

Approval of the Department of Electronic & Telecommunication Engineering

.....
Head, Department of Electronic &
Telecommunication Engineering

This is to certify that I/we have read this project and that in my/our opinion it is fully adequate, in scope and quality, as an Undergraduate Graduation Project.

Supervisors:

Dr. Ranga Rodrigo

Signature:

Dr. Ajith Pasqual

Signature:

Dr. Peshala G. Jayasekara

Signature:

Date:

Declaration

This declaration is made on October 5, 2017.

Declaration by Project Group

We declare that the dissertation entitled People Counting and Tracking with Xilinx ZC702 Evaluation Kits and the work presented in it are our own. We confirm that:

- this work was done wholly or mainly in candidature for a B.Sc. Engineering degree at this university,
- where any part of this dissertation has previously been submitted for a degree or any other qualification at this university or any other institute, has been clearly stated,
- where we have consulted the published work of others, is always clearly attributed,
- where we have quoted from the work of others, the source is always given. With the exception of such quotations, this dissertation is entirely our own work,
- we have acknowledged all main sources of help.

.....

Date

.....

R.V.C.N Abeyrathne (130008K)

.....

D.L Dampahalage (130093M)

.....

H.A.S.P Gunasekara (130183N)

.....

W.M.D.K Weerakoon (130633V)

Declaration by Supervisors

We have supervised and accepted this dissertation for the submission of the degree.

.....
Dr. Ranga Rodrigo

.....
Date

.....
Dr. Ajith Pasqual

.....
Date

.....
Dr. Peshala G. Jayasekara

.....
Date

Abstract

People Counting and Tracking with Xilinx ZC702 Evaluation Kits

Group Members: R.V.C.N Abeyrathne, D.L Dampahalage, H.A.S.P Gunasekara, W.M.D.K Weerakoon

Supervisors: Dr. Ranga Rodrigo, Dr. Ajith Pasqual, Dr. Peshala G. Jayasekara

Keywords: FPGA,ZYNQ-7000, people counting, multi-camera, multiple people tracking, particle filters, gaussian mixture models

In the contemporary society making correct decisions is vital for any business organization to stay on par with competitors. For that intent identifying and understanding the customers is a must. Tapping into the customers' subconscious is the preeminent way of making correct business decisions and for that an organization should track and analyze customer behavior.

As a solution to this problem, we propose a scalable multi-camera people tracking system for a large department store like environment. The main focus of this project is to implement an end to end system for multi-camera people tracking with business intelligence. Novelty of this system from contemporary multi-camera people tracking systems is the introduction of leaf node processing, where part of the processing is done in a hardware system close to the camera. By this way we could achieve the scalability required for a large department store like environment and reduce bandwidth usage in such systems.

This proposed system comprises of leaf nodes, central server and the web server. Leaf nodes are integrated close each camera. We use Zynq SoCs for leaf nodes and these devices read video frames from the cameras, detects people for each frame using GMM based background subtraction, calculate color histogram features for each detected person and sends the bounding boxes and features to the server using UDP.

A central server will receive this information from each leaf node and do people tracking by using particle filters and hungarian algorithm for data association. Multi-camera tracking is also achieved by means of homography based method and a graph algorithm. Real time tracking and heatmaps are generated in the server GUI.

Finally business intelligence software hosted in a web server, where real time people tracking for the whole store, heatmaps for desired time period, people count graphs for desired time period and real time people count graphs are implemented.

Therefore by the introduction of leaf node processing we achieve a scalable system since the processing power required at the central server is reduced. We achieve a low power system for leaf nodes by using Zynq SoCs and we reduce the bandwidth usage since we only send the bounding boxes and features for tracking to the server.

Acknowledgments

We would like to express our deepest appreciation and sincere gratitude to our project supervisors, Dr. Ajith Pasqual, Dr. Ranga Rodrigo and Dr. Peshala Jayasekara of Department of Electronic and Telecommunication Engineering of University of Moratuwa, for their guidance and constant supervision as well as for providing necessary information regarding the project.

We would like to thank ParaQum Technologies for their guidance and support given to us during the project. We are indebted to our final year project coordinator, Dr. Anjula Silva for his counsel regarding the project. We would also like to take this opportunity to thank Ms. Salgado for her advices on compiling a finer report. We would like to acknowledge with much appreciation the non-academic staff of the department who gave us their fullest support in providing necessary laboratory facilities.

Finally, our thanks and appreciations also goes out to our batch mates who helped in numerous ways to make this endeavor a success.

Contents

Approval	i
Declaration	ii
Abstract	iv
Acknowledgments	v
List of Figures	viii
List of Tables	ix
Acronyms and Abbreviations	x
1 Introduction	1
2 Literature Review	3
3 Methodology	5
3.1 Introduction	5
3.2 Scope	5
3.3 Overall Architecture	6
3.4 Embedded System (Leaf Node)	6
3.4.1 Xilinx ZC702 Board	7
3.4.2 ZYNQ-7000 SoC	8
3.5 Design Flow for ZYNQ-7000 based Embedded System	8
3.5.1 Vivado HLS	9
3.5.2 Vivado	9
3.5.3 Xilinx SDK	9
3.5.4 PetaLinux SDK	9
3.6 Installing Linux on ZYNQ-7000	9
3.7 People Detection and Feature Calculation Implementation on ZYNQ-7000 SoC	10
3.7.1 IP core implementation on ZYNQ PL	11
3.7.2 IP Core drivers, Morphological Operations and Communication Protocol Implementation on ZYNQ PS	14
3.7.3 False Positive Reduction	16
3.7.4 Communication Protocol among leaf nodes central server	20
3.8 Central Server	21
3.8.1 Receiving and Unloading Data Packets from Camera Nodes	21
3.8.2 Single Camera Tracking	22
3.8.3 Particle Filter Based Tracking	24

3.8.4	Business Intelligence Software	25
4	Results	27
4.0.1	Discussion	27
	References	28

List of Figures

3.1	Overall block diagram of the end-to-end system	6
3.2	Xilinx ZC702 Development Board	7
3.3	Internal architecture of ZYNQ-7000 SoC	8
3.4	Design flow block diagram for ZYNQ-7000 based embedded system . .	8
3.5	Architecture of the Leaf node	11
3.6	Structure of the 3D RGB histogram	14
3.7	Block diagram of the morphological operation flow to detect people from background subtracted binary mask	15
3.8	Results of morphological operations for a test image	17
3.9	Scatter plot matrix for the 6 features collected from running our object detection on a sample video	18
3.10	First Vs. second pricipal component	19
3.11	Density of the first principal component	19
3.12	Results of PCA based false positive elimination algorithm	20
3.13	Multiple People Tracking using Kalman Filter	22
3.14	Particle Filter based Tracking	24
3.15	Flask Logo	25

List of Tables

3.1	Percentage of variance along each principal component	18
-----	---	----

Acronyms and Abbreviations

FPGA	Field Programmable Gate Array
GMM	Gaussian Mixture Models
PCA	Principal Component Analysis
YOLO	You Only Look Once
IP	Intellectual Property
SoC	System on Chip
RTL	Register Transfer Level
AXI	Advanced eXtensible Interface
HLS	High Level Synthesis
UIO	Userspace Input/Output
REST	Representational State Transfer
API	Application Programming Interface
AJAX	Asynchronous JavaScript and XML
JSON	JavaScript Object Notation
UDP	User Datagram Protocol
TCP	Transmission Control Protocol
URL	Uniform Resource Locator
SDRAM	Synchronous Dynamic Random Access Memory
DDR	Double Data Rate

Chapter 1

INTRODUCTION

Nowadays retail stores are very popular among people. Even in Sri Lanka, retail stores like Cargills Food City and House of Fashion are visited by at least few hundreds of customers daily. In these type of retail stores, one of the factors that are used to measure the effectiveness of an advertisement or the popularity of store is the conversion rate. This number represents the percentage of the people actually bought something out all the people that visited the store. Usual way of calculating the conversion rate is by counting the people coming to the store and checking out at one of the counters. It is impractical to do this manually in a large retail store environment. Other than calculating the conversion rate, various business decisions could be made for the betterment of the store by analyzing customer behaviour.

Enhancing customer experience by rearranging the store structure, efficiently allocating staff to high density areas, understanding traffic trends and identify "opportunity hours" and schedule labor hours accordingly to optimize service levels are some of such business decisions. In order to assist making business decisions, large-scale retail stores need to count and track people using a number of ceiling-mounted cameras. In such environment it is hard to keep track of customer behavior separately using a manual system.

Identifying the history of customers and understanding the customer density inside the areas of the department store can be critical for financial and business decisions. However if a system is in place which counts and tracks the customers inside the store and generates basic intelligence to the managerial level, it can expedite the decision making process. This project is aimed at developing such a system. We have implemented a people tracking and counting system adaptable to any large scale store structure which is able to process multiple video streams obtained through separate cameras mounted on the store and generate business intelligence.

There have been several systems which handle multi camera people tracking introduced in the past decade. All these systems incorporate a central server which will receive all the frames from all the cameras to handle the people detection and tracking tasks real-time. Major drawback of these systems is the scalability. When the number of cameras increase, the bandwidth usage for sending video frames to the central server will be increased and the processing power required at the central server will increase. This raises the requirement for processing video frames close to the camera which is defined as leaf node processing.

The scope of this project is to implement a scalable end to end system for Multi-Camera People Tracking utilizing a Zynq SoC for leaf node processing. Thereby we use an

FPGA + ARM processor for leaf node processing, a central server for multi-camera people tracking and a web server for generating business intelligence.

In this proposed system by using a FPGA + ARM processor we have achieved leaf node processing of people detection at a lower power and also enable the scalability and lower bandwidth requirement.

People Detection in FPGA is an aspect that has been researched by many over the years. For example [2] suggests a FPGA based embedded platform for people detection using Gaussian Mixture Models and [3] suggests a similar system for people detection using Histogram of Features.

Multi-Camera People Tracking systems is also a field which has been researched over the years. [4] suggests a multi-camera surveillance system implemented with background subtraction for detection and color features for tracking. [6] suggests another multi camera people tracking system which uses particle filters for tracking on the ground plane.

Chapter 2

LITERATURE REVIEW

In this project main focus was to implement an end to end system for Multi-Camera People Tracking System and to use this system to generate business intuition. As explained earlier we use the technique of leaf node processing to achieve a scalable, low power, low bandwidth required system.

There are several challenging tasks in implementing this system. One is people detection in FPGA. There are various ways to achieve this. One way is to use a blob detection algorithm as used by Vicente, Alfredo Gardel, et al. [1] in 2009. They have done background subtraction followed by contour detection to select head candidates for videos obtained using overhead cameras and they have implemented people detection part in a low cost FPGA (spartan3).

Feature based methods are also another way to do people detection. In these methods a set of features are calculated around a window and some kind of machine learning algorithm is used to train a classifier that can classify each window into "contains a person" and "not contains a person". One such method is Histogram Oriented Gradient (HoG) based method suggested by Dalal, Navneet, and Bill Triggs. [2] in 2005. A variation of HoG is implemented on FPGA by Negi, Kazuhiro, et al. [3] in 2011.

Another approach that can be used for people detection is to utilize a neural network. If we can come up with a suitable neural network architecture, this method has the potential to outperform all the other methods mentioned earlier. There are several open source implementations of such architectures available. One such architecture is given by Redmon, Joseph, et al. [4] in 2016. Their work has made a huge leap in object detection space. Due to its success their architecture is utilized by many people. Challenge though in using this approach is the complexity of implementing a convolutional neural network in an FPGA. Implementing a CNN architecture in FPGA with the available resources and to achieve real time processing is an immense challenge.

We used a background subtraction based method for people detection in the proposed system. Though we considered using either HoG or CNN, as our focus was to implement the end end architecture for multi-camera people tracking with leaf node processing, background subtractions deemed to be a feasible solution. We could improve the existing system by implementing either HoG or CNN as a future improvement.

Other major challenge we have is people tracking based on our detections. There are many traditional approaches to this such as tracking using kalman filter and assigning detections to tracks using the Hungarian algorithm. But multi target tracking also can be formulated as a discrete-continuous optimization problem [5]. Andriyenko, Anton, Konrad Schindler, and Stefan Roth [5] in 2012 have considered data association as a discrete optimization problem with label costs. And they have posed trajectory estima-

tion as a continuous fitting problem with closed form solutions. And there method has performed well on some known data sets.

Other major challenge that we face is, how to extend multiple target tracking into a multi camera framework. There are some research papers discussing this issue. Tang, Nick C., et al. [6] in 2015 suggest a two pass regression framework to solve this challenge. First pass regression predicts the people count based on the features calculated from intra-camera video frames. And then the second pass is based on the conflicts between the prediction derived from multiple views. They have formulated this as a transfer learning problem. Yang, Tao, et al. [7] in 2007 discusses another method to do this. In this approach first they do single camera tracking and they transform these tracked paths into a global coordinate system. Then by using their multi camera handoff algorithm, they can track people across multiple cameras. Here they calculate a match score for an object appearing in a camera under overlapping or non overlapping conditions for all tracks under all cameras. The track having the maximum score is selected.

By going through the available literature we were able to get an insight of the scope of our project. We identified some potential solutions for the challenges we have. We were also able to identify that implementing the overall system architecture for multi-camera people tracking with leaf node processing will be a novel contribution.

Chapter 3

METHODOLOGY

3.1 Introduction

In this project we have designed and implemented an end to end multi camera people counting and tracking system based on Xilinx ZC702 development board. In this implementation we have considered the possibility of using leaf node processing to develop a scalable solution for multi camera people tracking. In this chapter we present the underline methodology of our system in detail. In order to give the reader a gradual walkthrough of our project we will explore each component of our system separately. First we will give an overview of the basic system. Then we will explore in detail the each component of our system.

3.2 Scope

The main components of the project that makes the scope is as follows.

- People Detection feature extraction
This involves reading image frames from the camera and extracting the features required for people detection. This computation will be done in the Programmable Logic (PL) of Zynq ZC702 Evaluation Kit.
- Writing Linux Kernel drivers for controlling PL IP cores.
Writing Linux Kernel drivers are required to control the IP cores from the software end.
- Communication with the backend server
Calculated features for each image frame will be sent to the backend server through Ethernet for the rest of the processing.
- Multi- Camera People Tracking
Multi-Camera people tracking will be done in the backend server using the features sent through each of the cameras.

3.3 Overall Architecture

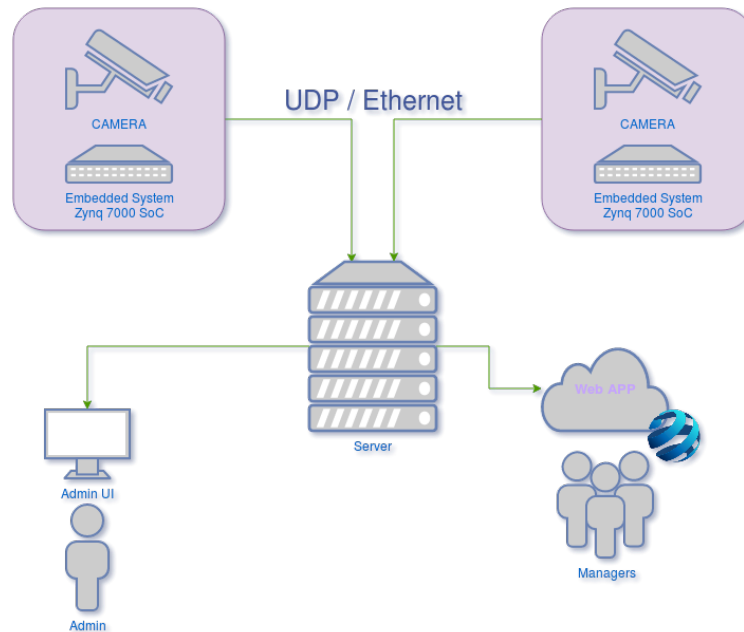


Figure 3.1: Overall block diagram of the end-to-end system

Figure 3.1 shows the overall architecture of our system. It mainly consist of 3 components, namely, An embedded system closer to the camera (Leaf node) People tracking applications running on a server Business intelligence generation and displaying interfaces.

We will explore each of the above components separately in the following sections.

3.4 Embedded System (Leaf Node)

Generally the term 'leaf node' is used for a system placed at the edge or bottom of a hierarchical network (of systems). In our system, leaf node is the embedded system connected to the camera.

This embedded system is responsible for capturing live video frames from the camera and preprocessing before sending into the central server. Objective of preprocessing on the leaf node is to reduce the required bandwidth of the network to send information (video in this case) to the central server and to reduce the processing done one the central server in a scalable manner. In a large system consisting of few hundreds of cameras, reduction of the bandwidth required and the processing power of the central system for each camera makes a big difference in the central system processing and bandwidth requirements, making the system more scalable.

In our prototype system, leaf node consists of a web camera connected to a Xilinx

ZC702 development board and its responsible for doing people detection of the live video feed and feature calculation of the detected people. Implementations of the people detection and feature calculation algorithms are explained in detail in later sections. These calculated features of the detected people are then sent to the central server via a custom application layer protocol based on User Datagram Protocol (UDP) / Internet Protocol (IP) over Ethernet.

Following sections provide an overview about the feature of Xilinx ZC702 board and ZYNQ-7000 All Programmable (AP) System on Chip (SoC).

3.4.1 Xilinx ZC702 Board

Xilinx ZC702 board, shown in Figure 3.2, is a development board manufactured by Xilinx targeting embedded hardware (FPGA) design development. This board consists of a ZYNQ-7000 AP SoC with various other peripherals. Key features and peripherals of Xilinx ZC702 is listed below.

- Zynq-7000 SoC
- 1GB DDR3 Component Memory
- Enabling serial connectivity with USB OTG, UART, IIC, CAN Bus
- Ethernet which supports 10-100-1000 Mbps transfer rates
- HDMI interface
- FPGA Mezzanine Card (FMC) interface
- Onboard Secure Digital (SD) card reader



Figure 3.2: Xilinx ZC702 Development Board

Features of the onboard ZYNQ-7000 SoC and usage of it in our project is explained in following sections.

3.4.2 ZYNQ-7000 SoC

ZYNQ-7000 SoC consists of a Xilinx XC7Z020-CLG484-1 FPGA and two ARM Cortex-A9 core processors. Block diagram in figure 3.3 shows how the Programmable Logic and the Processing System is connected inside the ZYNQ-7000 SoC.

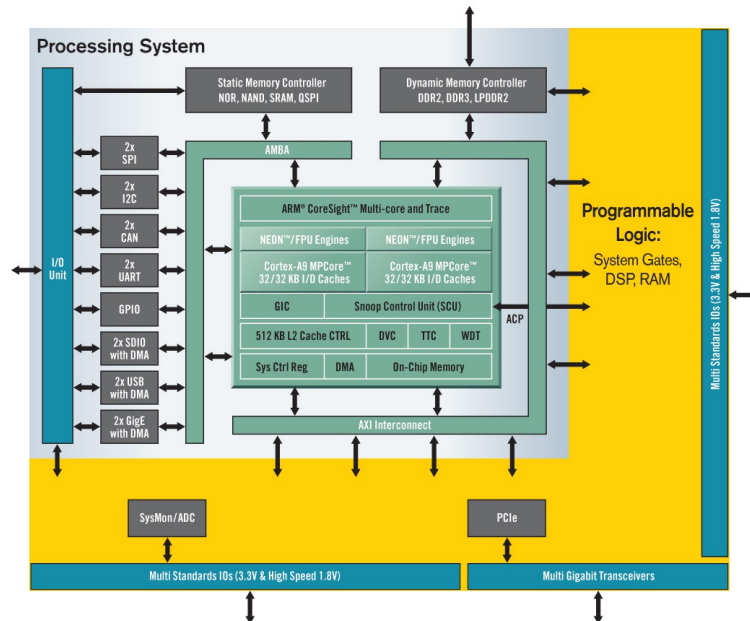


Figure 3.3: Internal architecture of ZYNQ-7000 SoC

In our implementation, we have installed a lightweight Linux distribution on dual ARM processors in ZYNQ-7000 SoC. Reasons for installing Linux and its usage is explained in a latter section.

3.5 Design Flow for ZYNQ-7000 based Embedded System

We used the standard Xilinx toolchain for developing an embedded system for designing and implementing our leaf node system. Figure 3.4 shows a block diagram of the design flow with the design tools we used.

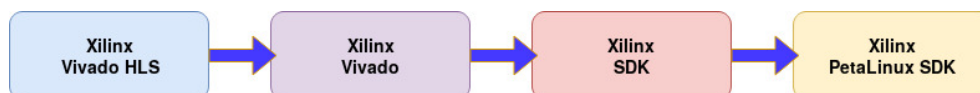


Figure 3.4: Design flow block diagram for ZYNQ-7000 based embedded system

Brief description about the usage of each tool in the design flow is included in the following subsections.

3.5.1 Vivado HLS

HLS stands for High Level Synthesis. Task of Vivado HLS is to synthesize a function developed using a high level language (C, C++ and SystemC) into a RTL design (Hardware Descriptive Language Code) of an IP Core.

3.5.2 Vivado

IP cores designed using Vivado HLS are then imported to Vivado for designing the overall architecture. This is the tool where we decide how the PL and PS parts of ZYNQ SoC is connected (which type of AXI protocol). After overall system is designed, its synthesized and implemented on the target device. As an output of this process, bit-stream file and hardware description (HDF) file is generated.

3.5.3 Xilinx SDK

HDF file generated is then imported to Xilinx SDK for testing the functionality of implemented hardware. Xilinx SDK has options for testing the hardware on the target device by writing a simple C/C++ application code for testing the functionality of the implemented hardware.

3.5.4 PetaLinux SDK

After confirming the functionality of the hardware on the target device, HDF file is then imported to a PetaLinux SDK project for generating the Linux boot files for the custom hardware. Enabling the required kernel modules and editing the device tree to enable hardware peripheral access is done here. More details about running Linux on ZYNQ PS is explained in the next section.

3.6 Installing Linux on ZYNQ-7000

Lightweight customized Debian Linux distribution was installed on ARM processors of ZYNQ SoC. When it comes to running applications on PS of ZYNQ, there are two primary options,

1. Running a Linux distribution
2. Running the application in baremetal mode (Using the basic drivers provided by Xilinx)

We choose to install and run a Linux distribution based on the following reasons,

- Complete and working protocol stacks for communication protocols (Ethernet, UART)
- Readily available device drivers for interfacing various hardware (USB web camera, USB tethering)
- Root file system stored in the SD card (Usually bare-metal application use on-board RAM to store data at run time, which will be destroyed after turning off the board)
- Linux comes with a package manager that can be used to install almost any third party Linux compatible software package (Python, C++ compilers).

Installing a third-party popular Linux distribution on ZYNQ PS, for a custom hardware, was a challenging task due to the lack of up to date documentation and community support. We managed to gather information from various Xilinx sources and we figured out the correct way of installing a custom Operating System(OS) on ZYNQ PS. Basic steps for installing a Linux distribution on the ZYNQ PS is listed below.

1. Generating Linux boot files for custom hardware using PetaLinux SDK with the required kernel modules and the device-tree configurations
2. Partitioning a SD card (at least 2GB) to a 'boot' partition and a 'rootfs' partition
3. Copying custom Linux root file system into 'rootfs' partition of the SD card
4. Copying Linux boot files into 'boot' partition of SD card
5. Changing the boot configuration of the board to boot from SD card

3.7 People Detection and Feature Calculation Implementation on ZYNQ-7000 SoC

As mentioned before, people detection and feature calculation is done on the leaf node, in other words, on ZYNQ-7000 SoC. This implementation utilizes both the processing system (PS) and Programmable Logic (PL) parts of ZYNQ-7000. While the main application runs on top of the Debian Linux OS on ZYNQ PS, it also contains code snippets for capturing video frames from USB web camera, controlling the IP cores and send detections and calculated features to the central server. Architecture of our implementation within ZYNQ SoC can be seen from figure . Each of the task implemented within PS part and the PL part of ZYNQ SoC is explained in next sections.

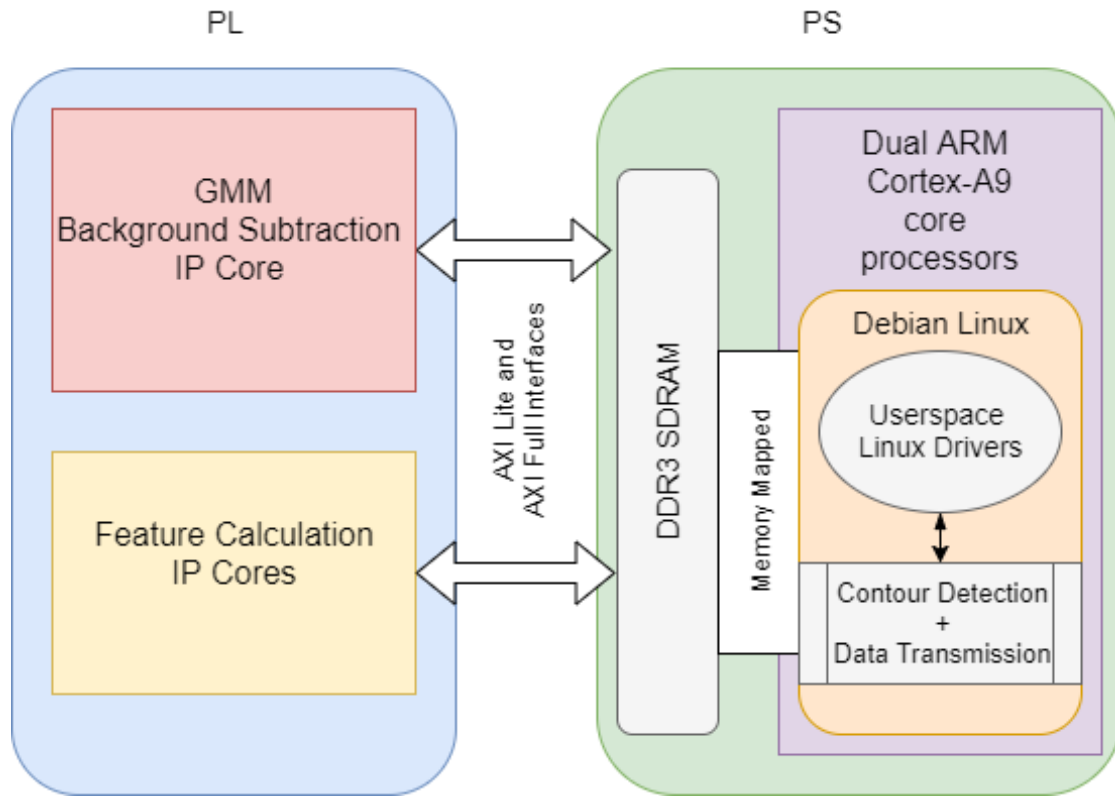


Figure 3.5: Architecture of the Leaf node

3.7.1 IP core implementation on ZYNQ PL

Gaussian Mixture Model Background Subtraction IP Core

Gaussian Mixture Model Background Subtraction IP Core is implemented using Vivado HLS tool. Algorithm is explained as follows.

Gaussian Mixture Models uses a statistical model in which the probability distribution of the luminosity of each pixel is modeled with a mixture of Gaussian distributions. GMM performs really well in identifying the background even in luminosity changes. Our implemented system contains GMM IP core based on the OpenCV GMM algorithm. The algorithm is optimized for hardware implementation by sectoring the received frame into 60 parts and processing it separately such that the latency and the BRAM utilization is minimized.

The proposed system models a pixel as a mixture of 2 gaussians. We have limited the system to two models, as a considerable increase of accuracy was not observed by increasing the no of models compared to the increase in resources and latency. A very low latency is desirable for this application as the process needs to be real time. Each model is represented by 4 parameters: mean($\hat{\mu}$), variance($\hat{\sigma}^2$), weight(w) and Fitness value (F). The weight can be defined as the confidence that a certain gaussian occurs.

These parameters differ for each gaussian of each pixel and are updated for every incoming frame of the video.

The fitness value is defined as following,

$$F(p, k, t) = \frac{w(p, k, t)}{\sigma(p, k, t)} \quad (3.1)$$

Where p - index of the pixel, k - index of the gaussian distribution and t - frame number. As we employ only two models we compute 2 F values for a single pixel in a frame. The matched condition is checked with the following condition.

$$M(p, k) = 1 \quad \text{if} \quad |pix - \mu(p, k, t)| < 2.5\sigma(p, k, t) \quad (3.2)$$

The above condition checks whether the said pixel can be categorized as a background pixel. As can be seen this condition could be true for more than one gaussian. Thus, the matched distribution is taken as the distribution with the highest F value. The parameters of the matched distribution are updated as follows.

$$\mu(p, k, t + 1) = \mu(p, k, t) + \alpha(pix - \mu(p, k, t)) \quad (3.3)$$

$$\sigma^2(p, k, t + 1) = \sigma^2(p, k, t) + \alpha[(pix - \mu(p, k, t))^2 - \sigma^2(p, k, t)] \quad (3.4)$$

$$w(p, k, t + 1) = w(p, k, t) - \alpha w(p, k, t) + \alpha \quad (3.5)$$

pix is the value of the pixel. The learning rate α defines the rate that mean variance and weight are updated. parameters of the unmatched gaussians are not updated. For the first frame $\alpha = 1$, so that the background is initialized from the first frame. Therefore this algorithm heavily depends on a correct background image being achieved through the first frame.

If no appropriate gaussian mixture found at all, that is no gaussian distribution for a certain pixel confirms the first condition then the weakest mixture is removed and a new one is created. The weakest mixture is defined as the mixture with the lowest F value.

$$\mu(p, k, t + 1) = pix \quad (3.6)$$

$$\sigma^2(p, k, t + 1) = \sigma^2 \quad (3.7)$$

$$w(p, k, t + 1) = w \quad (3.8)$$

σ^2 is the default variance, w is the default initial weight initialized at 0.05. If no appropriate gaussian is found for a certain pixel then that pixel is categorized as a foreground pixel for that frame. However, if a gaussian satisfies matched condition for a pixel, then that pixel is not straight away assigned as a background pixel. A thresholding

condition is checked to see whether the matched gaussian distribution is a background gaussian or not. Fitness values are sorted from maximum to minimum and weights of the respective gaussians are added in succession according to the following equation.

$$\min \sum_{k=1}^b w(p, k, t) > T \quad (3.9)$$

Where $T = 0.7$ is the default background ratio. The minimum number of gaussians that satisfies the above equation are categorized as background distributions. If the matched condition of the pixel is satisfied by a gaussian which also satisfies the thresholding condition then it is a background pixel, else its a foreground pixel.

Feature Calculation IP Core

An unique feature for each detected person was needed for tracking each person. We have selected 3D RGB color histogram as the feature based on following criteria,

- A simple feature with low computational complexity was needed to preserve the real-time nature of the end-to-end system.
- In the software simulation using above mentioned feature for people tracking, fairly good accuracy was observed.

Figure x shows how the 3D histogram is calculated. Input to the feature calculation IP core is a image frame with each pixel's RGB components, detected bounding box and background subtracted binary mask. Algorithm of feature calculation IP core is implemented such that, while iterating through each pixel, it checks two conditions,

1. Whether current pixel is inside the provided bounding box
2. Whether current pixel is a foreground pixel

If any of the above conditions fail, algorithm moves on to the next pixel. Otherwise that pixel is proceeded through to the next step.

In the next step, for each pixel, RGB components are divided by 8. Then these scaled down RGB components are used as the indexes for finding the relevant location of the histogram for that pixel. After relevant coordinate is located on the histogram, it will be incremented by 1.

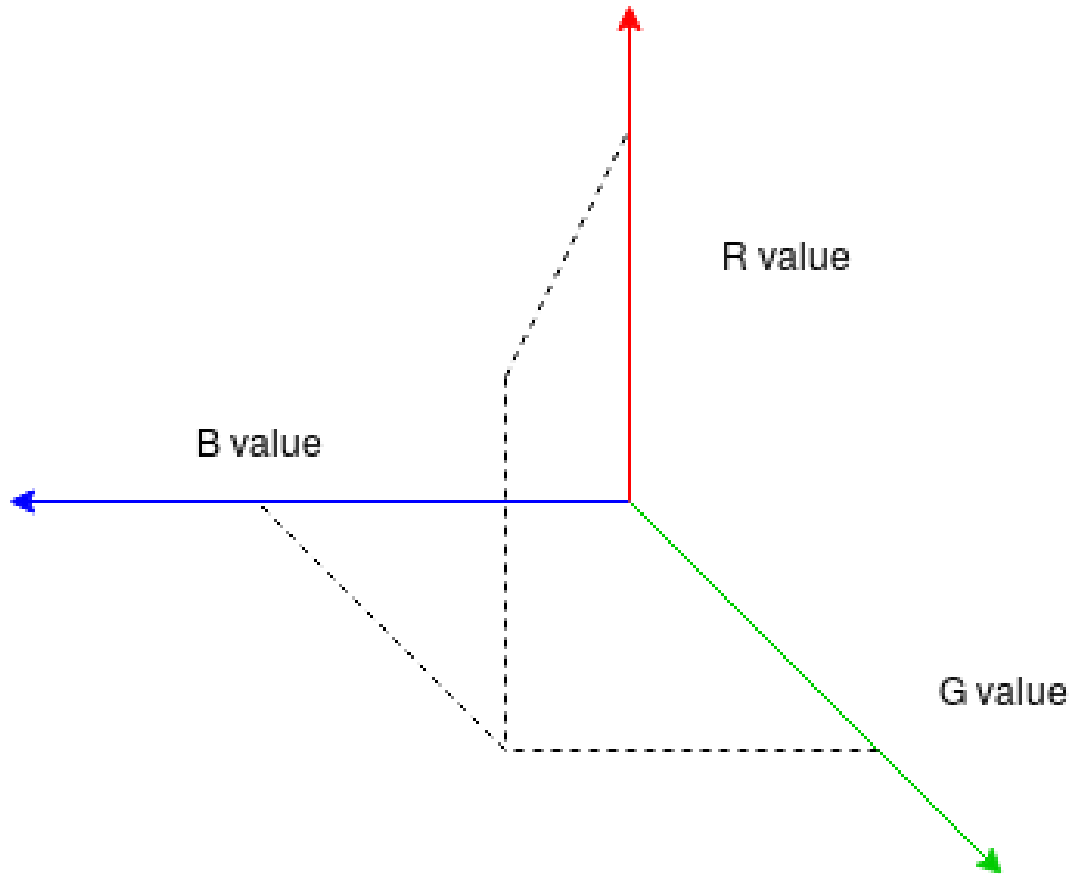


Figure 3.6: Structure of the 3D RGB histogram

Since RGB components of each pixel scaled down by 8, each dimension of the 3D histogram consists of 8 bins ($256/8$). Therefore result of the above process is a $8 \times 8 \times 8$ 3D histogram for each detected person. Output of the feature calculation IP core is a 512 element linear array which contains the flattened elements of 3D histogram for the detected person.

3.7.2 IP Core drivers, Morphological Operations and Communication Protocol Implementation on ZYNQ PS

Linux Drivers for controlling IP cores

This was also a challenging task due to the lack of upto date documentation and the complexity of physical memory accessing through Linux kernel. But we managed to overcome this task by using a simple but effective method.

When IP cores are generated using Vivado HLS, tool also generates the following functions, in C language, for controlling the IP core.

- Starting the IP Core

- Check whether the IP core is busy or idle
- Setting input and output addresses

But these functions cannot be used directly in the Linux side of ZYNQ PS for controlling the IP cores in the PL, because physical DDR3 memory cannot be directly accessed from Linux OS. Linux only allows the usage of virtual memory in the userspace of Linux OS. 'mmap' function in C++ language was used to overcome this barrier.

What 'mmap' function does is, it maps a specified region in the physical DDR3 memory onto virtual memory, which can be accessed from the Linux OS userspace level.

Input/output AXI busses of the IP cores implemented in the ZYNQ PL is mapped onto specific memory locations in the DDR3 RAM. An userspace driver was developed to map these control memory locations and input/output memory locations into virtual memory. Then C functions generated by Vivado HLS were adapted to use these virtual memory to access the physical memory locations relevant to IP cores in an indirect manner. This is shown in figure x as the mapped DDR memory to Debian Linux.

People Detection from Background Subtracted Binary Mask

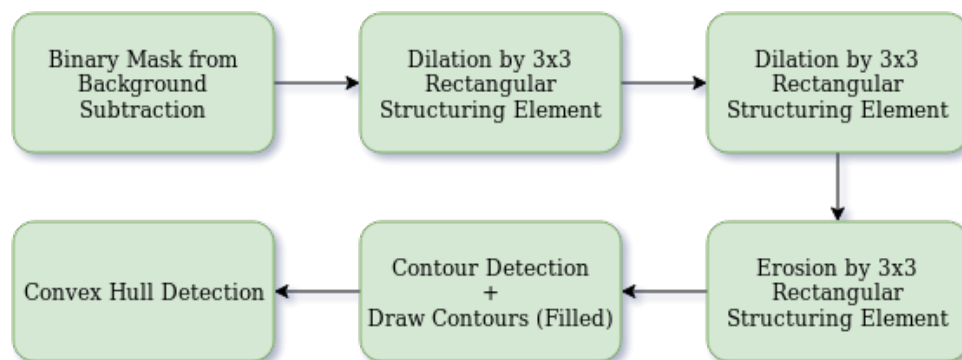


Figure 3.7: Block diagram of the morphological operation flow to detect people from background subtracted binary mask

Once the background subtracted binary mask is obtained a set of operations are performed as shown in figure 3.7 to enhance the binary image for object detections. First a set of morphological operations are performed to clean out the image as necessary. Since the quality of the background model can vary due to various environmental conditions, background subtraction will not give a cleanly filled blob for the shape of a person. There will be a lot of black pixels on the boundary as well as inside the shape of a person. We can use the morphological operation ‘closing’ to fill out these holes. It is essentially two dilations followed by an erosion. Structuring element used here is a 3x3 rectangle. We have tested out several alternatives for the kernel but this

kernel size and shape gave the best results for the sequences of images we have tested. Therefore we used it in our final implementation. The results of morphological operations for a test images is shown in figure 3.8c, 3.8d and 3.8e.

Morphological operations will give a cleaner image compared to the original background subtracted image. Next contour detection is applied to the output image resulting from morphological operations. We have used inbuilt functions of opencv for this. They have implemented contour detection based on Suzuki's [7] border following algorithm. Then again using inbuilt functions of opencv we fill out these contours. It will completely fill the shape of the person. The result of this operation for a test image is shown in figure 3.8f.

Next step in this process is to detect the bounding box surrounding the blob of a person. For this we use opencv inbuilt function for convex hull detection. A convex hull is the smallest polygon enclosing a set of points. Opencv's implementation is based on the algorithm proposed by Sklansky [8]. After this process we obtain a set of bounding boxes corresponding to people detected. Resulting bounding boxes for a test image is shown in figure 3.8g.

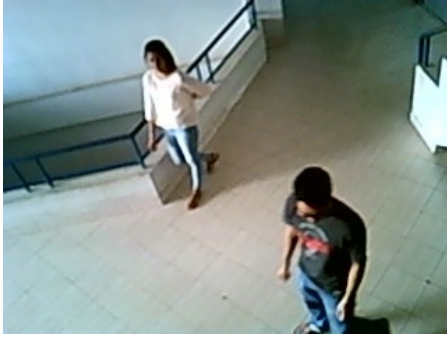
3.7.3 False Positive Reduction

Although the above process is able to detect people from an image, it also detects a lot of unwanted blobs due to various reasons such as reflections caused by background surfaces, imperfections of background modeling, etc. Therefore we need to utilize a method to reduce number of these false positives. Our initial approach was to select the bounding boxes satisfying a set of conditions as follows and reject all the other detections.

- Width of bounding box(w) > TW
- Height of bounding box(h) > TH
- $TAR1 > \text{Aspect ratio}(w/h) > TAR2$
- $TCA1 > \text{Contour Area} > TCA2$

Here TW, TH, TAR1, TAR2, TCA1 and TCA2 are constants that are chosen arbitrarily. Initially we tuned up these constants manually by looking at the results. But this was a tedious task and this has to be redone once the camera position is changed. Therefore we explored the possibility of machine learning based automatic way to do this. We selected a few features of the bounding boxes into consideration as follows:

1. Width
2. Height



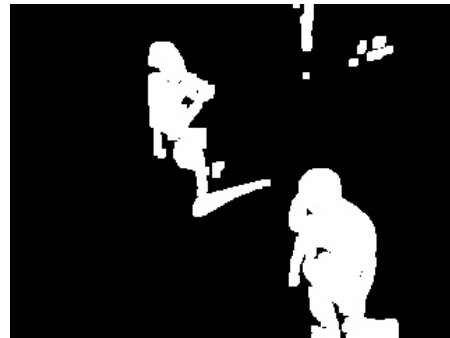
(a) Original image



(b) Background subtracted binary mask



(c) Dilation performed on the binary mask



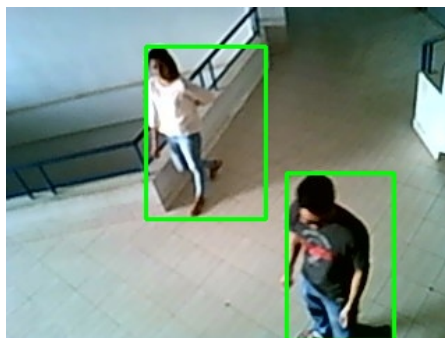
(d) Twice dilated



(e) Twice dilated and eroded



(f) Contour detection and drawing the filled contours taking 3.8e as input



(g) Detecting bounding boxes using convex hull detection

Figure 3.8: Results of morphological operations for a test image

Table 3.1: Percentage of variance along each principal component

Principal Component	Percentage of Variance Explained
1	99.4509
2	0.5401
3	0.0063
4	0.0019
5	0.0000
6	0.0000

3. Bounding box area
4. Aspect Ratio (Width / Height)
5. Contour Area
6. Diagonal Length

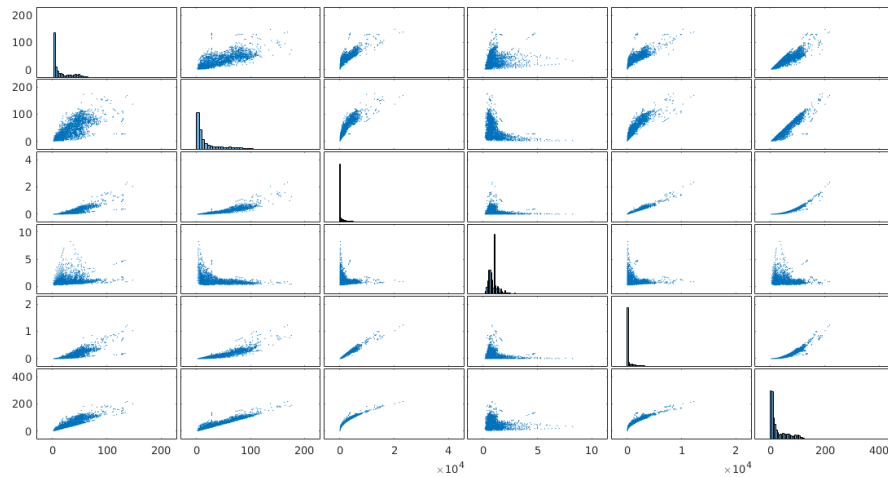


Figure 3.9: Scatter plot matrix for the 6 features collected from running our object detection on a sample video

We can see from the scatter plots above that most of the above features are highly correlated. Therefore we can use Principal Component Analysis to perform dimensionality reduction of the feature space. The percentage of variance explained by each principal component was obtained as follows.

We can see that only the first two principal components give a significant contribution. Even the 2d principal component is less significant compared to the 1st principal component. We can further see this in figure 3.10.

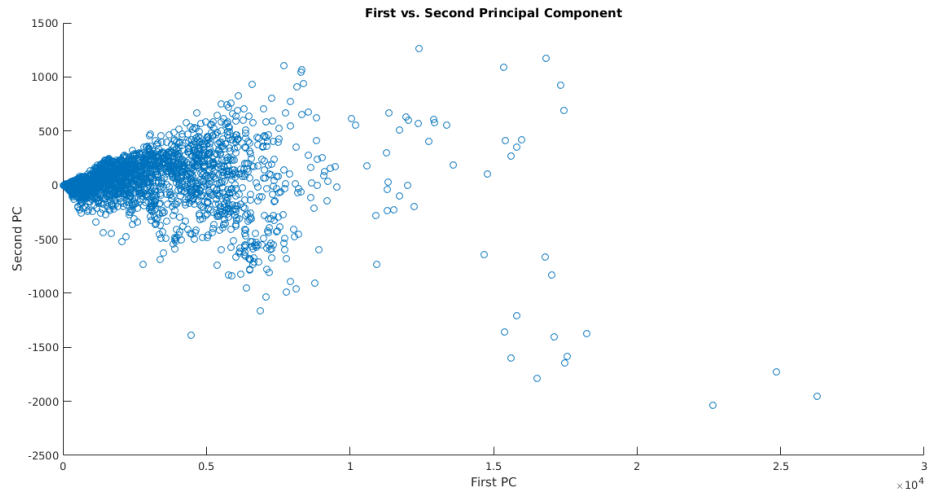


Figure 3.10: First Vs. second principal component

Therefore it is justifiable only to select the first principal component into consideration.

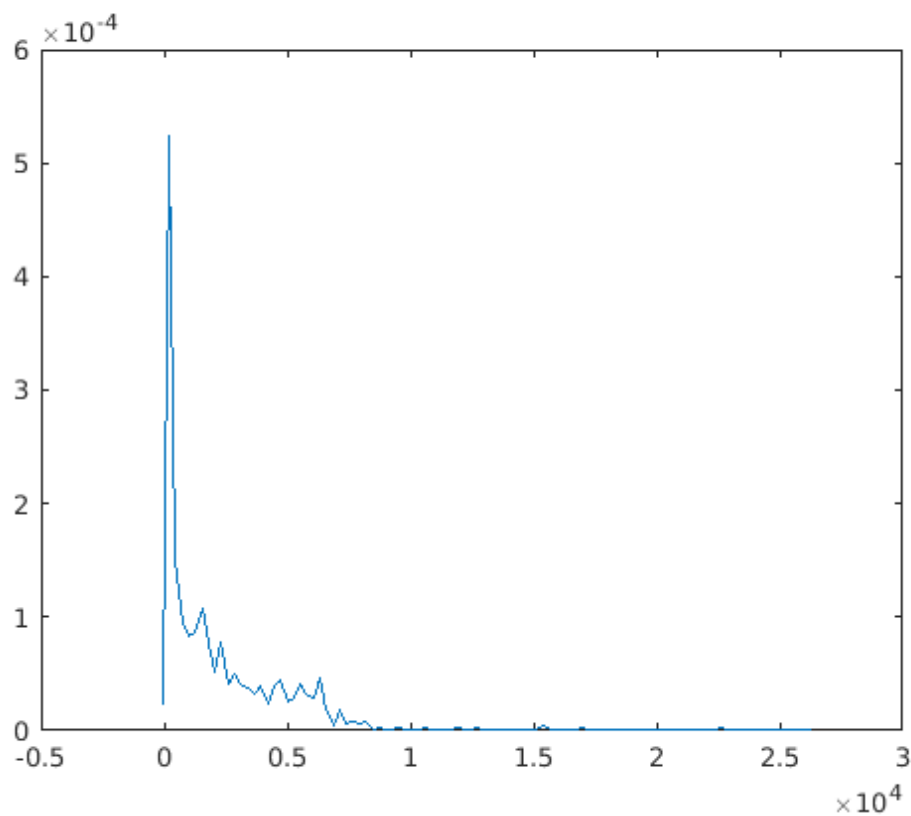


Figure 3.11: Density of the first principal component

Next we looked at the kernel density estimate of the first principal component as shown in figure 3.11. We can approximate this with a bimodal distribution. Therefore the problem reduces to finding the optimum threshold for this bimodal distribution. We

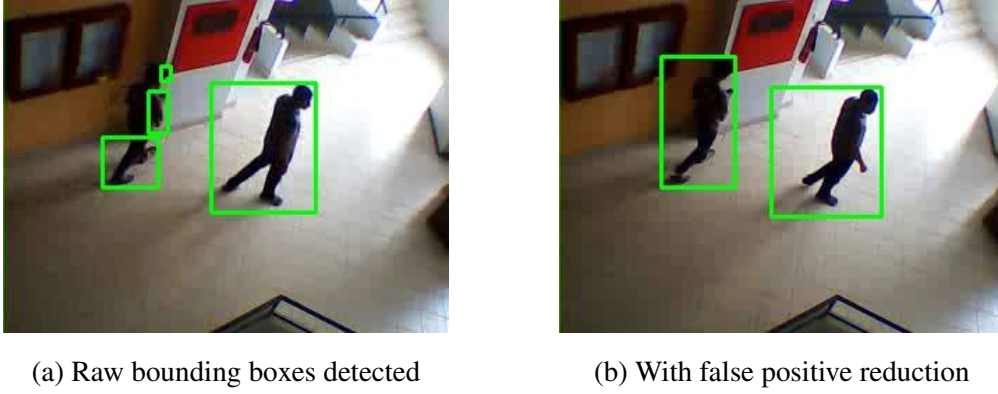


Figure 3.12: Results of PCA based false positive elimination algorithm

used Otsu’s [9] thresholding algorithm to find this threshold.

We utilized this framework for false positive reduction of the detections. Once we tested out this with several video sequences taken from different camera angles we found out that this method indeed gives better results as analysis suggested. Figure 3.12 shows the results of PCA based false positive reduction.

When applying this method the system should first undergo a training stage where raw detection features are collected. Then coefficients of the first principal component is determined using this dataset. Then the optimum threshold is determined. This whole process was implemented in code so that system can automatically learn to reduce false positives. In the normal mode of operation the system directly uses these calculated coefficients and threshold for false positive reduction.

3.7.4 Communication Protocol among leaf nodes central server

Since our system is a multi camera system we need to connect each camera to a central server where people tracking is done. we need to design a communication algorithm to handle the communication between the server and camera nodes. We identified the following requirements for designing such a communication protocol.

1. Low bandwidth utilization
2. Low latency
3. Scalable with increasing number of cameras

Taking these requirements into consideration we designed a communication algorithm. We used User Datagram Protocol (UDP) to send the data packets, since it has lower overhead and it is also connectionless. In our architecture the backend server will be listening to packets coming from cameras over the network. Each camera node can

just start sending data packets to the server at anypoint in time. There is no connecting phase. Server will distinguish packets from different cameras and act accordingly. The advantage of this kind of a setup is we can plugin any camera to the network quite easily and switch it on and start sending data.

Each datagram sent from a camera will contain the following information.

Field	Description
Camera ID	A number between 0-255 identifying the camera.
Time Stamp	Time at which this frame was generated.
Detections	Bounding boxes of the detections.
Histograms	512 bin Histograms corresponding to detections.
Binary Mask	Binary mask obtained from background subtraction is also sent for display purposes. This can be completely eliminated once the system is deployed.

With this information the server can distinguish packets coming from different camera nodes. And also it can identify any missing frames. But in our current implementation we discard missing frames completely. This information is converted into a byte stream using serialization libraries of boost (a c++ library). Then this byte stream is sent to the server through an UDP packet of 49152 bytes. We have chosen this value such that each UDP packet is self contained (i.e. it contain all the information necessary to process the corresponding frame), therefore no fragmentation of UDP packets is necessary.

3.8 Central Server

In our system, people detection is done at each camera node. Then each camera send this information to the central server. Multi camera people tracking and counting is performed at the central server. This central server perform the following major tasks.

1. Receiving and unloading data packets from camera nodes
2. Single camera tracking
3. Global multi camera tracking
4. Communication with the web server

3.8.1 Receiving and Unloading Data Packets from Camera Nodes

As mentioned previously we have designed a scalable communication protocol to connect camera nodes to the server. In our implementation we have a UDP server running

on a dedicated thread. It listens to any incoming packets through the network. Packets are received, deserialized and pushed into a queue. All the frames in the queue are processed at a interval of 100 ms.

3.8.2 Single Camera Tracking

When it comes to object tracking, detection based tracking methods are the most popular. Multiple object tracking can be decomposed into two parts as, data association and target tracking. These multi object tracking algorithms can be divided into two categories. The first category relies on past frames to estimate the current state recursively. The second category allows for a certain latency and globally solves for all trajectories within a given time window.

Since our system is to be run in real time, methods under first category are more appropriate for us. Under the first category we tried out Hungary algorithm followed by Kalman Filter based tracker for each object being tracked. As an alternative we tried out particle filter based tracking also. Out of these two we decided to use particle filter based tracking in our implementation due to reasons we will discuss later.

Kalman Filter based Tracking

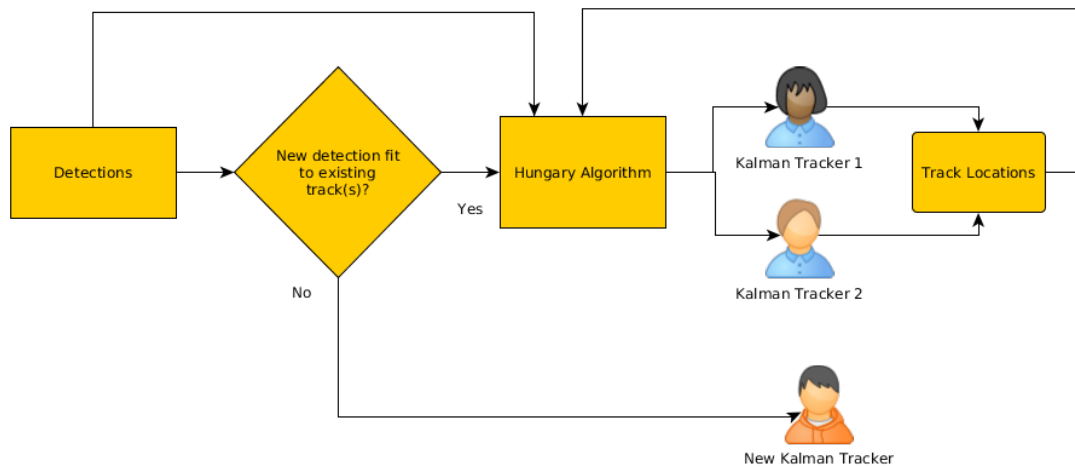


Figure 3.13: Multiple People Tracking using Kalman Filter

Here detection coordinates are fed to the backend tracking system by individual camera nodes (camera + zynq zc702). The next part of the algorithms calculates a cost matrix for detections vs existing tracks. For m detections and n tracks an entry of the cost matrix is given by equation (3.10).

$$c_{ij} = \text{cost}(\text{detection } i, \text{track } j) \quad \text{where } i = 1, \dots, m; j = 1, \dots, n \quad (3.10)$$

Then detections having the minimum track cost lower than a threshold value are kept and detections exceeding the threshold value are initialized as new tracks.

Here each track is a kalman filter. Dynamics model used in the kalman filter is a constant velocity model consisting of 6 state variables, namely {x coordinate, y coordinate, width, height, x velocity, y velocity} and 4 measurement variables, namely {x coordinate, y coordinate, width, height}.

Detections to track assignments are done through the Hungarian algorithm. The Hungarian method is a combinatorial optimization algorithm that solves the assignment problem in polynomial time. Here the individual cost is the euclidean distance between x,y coordinates of the detection and the track. However this can be modified to include the width and height as well.

Tracks which have not been assigned to a detection within consecutive frames greater than a threshold value are simply deleted.

Although the kalman filter based tracking gives somewhat satisfactory results, we decided to try out particle filter based tracking also.

When we look at results we identified several weak points in our current tracking scheme. So we applied some modifications.

Improvements for Kalman Filter based Tracking

In the data association step we have only considered the euclidean distance between detection and tracker coordinates for cost estimation. We improved the accuracy of data association by applying following modifications.

Gray Level Intensity Histogram as a Parameter in Cost Estimation

We obtained the gray level intensity histogram for all the detection locations. And we included a histogram parameter in each tracker, where the histogram of the tracker is modified as follows at each assignment of a detection.

$$histogram_{tracker} = \alpha * histogram_{detection} + (1 - \alpha) * histogram_{tracker} \quad \text{where } 0 < \alpha < 1 \quad (3.11)$$

Then the correlation between the detection histogram and the tracker histogram is calculated and its inverse ($1/correlation\ coef.$) is added to the cost. This makes our cost sensitive to the similarity of detection and tracker regions.

Adding a Penalty for the Cost for High Velocities in Kalman Filter State

When the kalman filter tends to diverge the velocity value in the state vector becomes high. We can discard the kalman filter quickly by adding a penalty to the cost when the velocity is greater than a certain threshold.

Tuning model parameters

We have the following model parameters that must be properly tuned. Currently this is done in an ad hoc manner.

- Track Initialization Threshold - When the cost exceed this value a new tracker is initialized
- Rejection Tolerance - When the count of a detection not being assigned to a tracker exceed this value the tracker is deleted
- Velocity Threshold - Penalty is added to the cost when tracker velocity exceeds this value

3.8.3 Particle Filter Based Tracking

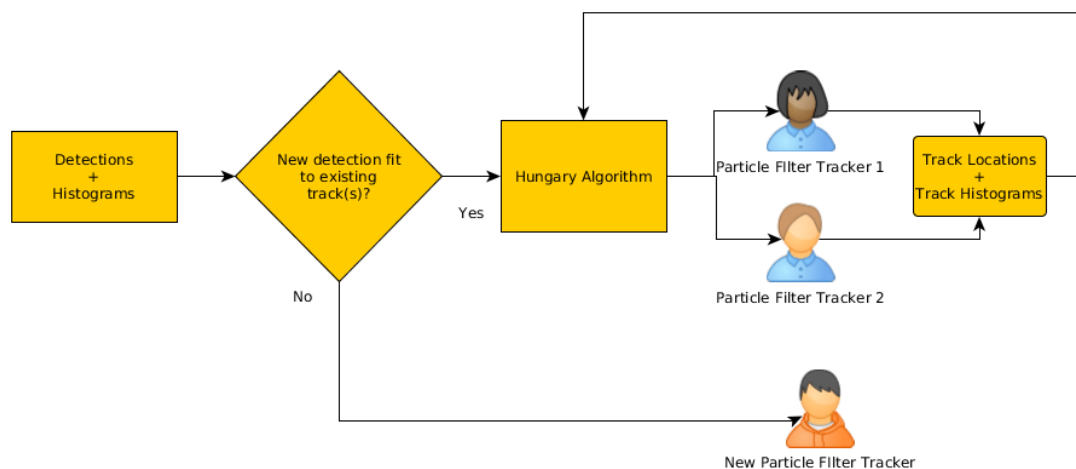


Figure 3.14: Particle Filter based Tracking

Target tracking system comprising of particle filter trackers are shown in figure 3.14. Arrangement is the same as the one with kalman filter trackers. Only difference is in using particle filter based trackers instead of kalman filter based ones.

In this setting the particle filter tracks the lower center point of a bounding, which corresponds to the foot of a person. We are using a constant velocity model for particle propagation. State of a particle consist of 4 states (x , y , v_x - velocity in x direction, v_y

- velocity in y direction). Out of these v_x and v_y are hidden states. A particle tracker undergoes following 3 stages.

1. Particle initialization
2. Particle propagation
3. Updating weights with a measurement

In particle filtering state space is represented as a distribution of particles. Each particle has a weight which corresponds to a probability of this particle being the through location. We will briefly go through his 3 stages of particle filtering.

Particle Initialization

A new particle tracker is initialized when a new object is detected for the first time. At this stage x , y are initialized to that of the detection location and v_x , v_y are initialized to 0. Weight of all the particles are initialized to $1/(no. of particles)$.

3.8.4 Business Intelligence Software

Business intelligence software provides the functionality as a website in order to allow access on the go. The web interface is based on Flask web Framework. Flask which is based on Python is ideal for interactive real time user interface which showcases graphical data because Python provides a large graphics library ensuring an aesthetic design pleasing the eyes of the user and conveying the necessary information at the same time effectively.



Figure 3.15: Flask Logo

To retrieve real time data to the web server we have taken the RESTfull approach, which is based on REST API (Representational State Transfer Application Programming Interface). REST is any interface between two systems which uses HTTP to obtain data and generate operations on those data in a variety of formats, such as XML and JSON. For our design, the web page needs to be constantly updated without reloading the page. That is, the page that generates business intelligence (various graphs) needs to be drawn realtime. To ensure that the data is received realtime without a special user interaction such as pressing a button or refreshing the page we have integrated REST API and AJAX calls which provides the needed functionality.

Another advantage of the RESTfull approach is that it enables data transfer in JSON object format. JSON is a text-based data format following JavaScript object syntax, which exists as a dictionary which is useful in transmitting data across networks. The data is retrieved to the javascript with AJAX(Asynchronous JavaScript and XML) calls. AJAX calls are important in this scenario because data needs to be sent to the server in the background after the page has been loaded and update it without reloading the page.

The received data from the server is written to a database to facilitate data storing and retrieval at future time. This is developed with MySQL database platform. Received data is written to two separate tables one to showcase heatmap information and one to provide historical density information.

All the user credentials are stored in the database in a separate table as well. To access information about his/her store the manager should have user credentials created by the admin team. After successfully logging in only can he access the needed information. Successfully logged in user would be redirected to a URL with generated information about his office.

Web version of the business intelligence supports 4 formats.

1. Density - real time store density graph showing the number of customers at the moment.
2. Tracking - real time store map being updated showing the position of each customer at the moment.
3. Heatmap - an interface is provided to enter the starting time and ending time. Heatmap is drawn showing the areas most utilized by the customers during the time period specified. The red color shows high customer density and blue shows relatively low customer density.
4. Historical - when the starting time and ending time is specified customer density over time is drawn for the specified time period.

This information is vital for an organization in various ways. The number of customers currently browsing the store would be a good indicator to understand how much sales staff should be reassigned to the cashier. Tracking maps could be beneficial to comprehend the customer behavior. Paths customers take while browsing through the store can be recognized and the store can be restructured to provide the ultimate customer experience. If a higher number of customers search for a specific product after selecting another product this could be taken as an indicator that these two products should be placed near to each other. With the heatmap information, a manager would be able to understand customer behavior at different time periods. Which places tend to

attract more customers at morning and how should the sales staff be dispersed throughout the store at that time period could be identified. The places that is ideal to showcase advertisements could be identified by looking through these maps. Also at what time advertisements should be shown or promotional activates should take place could be recognized. Historical density data could be utilized into forecasting the number of customers that the store would attract in following days or the same period in the next year. Trends could be analyzed from this data to formulate plans for future.

Chapter 4

RESULTS

4.0.1 Discussion

BIBLIOGRAPHY

- [1] Vicente, Alfredo Gardel, et al. "Embedded vision modules for tracking and counting people." *IEEE Transactions on Instrumentation and Measurement* 58.9 (2009): 3004-3011.
- [2] Dalal, Navneet, and Bill Triggs. "Histograms of oriented gradients for human detection." *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. Vol. 1. IEEE, 2005.
- [3] Negi, Kazuhiro, et al. "Deep pipelined one-chip FPGA implementation of a real-time image-based human detection algorithm." *Field-Programmable Technology (FPT), 2011 International Conference on*. IEEE, 2011.
- [4] Redmon, Joseph, et al. "You only look once: Unified, real-time object detection." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016.
- [5] Andriyenko, Anton, Konrad Schindler, and Stefan Roth. "Discrete-continuous optimization for multi-target tracking." *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012.
- [6] Tang, Nick C., et al. "Cross-camera knowledge transfer for multiview people counting." *IEEE Transactions on image processing* 24.1 (2015): 80-93.
- [7] Yang, Tao, et al. "Robust people detection and tracking in a multi-camera indoor visual surveillance system." *Multimedia and Expo, 2007 IEEE International Conference on*. IEEE, 2007.
- [8] Flask.pocoo.org. (2017). Tutorial – Flask Documentation (0.12). [online] Available at: <http://flask.pocoo.org/docs/0.12/tutorial/> [Accessed 6 Oct. 2017].
- [9] Pythonprogramming.net. (2017). Python Programming Tutorials. [online] Available at: <https://pythonprogramming.net/flask-registration-tutorial/> [Accessed 6 Oct. 2017].
- [10] Flask-restful.readthedocs.io. (2017). Flask-RESTful – Flask-RESTful 0.3.6 documentation. [online] Available at: <https://flask-restful.readthedocs.io/en/latest/> [Accessed 6 Oct. 2017].