bike Model

The bike is deemed to be the most important model and can therefore predictions should take into account higher order features, for example the current state of the bicycle in the frame and perhaps external factors such as the surrounding environment. The problem is therefore split into two branches, one for sequence prediction of Okvis pose outputs and the other of predictions using image inputs.

Okvis Sequence Inputs: Since the model is required to learn a sequence a recurrent neural network is the obvious choice. This mimicks memory from one-frame to the next, providing n previous frames to predict the $(n + 1)^{th}$ frame. The proposed architecture for this model, by no means fine-tuned but merely to provide a proof-of-concept is as follows:

Pose RNN Model: The main problem is to generalise the model to any world point. For example, the bike could be at any location or offset from the world origin and thus a huge amount of data would be needed for learning if the world location was an input. Instead, each sequence contains an origin point at sequence number zero. The model simply learns the difference. Note however, things such as orientation can make a large difference, if the bike is at a big angle in the world frame it will be turning. This information cannot be transferred from the difference between an original frame and the next. Hence the rotation part of the transformation matrix remains in the world frame while the translation has an origin equal to that of the 1st frame in the sequence Equation 3.25.

$$\begin{aligned} X_{frame_n} &= [C_{WC_n} \quad r_{WC_n} - r_{WC_0}] \\ y &= r_{WC_{n+1}} \end{aligned} \tag{3.25}$$

The model is formed of the layers detailed on the right of Figure 3.4. Given a limited input space $X = (len_{sequence}, 3 \times 4)$ where the transformation matrix is flattened for use in the LSTM layers to a matrix of size 12, discarding the last row $[0_{1 \times 3}, 1]$. The model is trained in a regression setting with labels $y = [x, y, z]$ which are coordinated representing the displacement from first frame of the sequence. As it would be too granular to learn from every frame, a larger time-gap between frames can be used and hence the model sequence can consist of every k^{th} frame from the camera. The loss function used is mean squared error and is formulated as follows:

$$mse = \sum_{i=0}^n \sum_{j=0}^m (\mathbf{y}_{i,j} - \hat{\mathbf{y}}_{i,j})^2 \tag{3.26}$$

where i is the sequence number, j is the variable prediction and therefore $m = 3$ for $[x, y, z]$.

Images RCNN Model: For the image section of the model, a standard convolutional learning approach is adopted as described in Section 2.1.2. The loss function for the RCNN model is the same as Equation 3.26, as once again this is a regression problem. The images are fed into the network and the difference predicted much like before. The only variation is the

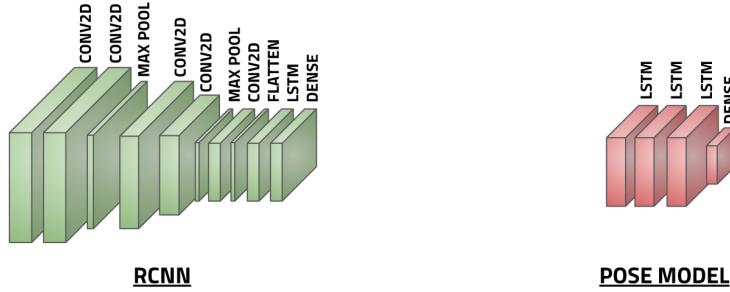


Figure 3.4: RCNN Architecture: Left, Pose Architecture: Right.

data input and the flattening layer before the LSTM. A convolutional layer outputs a number of 3D tensors (w, h, d) which is not compatible with the LSTM or dense matrices. Hence a flattening layer converts the tensors into a shape $(sequence_length, w \times h \times d)$ for input into the final network layers.

Optimising Learning Rate and Decay: Adam optimiser, a common solver for deep learning problems requires a number of parameters which are all set as the default recommended from the paper [?] $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$ these are adaptive moment parameters that control the algorithm's descent into minima. However learning rate and learning rate decay are problem specific and can vastly impair the ability of a model to learn correctly. Hence a logical approach to selection of the right values is required. A grid search method was chosen and loss characteristics at epoch 200 compared to get a reasonable approximate for the correct learning rate.

Bayesian Optimisation Layer Dimensions:

Total Model: Both models discussed in this section can be trained separately to obtain predictions of the future pose. A naive estimation would take weighted predictions and sum for a final prediction result. The final layers of these networks however could infer details from each-other to inform the pose estimation. This can only be done in a branched network that uses both inputs effectively and concatenates the outputs for some final layers. The proposed model is graphed in Figure 3.5, it has two inputs tensors of the images and pose data and three output variables which is the predicted $[x, y, z]$.

Training Regime: The network is large and so a careful training regime is required to ensure convergence to a reasonable minimum. The two initial models, the pose and RCNN images are first pre-trained to predict poses. These are the complete models described in Figure 3.4. Once this training process is complete the weights are used in the larger branched model. These layers are then frozen and the final dense layers are trained. Finally all layers are unfrozen and the whole model is trained with a lower learning rate to fine-tune the weights.

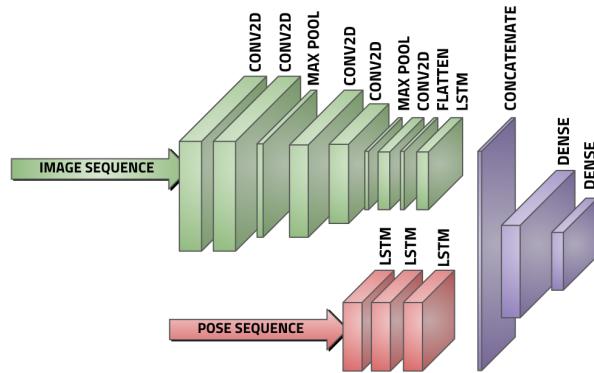


Figure 3.5: Complete Network Architecture.

3.4 System Evaluation

3.4.1 Mask-RCNN

In applying Mask-RCNN to the specific urban environment there have been a number of interesting outcomes. On the whole, the system is robust, often detecting images well and so there are rarely gaps when an object is detected in one frame and not the next. A number of incorrect classifications often provide a worry however, for example, the classification of road paint can cause a problem to the system (Figure ...). This often happens with road markings of bicycles that are relatively convincing examples. After experiencing this phenomenon, a number of examples were sought out to determine how the segmentation performs.

3.4.2 Okvis with Dynamic Objects

3.4.3 Depth Camera

The intel realsense depth camera is typically used in indoor environments as it operates with a small laser. This, as is expected works less well in highly lit environments. The extent of the problem, however is often drastic and can cause okvis to perform poor estimations of pose and hence is unable to determine the bicycle location. An example of these images can be seen in Figure... . It could quite easily be solved, however, with the use of a stereo camera for depth perception. For example the realsense R... offers such capabilities which would mitigate the depth camera sensitivity to bright lighting conditions.

Chapter 4

Ethics

4.0.1 Modification for Military Purposes

The code presented in this report, although strictly for civilian applications does have the potential to be modified for military use. It obviously provides the ability to track people and objects from frame-to-frame. This however, is not in a real-time context at present and would require significant modification for use in such an environment that would be conducive to a military application. However, that being said little modification is needed for a simple tracking application from CCTV data for example. This would not need to run at real time and could be a potential source of misuse.

4.0.2 Other Ethics Issues

The code developed is intended to offer a basis for an early warning system to cyclists in urban environments. Hence a number of the inherent ethical problems with, for example the autonomous vehicle space are also applicable to this project. The potential over-reliance on such a system of cyclists could be dangerous and allow for lack of attention and potential accidents when the system is unable to detect a collision. That having been said, it is evident that there is a need for such a system for cyclists and so the benefits would clearly outweighs the disadvantages in such a case.

4.1 Ethics Checklist

	Yes	No
Section 1: HUMAN EMBRYOS/FOETUSES		
Does your project involve Human Embryonic Stem Cells?		✓
Does your project involve the use of human embryos?		✓
Does your project involve the use of human foetal tissues / cells?		✓
Section 2: HUMANS		
Does your project involve human participants?		✓
Section 3: HUMAN CELLS / TISSUES		
Does your project involve human cells or tissues? (Other than from Human Embryos/Foetuses i.e. Section 1)?		✓

Section 4: PROTECTION OF PERSONAL DATA		
Does your project involve personal data collection and/or processing?		✓
Does it involve the collection and/or processing of sensitive personal data (e.g. health, sexual lifestyle, ethnicity, political opinion, religious or philosophical conviction)?		✓
Does it involve processing of genetic information?		✓
Does it involve tracking or observation of participants? It should be noted that this issue is not limited to surveillance or localization data. It also applies to Wan data such as IP address, MACs, cookies etc.		✓
Does your project involve further processing of previously collected personal data (secondary use)? For example Does your project involve merging existing data sets?		✓
Section 5: ANIMALS		
Does your project involve animals?		✓
Section 6: DEVELOPING COUNTRIES		
Does your project involve developing countries?		✓
If your project involves low and/or lower-middle income countries, are any benefit-sharing actions planned?		✓
Could the situation in the country put the individuals taking part in the project at risk?		✓
Section 7: ENVIRONMENTAL PROTECTION AND SAFETY		
Does your project involve the use of elements that may cause harm to the environment, animals or plants?		✓
Does your project deal with endangered fauna and/or flora /protected areas?		✓
Does your project involve the use of elements that may cause harm to humans, including project staff?		✓
Does your project involve other harmful materials or equipment, e.g. high-powered laser systems?		✓
Section 8: DUAL USE		
Does your project have the potential for military applications?		✓
Does your project have an exclusive civilian application focus?	✓	
Will your project use or produce goods or information that will require export licenses in accordance with legislation on dual use items?		✓
Does your project affect current standards in military ethics e.g., global ban on weapons of mass destruction, issues of proportionality, discrimination of combatants and accountability in drone and autonomous robotics developments, incendiary or laser weapons?		✓
Section 9: MISUSE		
Does your project have the potential for malevolent/criminal/terrorist abuse?		✓

Does your project involve information on/or the use of biological-, chemical-, nuclear/radiological-security sensitive materials and explosives, and means of their delivery?	✓	
Does your project involve the development of technologies or the creation of information that could have severe negative impacts on human rights standards (e.g. privacy, stigmatization, discrimination), if misapplied?		✓
Does your project have the potential for terrorist or criminal abuse e.g. infrastructural vulnerability studies, cybersecurity related project?		✓
Section 10: LEGAL ISSUES		
Will your project use or produce software for which there are copyright licensing implications?	✓	✓
Will your project use or produce goods or information for which there are data protection, or other legal implications?		✓
Section 11: OTHER ISSUES		
Are there any other ethics issues that should be taken into consideration?		✓

Chapter 5

Experimental Results

Chapter 6

Conclusion

Bibliography

- [1] D. F. Transport, “Table ras30001: Reported road casualties by road user type and severity, great britain, 2016,” report, Dft, 2016. pages 1
- [2] D. F. Transport, “Ras50001: Contributory factors in reported accidents by severity, great britain, 2016,” report, Dft, 2016. pages 1
- [3] D. F. Transport, “Focus on cycling in reported road casualties great britain 2013,” report, Dft, 2013. pages 1
- [4] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. pages 2
- [5] O. R. et al.and Jia Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge,” *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015. pages 2
- [6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” vol. 25, 01 2012. pages 2
- [7] A. Karpathy, “Networks,” 2017. pages 3
- [8] A. Karpathy, “Networks,” 2017. pages 3
- [9] A. Karpathy, “Networks,” 2017. pages 3
- [10] S. Sabour, N. Frosst, and G. E. Hinton, “Dynamic routing between capsules,” *CoRR*, vol. abs/1710.09829, 2017. pages 3
- [11] R. B. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” *CoRR*, vol. abs/1311.2524, 2013. pages 4, 5
- [12] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” *CoRR*, vol. abs/1506.02640, 2015. pages 4
- [13] J. R. R. ”Uijlings, K. E. A. van de Sande, T. Gevers, and A. W. M. Smeulders, “selective search for object recognition”, ”*International Journal of Computer Vision*”, ”2013”. pages 4, 6
- [14] P. F. ”Felzenszwalb and D. P. Huttenlocher, “”efficient graph-based image segmentation”, ”*International Journal of Computer Vision*”, ”2004”. pages 4

- [15] R. B. Girshick, “Fast R-CNN,” *CoRR*, vol. abs/1504.08083, 2015. pages 5, 6
- [16] K. He, X. Zhang, S. Ren, and J. Sun, “Spatial pyramid pooling in deep convolutional networks for visual recognition,” *CoRR*, vol. abs/1406.4729, 2014. pages 5
- [17] S. Ren, K. He, R. B. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *CoRR*, vol. abs/1506.01497, 2015. pages 6, 9
- [18] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” *CoRR*, vol. abs/1411.4038, 2014. pages 7, 8
- [19] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, “The cityscapes dataset for semantic urban scene understanding,” *CoRR*, vol. abs/1604.01685, 2016. pages 8
- [20] L. Wan, M. Zeiler, S. Zhang, Y. LeCun, and R. Fergus, “Regularization of neural networks using dropconnect,” in *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ICML’13, pp. III–1058–III–1066, JMLR.org, 2013. pages 8
- [21] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014. pages 8
- [22] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006. pages 8
- [23] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, “Mask R-CNN,” *CoRR*, vol. abs/1703.06870, 2017. pages 9
- [24] M. Deisenroth, “Gaussian processes and bayesian optimisation,” 2018. pages 10
- [25] N. Srinivas, A. Krause, S. Kakade, and M. Seeger, “Gaussian process optimization in the bandit setting: No regret and experimental design,” 2010. pages 11
- [26] O. docs, “Images,” 2012. pages 12
- [27] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, and P. Furgale, “Keyframe-based visual-inertial odometry using nonlinear optimization,” *The International Journal of Robotics Research*, vol. 34, no. 3, pp. 314–334, 2015. pages 14, 17, 18, 21
- [28] H. F. Durrant-Whyte, “Uncertain geometry in robotics,” vol. 4, pp. 23 – 31, 03 1988. pages 17
- [29] T. Bailey and H. Durrant-Whyte, “Simultaneous localization and mapping (slam): Part i,” 10 2006. pages 18
- [30] S. Leutenegger, M. Chli, and R. Y. Siegwart, “Brisk: Binary robust invariant scalable keypoints,” in *2011 International Conference on Computer Vision*, pp. 2548–2555, Nov 2011. pages 18, 20
- [31] O. docs, “Images,” 2012. pages 19

- [32] P. Furgale, “Extensions to the visual odometry pipeline for the exploration of planetary surfaces,” 01 2011. pages 20
- [33] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, “High-speed tracking with kernelized correlation filters,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, pp. 583–596, March 2015. pages 21, 22