

UNIVERSIDAD AUTÓNOMA DE MADRID
ESCUELA POLITÉCNICA SUPERIOR



**ONLINE CONTEXTUAL UPDATING
IN MULTI-CAMERA SCENARIOS**

Alejandro López Cifuentes
Director: Marcos Escudero Viñolo
Supervisor: Jesús Bescós Cano

-MASTER THESIS-

Departamento de Tecnología Electrónica y de las Comunicaciones
Escuela Politécnica Superior
Universidad Autónoma de Madrid
June 2017



PÁZMÁNY PÉTER CATHOLIC UNIVERSITY
Faculty of Information Technology and Bionics



ONLINE CONTEXTUAL UPDATING IN MULTI-CAMERA SCENARIOS

Alejandro López Cifuentes

Director: Marcos Escudero Viñolo

Supervisor: Jesús Bescós Cano



Video Processing and Understanding Lab

Departamento de Ingeniería Informática

Escuela Politécnica Superior

Universidad Autónoma de Madrid

June 2017

This work has been partially supported by Ministerio de Economía, Industria y Competitividad of the Spanish Government and Fondo Europeo para el Desarrollo Regional of the European Union under the project TEC2014-53176-R (HAVideo)



Unión Europea

Fondo Europeo
de Desarrollo Regional
"Una manera de hacer Europa"

Resumen

Este proyecto describe un método para realizar detección de personas y segmentación semántica en un escenario multicámara. La fusión de estas dos disciplinas lleva a detecciones de personas contextualmente refinadas. El entorno multicámara es usado para reproyectar detecciones de una cámara a otra y mejorar la precisión. Como ejemplo de uso, datos estadísticos de áreas semánticas concretas en la escena han sido también extraídos. Para lograr una interacción completa con el usuario se ha diseñado una Interfaz de Usuario Gráfica. Esta interfaz permitirá al usuario definir el entorno de detección de personas, así como mostrar resultados en tiempo real de ejecución.

Para poder llevar a cabo estas tareas un estudio del estado del arte ha sido realizado. Se ha analizado los diferentes detectores de personas, haciendo énfasis en aquellos que utilizan propuestas de objetos. Además, se han estudiado nuevos métodos en el entorno de la detección de personas tales como las nuevas redes neuronales. Se ha revisado el estado del arte actual sobre la extracción de información contextual, y en concreto, en el uso de la segmentación semántica. Finalmente, los escenarios con configuración multicámara han sido descritos.

Para presentar los resultados y configurar los diferentes algoritmos se ha diseñado y desarrollado una aplicación multihilo.

Un nuevo sistema ha sido propuesto para lograr los objetivos mencionados y para lograr detecciones de persona bajo diferentes condiciones de filtrado. Se han integrado en el sistema detectores como HOG, DPM, o PSP-Net y además se ha realizado segmentación semántica. Ambas fuentes de información han sido combinadas en un entorno común representado mediante un plano cenital de la escena.

Finalmente, el rendimiento del sistema ha sido probado en un data-set generado y manualmente anotado para generar gráficas de rendimiento.

Palabras Clave

Detección de personas, multi cámaras, segmentación semántica, plano cenital.

Abstract

This project describes a method to perform pedestrian detection and semantic segmentation in a multi-camera recorded scenario. The fusion of both disciplines leads to a contextually refined pedestrian detection. The multi camera system is used to reproject detections from one camera to the others to improve accuracy. As an use example, statistical data usage of specific semantic areas in the scene is also extracted. For user interaction a Graphical User Interface (GUI) is designed. The GUI allows the user to define the method setup as well as display results in execution time.

In order to carry these tasks out a study of the state of the art has been done. Pedestrian detection approaches are reviewed emphasizing in those that rely on object proposals. Also, recent trends in the task of Pedestrian Detection are analyzed. In addition, current state of the art in the extraction of contextual information and - specifically- on the sue of, semantic segmentation is studied. Finally, multi camera scenarios configurations are also described.

A multithread application has been developed and analyzed to, as mentioned before, present results and tune different algorithms.

A new system has been proposed in order to achieve the objectives and perform pedestrian detections under different filtering and fusion conditions. Detectors such as HOG, DPM or PSP-Net have been integrated and a complete semantic segmentation has been performed. Both information has been combined in a common developed frame.

Finally, system performance has been tested in a generated dataset with manually annotated ground truth.

Keywords

Pedestrian detection, multi camera, semantic segmentation, cenital plane.

Acknowledgements

*Alejandro López Cifuentes.
2017.*

Contents

Resumen	v
Abstract	vii
Acknowledgements	ix
1 Introduction	1
1.1 Motivation	1
1.2 Objectives	2
1.3 Thesis Structure	3
2 State of the Art	5
2.1 Pedestrian Detection	5
2.1.1 Classical Pedestrian Detection Approaches	7
2.1.2 Recent Trends in Pedestrian Detection	9
2.1.3 Next Steps Towards Generic PD	11
2.2 Contextual Information	11
2.2.1 Global Context	12
2.2.2 Local Context	12
2.2.3 Offline	12
2.2.4 Online	12
2.2.5 Semantic Segmentation	13
2.3 Multi-camera scenarios	15
3 Proposed System	19
3.1 Contextual Model Generation	19
3.1.1 Pyramid Scene Parsing Network	19
3.2 Pedestrian Detection	23
3.2.1 Histogram of Oriented Gradients	23
3.2.2 Deformable Part Model	24
3.2.3 Aggregated Channel Features Detector	25
3.2.4 Fast Region-Based Convolutional Network	26
3.2.5 PSP-Net	27
3.3 Fusion and Filtering	28
3.3.1 Multi-camera Scenario	28

3.3.2	Pedestrian Filtering and Constraining	34
3.4	Reference RGB and Semantic Planes	37
3.4.1	Semantic and RGB Reference Plane Generation	40
3.4.2	Semantic Fusion in the Reference Plane	41
3.5	Statistical Usage Data	43
3.5.1	Statistical Semantic Map Generation	43
3.5.2	Usage Curves and Paths	43
4	Developed Application	47
4.1	Single-thread Developer Application	49
4.2	Multi-thread Developer Application	51
4.2.1	Main Application Window	53
4.2.2	Classes Distribution	56
4.3	Multi-thread User Application	57
5	Results	61
5.1	Used Hardware	61
5.1.1	Camera Specifications	61
5.1.2	Camera User Web Interface (GUI)	63
5.2	Experiment Setup	64
5.2.1	Data-set	64
5.2.2	Ground Truth Generation	66
5.2.3	Evaluation Framework	68
5.3	Application Performance	69
5.3.1	Application Performance Results Discussion	70
5.4	Homography Calculation	70
5.4.1	Homography Results Discussion	70
5.5	Semantic Segmentation	70
5.5.1	Semantic Segmentation Results Discussion	72
5.6	Pedestrian Detection	72
5.6.1	Mono-Camera	74
5.6.2	Multi-Camera	75
5.6.3	Mono-Camera with Semantic Constraining	76
5.6.4	Multi-Camera with Semantic Constraining	77
5.6.5	Pedestrian Detection Results Discussion	78
5.7	Statistical Usage Data	81
5.7.1	Statistical Usage Data Results Discussion	82
5.8	Overall Discussion	86
6	Conclusions and Future Work	87
6.1	Conclusions	87
6.2	Future Work	88
A	Cenital Plane Design	89
B	AKAZE	93

C Parametric Homographies Between Inertial Planes	95
Bibliography	97

List of Figures

2.1	Examples images of PD data-sets	6
2.2	Generic Pedestrian Detector Example Diagram.	8
2.3	Semantic segmentation on ADE20K data-set.	13
2.4	Semantic Segmentation result examples on Cityscapes Data-set.	16
3.1	Flowchart of the proposed method	20
3.2	Overview of the proposed PSPNet	21
3.3	HOG Pedestrian Descriptor	24
3.4	Example detection obtained with the DPM person model.	25
3.5	Aggregated Channel Features architecture.	26
3.6	Fast-RCNN architecture.	27
3.7	Multi-camera configuration	28
3.8	Initial Camera Views	29
3.9	Cenital plane with camera positions.	30
3.10	Multi-camera configuration with panning set up	31
3.11	View selection process and homography calculation.	33
3.12	Cylinder estimation for camera instance projections	35
3.13	Cylinder estimation for cenital view projections.	35
3.14	Gaussian representation examples.	36
3.15	Pedestrian detection reprojection.	38
3.16	Pedestrian semantic constraining examples.	39
3.17	RGB and semantic frames projected.	40
3.18	RGB Reference Planes	41
3.19	Semantic Median Average	42
3.20	Common semantic areas between pair of cameras.	44
3.21	Common semantic areas for all the cameras	45
3.22	Statistical semantic map.	45
4.1	QT Main Window Designer	48
4.2	Flow-chart legend	49
4.3	Flow-chart diagram for the single-thread application.	50
4.4	Flow-chart diagram for the multi-thread application.	52
4.5	Main application window.	53
4.6	Application menu bar	54
4.7	Options Menu.	55

4.8	Information Display	55
4.9	Results Display Area	56
4.10	Hierarchical representation of the code.	58
4.11	User version application main window.	59
5.1	Camera Sony SNC-RZ50P Pan/Tilt Range diagram	62
5.2	Visualization and control menu	63
5.3	Preset position setting menu	63
5.4	Tour setting menu	64
5.5	Data-set example frames	65
5.6	Via Annotation Tool software main window	66
5.7	Annotated ground truth frames	67
5.8	Intersection over Union or Jaccard Index	69
5.9	Projected views with points	71
5.10	RGB frames and PSP-Net Semantic segmentation results	73
5.11	Recall / Precision graphs mono camera	74
5.12	Recall / Precision graphs multi camera	75
5.13	Recall / Precision graphs mono camera with semantic constraint . . .	76
5.14	Recall / Precision graphs multi camera with semantic constraint . . .	77
5.15	Differences in number of detections for PD approaches	80
5.16	Misclassification of people in semantic segmentation. Camera 3.	80
5.17	Pedestrian detection error leads to a reprojection displacement.	81
5.18	Pedestrian reprojection error due to height.	82
5.19	Floor usage graph	83
5.20	Doors usage graph	83
5.21	Pedestrians paths usage along the Hall. Region of 40 pixels	84
5.22	Pedestrians paths usage along the Hall. Region of 70 pixels	85
A.1	First cenital plane approach	89
A.2	Second cenital plane approach with real measures	90
A.3	Second cenital plane approach with camera positions	91
C.1	Extending homography for parallel planes.	96
C.2	Set of inertial planes	96

List of Tables

2.1	Pedestrian Detection Performance.	9
2.2	Cityscapes Data-set Class Definitions	14
2.3	Cityscapes Data-set Challenge Results (Intersection over Union metric)	14
3.1	Final classes list from ADE20K to use in PSP-Net	23
5.1	Camera Sony SNC-RZ50P Specifications	62
5.2	Comparison between the use of single thread or multi threads	69
5.3	Comparison between the use of single thread or multi threads	69
5.4	F-Score comparison for pedestrian detector approaches.	78

Chapter 1

Introduction

1.1 Motivation

Nowadays, we live surrounded by electronic devices which claimed objective is to ensure the safety and security of the global population and to ease our lives on everyday tasks. These range from biometric systems [1] to all kind of different electrical sensors, including video surveillance cameras. These cameras are the ones that are of real interest when developing Computer Vision algorithms in the scope of video surveillance [2].

The combination of these veins of research may lead to the automation of high-level human semantic tasks such as people detection [3], object detection and recognition [4, 5, 6] and extraction of contextual information [7]. The automation of these processes permits end-users build on these information sources to define the latest stages of video surveillance systems. These are usually the critical ones, e.g. alarm raising when some predefined event occurs.

Usually, video surveillance systems are focused either on the analysis of a single-camera point of view -which will lead to a simple scenario in which the potential actions/events to detect will be observed from a single point in the scene- or, on the analysis of a multi-camera setup. This last configuration may provide multiple benefits when analyzing big spaces as it provides to the user different views of the scene, disambiguating occluded areas the mono-camera views.

Among Computer Vision applications running on a multi-camera scenario, a pivotal field of research is the analysis of public spaces. These are often crowd populated scenarios which analysis requires the combination of the data obtained by all recording cameras. It is of real interest to analyze people behavior patterns [8, 9, 10] and temporal usage of a given area in large-scale scenarios such as shopping malls, uni-

versities and, generally, public-use buildings. Analysis ranges from the extraction of statistical measures of behavior to the detection of anomalous unexpected events [11]. This results may come from a combination of complementary algorithms such as contextual and semantic area classification, people detection and crowd behavior analysis.

1.2 Objectives

The main objective of this thesis is to extract contextual descriptions from a large-scale populated multi camera scenario. A potential application of this task is illustrated by the extraction of temporal statistical usage data from relevant areas in the scene. The whole solution needs to be controlled through the use of a Graphical User Interface application.

To fulfill this objective, this work embraces two different blocks of objectives that complement each other. The first one targets the design of a graphical user interface (GUI), The second block deals with algorithm and research-related objectives.

Graphical User Interface

The GUI should be able to visualize and dynamically arrange –under a user-friendly environment– statistics from different areas of interest in a public space.

Algorithm

The algorithm related objectives are:

1. To integrate a semantic segmentation algorithm to perform contextual element in video sequences. The objective is to detect, classify and determine the spatial extend on each frame of the video of relevant elements such as doors, chairs, corridors and floor areas. We aim to:
 - (a) Combine semantic information coming from different cameras.
 - (b) Identify the usage rate of some important elements of the scene measured by number of people per time interval.
2. To globally integrate state-of-the-art pedestrian detection algorithms results per view. To this aim, we need to:
 - (a) Create a pedestrian detector fusion mechanism to take advantage of the multi-camera scenario.

- (b) Increase pedestrian detection algorithms performance by the use of semantic constraining information to suppress false detections.

1.3 Thesis Structure

The master thesis is divided into the following chapters:

- Chapter 1. Introduction.
- Chapter 2. State of the Art.
- Chapter 4. Developed Application.
- Chapter 3. Proposed System.
- Chapter 5. Results.
- Chapter 6. Conclusions and Future Work.
- Appendices
- Bibliography.

Chapter 2

State of the Art

As explained in Chapter 1, the analysis of a public crowded space embraces many different algorithms from Computer Vision disciplines.

This Chapter aims to study the State Of the Art in pedestrian detection approaches. In addition, it also covers the topic of contextual information and specifically, the algorithms in the field of semantic segmentation. The advantages and disadvantages of analysis in multi camera scenarios are also discussed.

2.1 Pedestrian Detection

Pedestrian detection (PD) has been a hot research topic in Computer Vision during the past few years due to its impact in several Computer Vision applications. Its main objective is to identify a potential object as a person by automatically detecting its position and relative size in the scene.

Nowadays, it can be consider a partially-solved problem. Although there are excellent PD in the literature, there is no algorithm able to effectively perform PD on a generic scenario. This is the main reason why PD is still one of the most researched areas in computer vision.

The complexity related to PD lies on the large amount of available data-sets with different video and people characteristics including challenges such as: people occlusions, poses and scales and scenes captured under extreme illumination conditions. Caltech [12] –recorded on a vehicle in an urban environment–, ETHZ [13] –recorded from a chariot which moves through pedestrian paths–, TUD [14] –static camera in a crossing campus scene– or INRIA [15] –collects precise people images both static and moving– are some of the PD data-sets that have been proposed through the years to train PD. Figure 2.1 gathers some examples of the images in these data-sets.

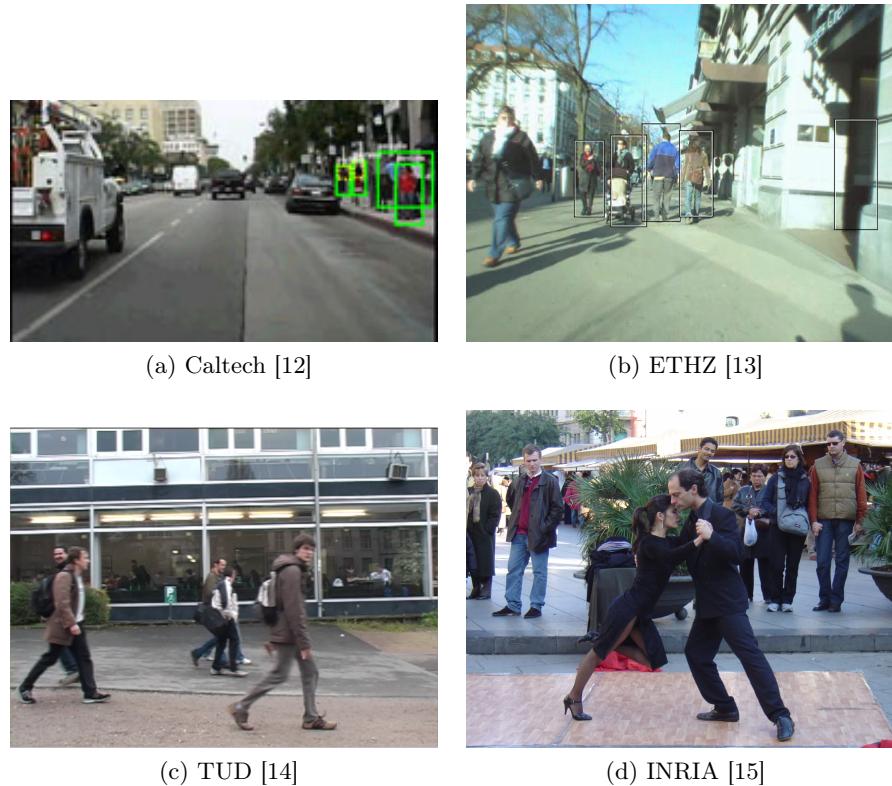


Figure 2.1: Examples images of PD data-sets.

2.1.1 Classical Pedestrian Detection Approaches

Several PD in this State of the Art are arranged under three different topics. First, different pedestrian descriptor schemes are explained. Second, person detection approaches are studied. Finally, various approaches to define person model are analyzed.

Person Description

An organization of existing PD approaches based on the descriptor may start with the Histograms of Oriented Gradients (HOG), which in combination with linear Support Vectors Machine (SVM) have been mainly used to describe pedestrian shape [15]. Differentially, discriminative Part Models (DPM) such as [16] propose to divide the human body into different parts (head, trunk, legs...) and search for their combination on the image to extract PD.

Others PD approaches are based on the use of the Aggregate Channel Features [17]. Algorithms based on ACF rely on a combination of different channels such as normalized gradient magnitude, histogram of oriented gradients (6 channels), and LUV color channels to achieve the final detection.

Person Detection

In [18] object detection is defined as the extraction of potential object candidates to be a person from a scene. Mainly object detections algorithms are:

- Sliding window –also known as– exhaustive search: uses an efficient classifier to test every possible image window. Parameters such as window size, or overlapping between them, are common tuning values that increase or decrease the performance of the detector. These methods usually need from 10^4 to 10^5 windows per image to perform decently. This number grows exponentially for multi-scale detection. If the complexity of the core classifier is increased in every window testing, the computational time will end up being not affordable.
- Segmentation: Uses segmentation as a preliminary step for PD. The use of algorithms such as background subtraction lead to the segregation of the image. Alternatively, color segmentation based on color skin detection can be used to restrict people search. By all means, segmentation severely reduces person candidates reducing the computational time.
- Segmentation + Exhaustive search: Alternatively, a combination of previous techniques. In this case the previous step of segmentation does not lead to

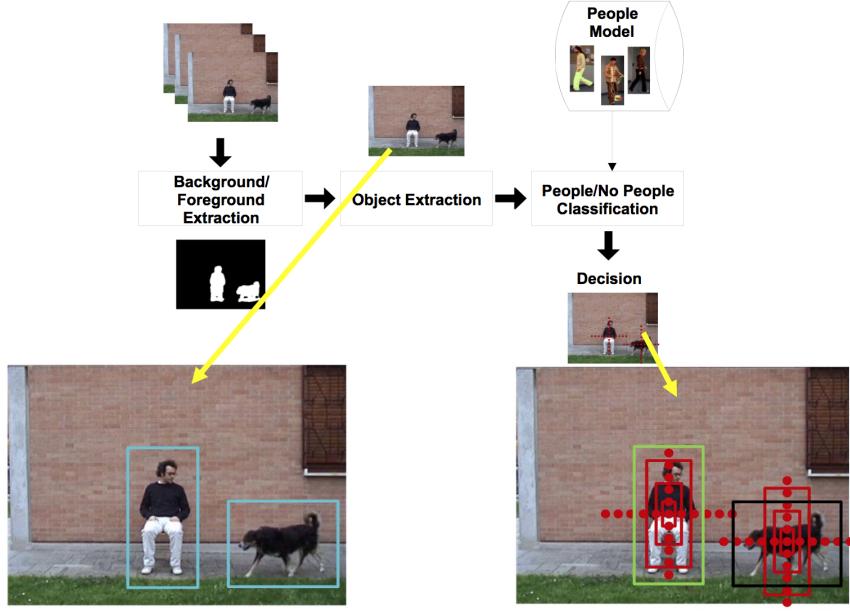


Figure 2.2: Generic Pedestrian Detector Example Diagram. Image extracted from [18]

final candidates but to a delimited small area that could contain some candidates objects. After the segmentation process, a sliding window technique is performed over the reduced scene area. In this case, improvements from both approaches are exploited as the computational cost of the exhaustive search (which is its main drawback) is reduced by the use of segmentation. Figure 2.2 depicts a flowchart for a generic PD approach which relies on the combination of segmentation and exhaustive search.

Person Model

The model of a person can be considered as the set of characteristics used to discriminate between people and any other object in the scene. In [18] three different types of person model are assumed:

- Based on appearance: Most of the available detectors in the State of the Art use appearance information to define the person model. In this group one can differentiate two approaches to describe the shape of a person.
 - ❖ Holistic: People are defined as a unique and indivisible region or shape.
 - ❖ Part-based: Rely on more complex characteristics where a person is defined as a combination of multiple shapes, regions or parts of its body.

Video	HOG	ISM	Fusion	Edge	DTDP	ACF	Faster-RCNN
1	89.3	71.4	34.9	84.9	96.7	99.3	99.7
2	69.2	82.9	92.5	90.2	77.1	77.1	98.2
3	55.6	75.7	64.3	71.7	68.9	68.9	82.9
4	10.1	1.0	0.5	5.4	33.9	33.9	37.5
Average AUC	56	57.5	48	63	69.1	69.8	79.5

Table 2.1: Pedestrian Detection Performance. Adapted from [22] (selected approaches). Metric for this evaluation is the average AUC.

Examples of appearance-based detectors are those that use silhouettes to classify people, either from an holistic or a part-based basis, or color distribution.

- Based on motion: Human appearance is likely to change due to environmental factors such as light conditions, clothes or camera settings. In addition, people variability in terms of height, weigh and poses make appearance likely to vary. For these reasons, some approaches try to get rid of these factors and detect pedestrians using only its motion information. For instance in [19], detections are based on a periodic motion analysis.
- Based on appearance + motion: Algorithms such as [20, 21] merge both appearance and motion information. Most of these algorithms combine people detection and tracking, targeting to improve people tracking rather than PD.

In [22] a comparison of the performance of PD on different data-sets is made. This comparison is partially included in Table 2.1. See [VPU Website](#)¹ for the complete table.

2.1.2 Recent Trends in Pedestrian Detection

During this section recent PD trends in terms of person detection and description are presented.

Person Detection

PD based on HOG, DPM and ACF generally rely on “sliding window” detectors , however, as mentioned before one of the main drawbacks of this approach is the high computational time needed to achieve good performance.

¹<http://www-vpu.eps.uam.es/DS/PDbm/results.html>

One of the most successful solutions to overcome this time consumption problem without losing detection quality is the use of object proposals [6].

Object proposals approaches perform a complete search over an image to detect potential object candidates. These candidates are detected as image areas with visual properties that distinguish them from the scene background.

In general, object proposal approaches reduce pedestrian candidates with respect to sliding-window like algorithms and generally outperform segmentation based methods. This advantages lead to a higher object recall and more efficient detection processes. Successful examples of using proposals to improve and speed-up detection include Faster R-CNN [5].

In [6] three set of proposal methods are analyzed:

- **Grouping proposal methods** attempt to generate multiple, and so, overlapping segments that are likely to correspond to objects. Here one can distinguish between three types of methods according to how they generate proposals. Methods can generate proposals by groping super-pixels (SP), solving multiple graph cut (GC) problems or finally, using edge contours (EC). Among those that use SP we can find Selective Search [23], Randomized Prim's [24], Rantalaikia [25] or Chang [26]. Those that use GC are CPMC [27], Endres [28] or Rigor [29]. Finally Geodesic [30] and MCG [31] use EC to obtain proposals.
- **Window scoring proposal** methods are an alternative approach to score each candidate window according to the probability to contain an object. Usually this methods tend to be faster and, in addition, they typically extract only bounding boxes. One can find among other approaches Objectness [32], Rahtu [33], Bing [34], EdgeBoxes [35], Feng [36], Zhang [37], RandomizedSeeds [38].
- **Alternative proposal methods.** Apart from the main groups a set of alternate approaches such as ShapeSharing [39] or Multi-box [40] are also used to extract object proposals.

Person Description

In recent years new schemes of PD have been proposed. Detectors based on deep Convolutional Neural Networks (CNNs) have notably improved the accuracy of all the previous analyzed algorithms.

Examples such as ImageNet [41] for image classification, CompACT [42], Fast-RCN [43] or Faster R-CNN [5] for object detection expose that deep convolutional networks usually improve the performance of aforementioned approaches. This fact is clearly presented Table 2.1 where Faster R-CNN outperforms every other approach.

METER MAS

2.1.3 Next Steps Towards Generic PD

PD is constantly in development and so, some future work lines can be set. In [3] some research directions are proposed that could be of interest in the scope of this work.

1. Use of context information. Starting from the hypothesis that a person is standing on the floor, the ground plane assumption can reduce errors if the detection for both the person and the floor are accurate. This could be achieved by extracting useful contextual information from the scene. In [18] contextual information is added to PD to increase performance. This is one of the main objectives of the work.
2. Occlusion treatment. Usually pedestrians, due to other scene elements such as columns or even other pedestrians appear occluded. When this happens PD performance is substantially degraded under even mild occlusions. Improvements in this area could increase the overall performance of PD in generic scenarios.

2.2 Contextual Information

One can describe contextual information as the set of additional circumstances or facts that can be extracted from a scene besides the target of analysis. Generally, this set of circumstances is a source of information that is not extracted by machine applications but constitutes key evidences for humans, which acquired this knowledge during their life. By just taking a look to an outdoor image a human can derive where the sky will be, what the weather conditions are or which time of the day is. Also, by knowing the place where a video was recorded, one would imagine which objects are more or less probable to be in the scene.

Dealing with computer vision disciplines, contextual information sources also include camera information (such as position, configuration, distance to an object and camera motion), the set of objects that one could detect in the scene and the number of available cameras.

One can divide contextual information into two different levels: global and local. Besides, we can also divide contextual information into two different categories: offline, and online. Finally, we focus on a set of specific methods to extract context –semantic segmentation–.

2.2.1 Global Context

Global context considers descriptions from an image as a whole. For instance, if the context of a scene is known –kitchen–, we can use this information to search for typical objects in this context –e.g. a stove–.

This kind of approaches are focused on psychology studies that suggests that human perceptual processes work following a hierarchically organized process [44, 45]. Our perception system goes from a global structure towards a more detailed analysis in a top-down scene interpretation.

Global context approaches aim to define a scene as an extra source of global information. The structure of a determinate scene image can be estimated by means of global image features as in [46].

2.2.2 Local Context

On the contrary, local context refers to contextual information related to a specific object, e.g. a kitchen table may help to predict the presence of a spoon.

The impact area of an object –in contextual terms– is defined as a set of neighboring objects, patches or even pixels. Algorithms dealing with the extraction of local contextual information aim to correctly define the area that surrounds an object to precisely detect other object instances.

In [47] the inclusion of local contextual regions such as facial bounding contour are used to improve face detection performance.

2.2.3 Offline

Offline contextual information is defined as the set of circumstances that are computed before starting any kind of analysis procedure. This information leads to external image information that may be used to constraint analysis algorithms. Approaches such as [18] use previously introduced contextual information to improve part-based PD over a scene.

2.2.4 Online

Online information, on the other hand, is to be extracted with the analysis. Online extraction entails a degradation of an algorithm efficiency, albeit allows to dynamically update the context.

Examples of online (and local) contextual information extractors are the algorithms in the semantic segmentation branch. Next section (2.2.5) discuss some of the approaches in this vein.

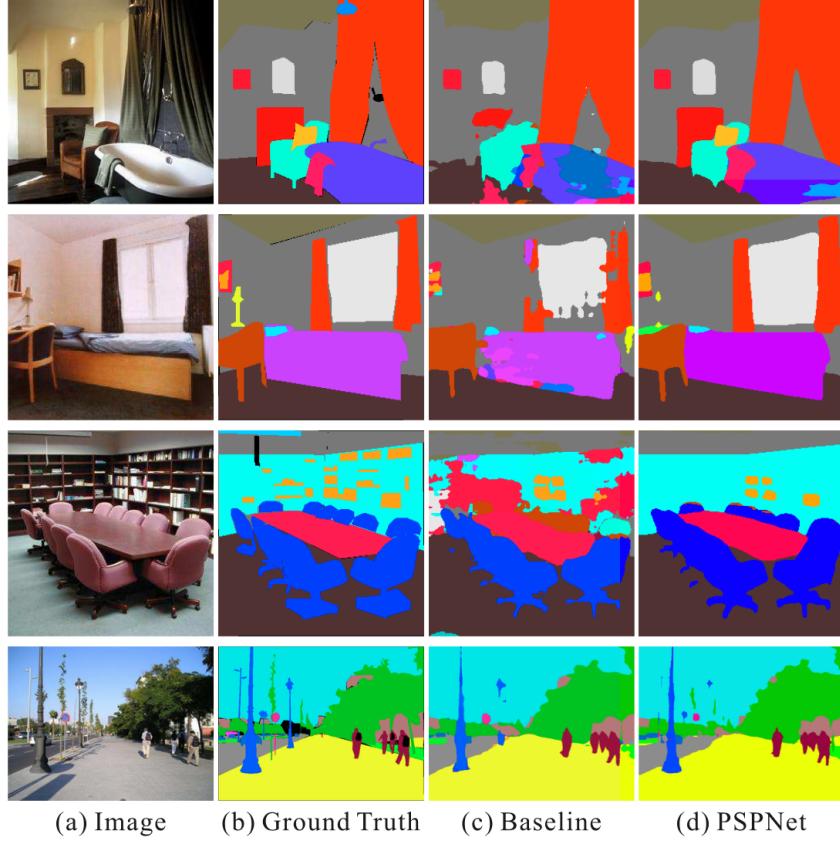


Figure 2.3: Semantic segmentation on ADE20K data-set [50] by the algorithm described in [7].

2.2.5 Semantic Segmentation

Semantic information is defined as the set of high-level elements from contextual information. This can often lead to characteristics such as global image color [48]. It can also describe an image by the position and status of relevant objects [49]. Or it can also be used to define specific image areas in a scene such as walls, corridors or walking paths [7].

Semantic segmentation targets to assign each image pixel a high-level label. If accurately performed, it provides fully semantic understanding –which in terms of Computer Vision–, means that the location of an object within an image is known. A potential result of a semantic segmentation method is depicted in Figure 2.3. In the Figure four different scenes are analyzed and divided into non-overlapping semantic areas.

Semantic segmentation is a recent trend and nowadays, and so, it remains a significant challenge for the computer vision community. Due to its short-life term there

Category	Classes
Flat	road · sidewalk · parking · rail track
Human	person · rider
Vehicle	car · truck · bus · on rails · motorcycle · bicycle · caravan · trailer
Construction	building · wall · fence · guard rail · bridge · tunnel
Object	pole · pole group · traffic sign · traffic light
Nature	vegetation · terrain
Sky	sky
Void	ground · dynamic · static

Table 2.2: Cityscapes Data-set Class Definitions

Algorithm Name	IoU Category	IoU Class	Available Code
motovis (Anonymous)	91.5	81.3	No
PSPNet [7]	91.2	81.2	Yes
ResNet-38 [52]	91.0	80.6	Yes
NetWarp (Anonymous)	91.0	80.5	No
TuSimple_Coarse [53]	90.7	80.5	Yes

Table 2.3: Cityscapes Data-set Challenge Results (Intersection over Union metric)

is not yet a complete survey available in which algorithms are deeply analyzed and compared. However, there are a set of benchmark where developers can upload their obtained results with a given dataset and so, algorithms performance can be compared.

An example of a popular benchmark is the Cityscapes Data-set [51]. It is a large-scale data-set that contains stereo video sequences recorded in street scenes from among 50 different cities around the world. This data-set presents categories annotations over pixels in more than 5000 frames. The set of categories is presented in Table 2.2, whereas a subset of the compared methods in the [Cityscapes Website](#) are depict in Table 2.3. Only those algorithms that have more than 80% on IoU class metric have been included in the Table.

PSPNet [7], ResNet-38 [52] and TuSimple_Coarse [53] are all based on convolutional networks. This fact reveals that deep convolutional neural networks have led

to significant improvement over previous semantic segmentation systems since the presentation of AlexNet [41] in 2012.

However, even when using CNNs, the main difficulty of scene parsing is related to the type of scene and to label variety.

TuSimple_Coarse [53] propose a combination between dense upsampling convolution (DUC) to generate pixel-level prediction and a hybrid dilated convolution (HDC) framework.

ResNet-38 [52] on the other hand, propose not only to not increase CNNs depth, but rather to ensemble many relatively shallow networks to increase performance. Their approach also improves usability reducing memory use and sometimes even training time.

Finally, PSPNet [7] deals with this problems assigning relationships between different labels, i.e. an airplane is likely to be in runway or flying in the sky while not over a road or in the water. This relationships reduce slightly the complexity of having large amounts of labels to predict and improve the general performance of the algorithm.

In Figure 2.4 some visual examples of how this algorithms perform on Cityscapes Data-set frames are displayed.

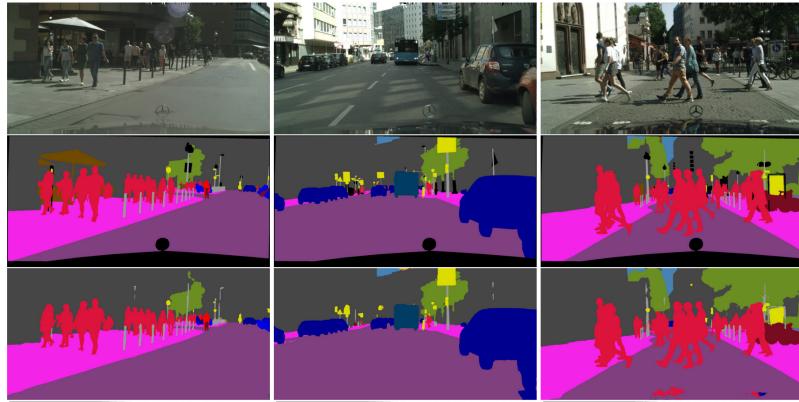
2.3 Multi-camera scenarios

The use of multiple cameras is a common setup when dealing with video surveillance problems. One can define a multi-camera scenario as a space that has more than one video camera recording. Ideally, the recordings for the different cameras are temporally aligned –synchronized–. Having N camera instances allows to observe the same event or object of interest from different points of view. This leads to a set of advantages in the scope of our work when dealing with PD and semantic classification and also, to some unavoidable disadvantages.

- Advantages

As discussed in Section 2.1 one of the research paths towards PD is generic occlusion handling. In this case, the use of a multi-camera scenario with relating camera views could help. This could be achieved by reprojecting detections from one camera to another whose miss rates are high as in [54].

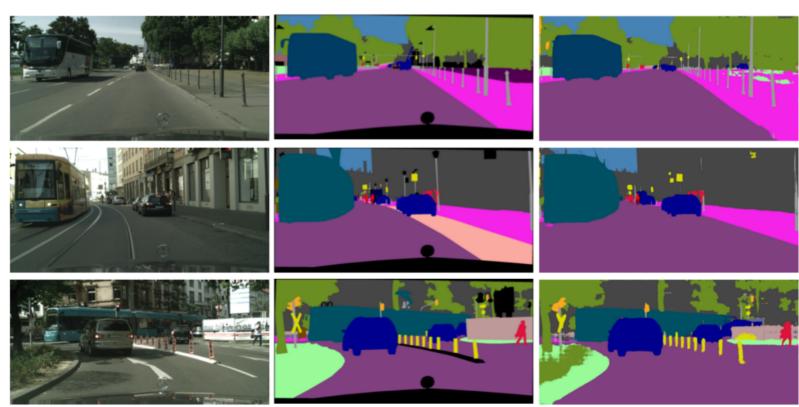
When dealing with semantic segmentation, the inclusion of a multi-camera system may arise some benefits. A single-camera system could lead to misclassification of labels in the image. In a multi-camera system one camera instance



(a) ResNet-38 Algorithm. From top to bottom: Input images, ground truth and results.



(b) PSPNet Algorithm. From left to right: Input images, ground truth and results.



(c) TuSimple_Coarse Algorithm. From left to right: Input images, ground truth and results.

Figure 2.4: Semantic Segmentation result examples on Cityscapes Data-set.

may help to refine the classes in another one provided that, evidentially, the views of both cameras partially overlap [55].

- Disadvantages

The main disadvantage when dealing with multi-camera systems is the exponential grow of computational time. Algorithms should be performed N times. This issue could be solved by the use of parallel coding to process cameras views simultaneously.

Besides, the use of multiple cameras entail two additional tasks: the temporal alignment of the views –synchronization– and the spatial arrangement of the different views –homographies–.

Chapter 3

Proposed System

During this Chapter our proposed system is analyzed. We start from the contextual model generation and pedestrian detectors. Then, the fusion of all the obtained elements in a multi camera system with semantic constraints is described. Finally, a case of example which generates statistical usage data from semantic areas is proposed. Figure 3.1 depicts the flowchart of the modular proposed method.

3.1 Contextual Model Generation

One of the main objectives in the scope of this work is to perform semantic segmentation, which –as explained in Chapter 2–targets to divide one frame into different semantic areas. The relative position in the scene for elements such as doors, walls, paths, and columns is required to achieve further objectives such as multi camera pedestrian constraint and statistical data extraction.

For this complex task the algorithm PSP-Net [7] presented in Chapter 2 is used. We choose PSP-Net because at the moment of this work was the one with available code achieving the best results (see Table 2.3). The goal of this algorithm is to assign each pixel in the image a category label.

3.1.1 Pyramid Scene Parsing Network

It uses a deep Convolutional Neural Network (CNN) called Pyramid Scene Parsing Network (PSPNet). This network is designed to improve performance for open-vocabulary object identification in complex scene parsing. The structure of the network is represented in Figure 3.2.

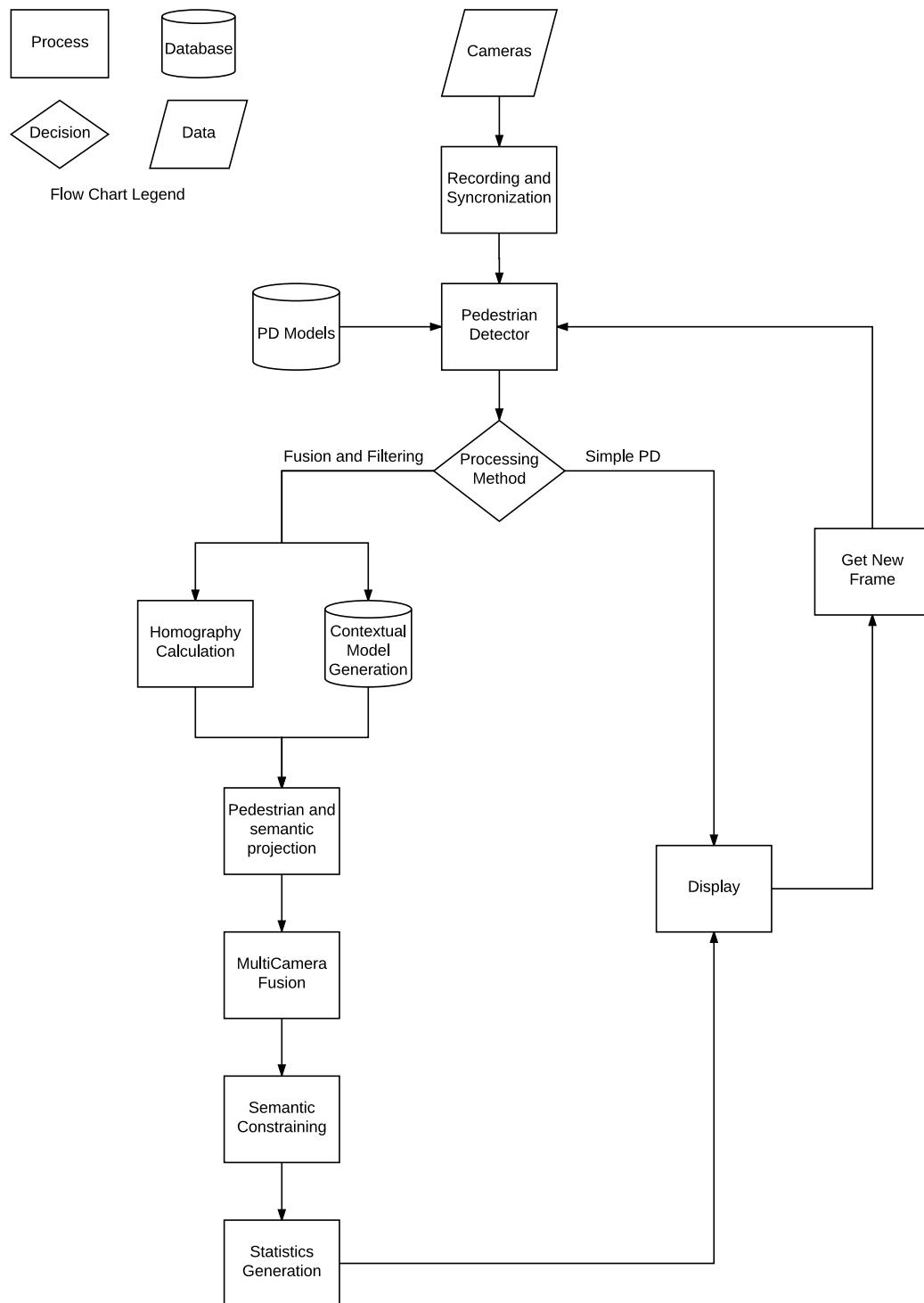


Figure 3.1: Flowchart of the proposed method

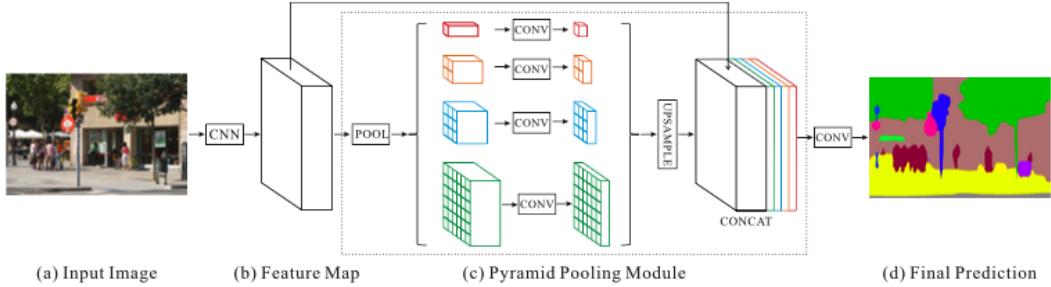


Figure 3.2: Overview of the proposed PSPNet

Algorithm Stages

The algorithm is divided into different stages:

1. Image in 3.2.(a) is processed it through a pre trained CNN called ResNet [56]. The objective is to get the full feature map 3.2.(b) of the last convolutional layer. The final feature map in this step is 1/8 of the size of the input image.
2. Apply the main contribution of [7], called the Pyramid Pooling Module 3.2.(c). The main objective of the module is to collect a few levels of information, much more representative than global pooling. It separates the feature map into different sub-regions and forms pooled representations for different locations. Here a set of pooling, convolutional and upsampling layers are applied to harvest different sub-region representation in N different scales.
3. Concatenation layers are used to form the final feature representation by fusing the feature map extracted in 3.2.(b) and the Pyramid Pooling Module output. This final feature carries both local and global context information.
4. The representation is fed into a convolutional layer which gets the final per-pixel prediction 3.2.(d).

Semantic Segmentation Particularization

PSP-Net¹ comes with a set of three different pre trained Caffe² models for three different datasets. The main difference between the models for our scope is the environment in which the network has been trained.

¹<https://github.com/hszhao/PSPNet>

²<http://caffe.berkeleyvision.org/>

- ADE20K: This dataset is the most challenging as it has up to 150 different labels in a wide range of scenes. The scenes go from interior room places to outdoor scenarios.
- VOC2012: It contains 20 object categories and one background class from diverse indoor and outdoor scenes.
- CityScapes: The last dataset defines 19 categories containing both stuff and objects. All the available sequences have been recorded from a driving car while driving in the street.

Model Particularization

As one can observe the three different models represent different object categories in different real spaces. In our case we select the model based on two main reasons:

1. The model should have been trained with indoor scenes. This leads to discard those models that represent only outdoors scenes as we would like our approach to be used in an interior scenario. This will exclude CityScapes dataset from our options as all the classes and sequences used for training are from outdoor scenes.
2. From the trained indoor models we have to choose between those whose categories best fit in our work. In this case VOC2012 dataset uses classes such as boat, airplane or table which are not interesting for our segmentation problem and it does not have classes such as door or wall which are really important for us.

Selected Model

Considering these two reasons, we have selected the model ADE20K because as said, it has elements such as walls, floor, person and column in its model.

However, we consider that most of the 150 label categories will be unused in our procedure, so, the number of classes from the model has been reduced to the 21 classes of our interest. Position and scores for those objects are the only ones obtained. This class limitation leads above all in a considerably hard drive space saving. In Table 3.1 the final 21 selected classes are exposed.

wall	building	floor	ceiling	road
window pane	person	door	table	chair
seat	desk	lamp	column	counter
path	stairs	screen door	stairway	toilet
poster	bag			

Table 3.1: Final classes list from ADE20K to use in PSP-Net

3.2 Pedestrian Detection

Along this section pedestrian State of the Art detectors that have been integrated in the proposed system are presented. Some of them have been chosen due to their efficiency, whereas some have been chosen due to its contrasted good performance (see Chapter 2).

Besides, algorithm source code of all the chosen approaches is available.

However, one of the main reasons of choosing the following five algorithms among others in the SoA is the possibility of either having the original source code or the practical implementation within an external OpenCV library.

3.2.1 Histogram of Oriented Gradients

Histogram of Oriented Gradients, i.e. HOG, is one of the main used detectors along pedestrian detection field. This fact is due to it's extremely simplicity in terms of the descriptor complexity.

Person Descriptor

Pedestrians are described as set of HOG. This means that its shape and appearance can be described by a set of gradients and intensities organized as orientation histograms. These histograms describe intensity distributions from local gradients or border directions. This descriptor can be observed graphically in Figure 3.3.

Model Generation

Once HOG have been used to describe the person shape, Support Vector Machines are used to train a person model and to classify potential candidates as people. SVM are a data classification method formed by a set of supervised training. The aim of

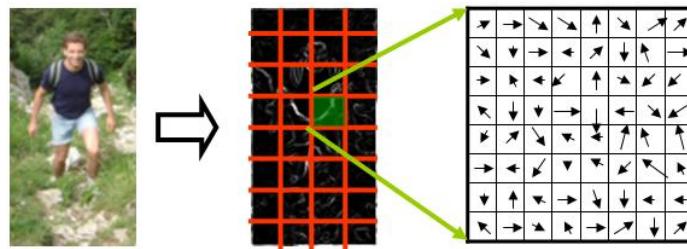


Figure 3.3: HOG Pedestrian Descriptor

this kind of approaches is to produce a model which is able to predict classification labels on a test set based only on the descriptors of the set and the model.

The idea behind SVM is that training vectors are mapped on to a bigger dimensional space in which the data separation, by means of one or many hyperplane, is much easier to divide than in the original dimensional space.

The combination between HOG and SVM leads to a fast detector that depending on the situation will perform decently, although it has some main drawbacks as its lack of occlusion treatment which does not make this algorithm usable when working with crowded spaces.

The main implementation of Histogram of Oriented Gradient Pedestrian Detector is in OpenCV library for C++.

3.2.2 Deformable Part Model

As was mentioned in the previous paragraph one of the main drawbacks when working with the simple HOG pedestrian detector is that it describes the person model as a hole which leads, inevitably, to the mentioned occlusion drawbacks.

Person Descriptor

Deformable Part Model tries to solve this problem, among others, by defining the model as first, a global coarse template, secondly, several higher resolution part templates and finally a spatial model for the location of each part. This description is the one that can be observed in Figure 3.4.

Model Generation

Both global and part templates are modeled with histogram of gradient features and the model is built by using an improving over SVM called latent SVM. In addition, scores for every detection are obtained by applying a root filter on the window plus

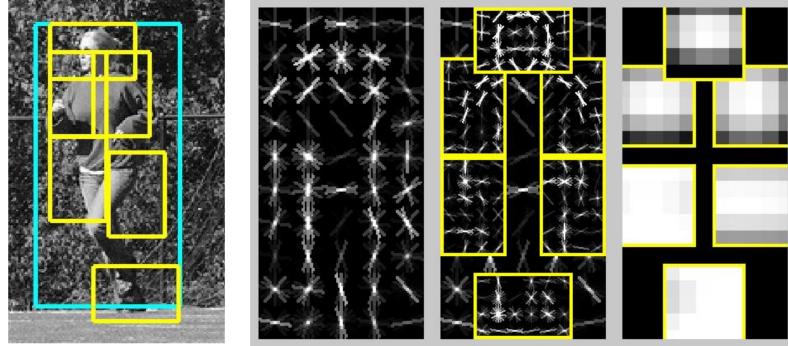


Figure 3.4: Example detection obtained with the DPM person model. From left to right: Original image + detections, global coarse model, part templates and spatial model for each part. Image extracted from [4].

the sum over parts of the maximum over placements of the part filter score on the resulting sub window minus the deformation cost.

Advantages over alternative PD Approaches

The use of this detector leads to a set of advantages than when working with others simpler approaches. In terms of pedestrian occlusion treatment we are able to detect those people that have been occluded by something in the scene just by detecting some visible part. This outperforms other detectors while working with crowded scenes in which holistic methods have problems that lead for instance to groups of people being detected as a unique detection while DPM is ideally able to separate them into different person instances.

However, its scanning window approach as well as the part based model lead to some computational cost that will increase the time needed to obtain detections.

The main implementation of Deformable Part Model Pedestrian Detector can be found in OpenCV library for C++.

3.2.3 Aggregated Channel Features Detector

The basic idea behind ACF detector is to increase other approaches performance by the use of many different channels to describe an input image I .

Algorithm Stages

In Figure 3.5 one can observe the working path of the mentioned detector.

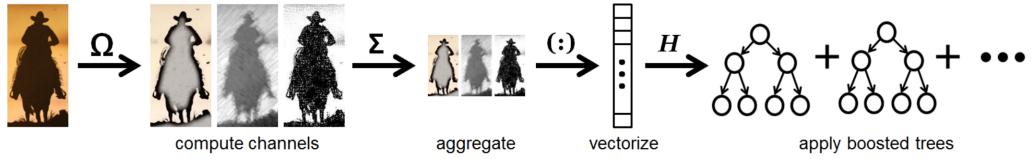


Figure 3.5: Aggregated Channel Features architecture. [17]

1. Given an initial image I , several channels are computed. These channels are named $C = \Omega(I)$. In [17] a set of 10 channels are used to achieve state-of-the-art performance in pedestrian detection:
 - (a) Normalized gradient magnitude (1 channel).
 - (b) Histogram of oriented gradients (6 channels).
 - (c) LUV color channels (3 channels).
2. After the computation, every block of pixels in C is summed and the resulting lower resolution channels are smoothed.
3. After a vectorizing process features are single pixel lookups in the aggregated channels. Boosting trees are then used to learn these features (pixels) in order to distinguish people from the background.

3.2.4 Fast Region-Based Convolutional Network

Fast Region-Based Convolutional Network (Fast R-CNN) is used to perform offline pedestrian detection in our proposed system, however, as said in Chapter 2 it is not only a pedestrian detector but an algorithm for object detection.

In 2.1.2 we mentioned that Fast-RCNN must use objects proposals for its usage. In our developed system MCG [31] grouping method is used.

Algorithm Improvements over R-CNN

Fast-RCNN method sets its contributions in a several number of innovations to improve training and testing speed over its fundamental base R-CNN:

1. Higher detection quality (mAP).
2. Training is single-stage, using a multi-task loss.
3. Training can update all network layers.
4. No disk storage is required for feature caching.

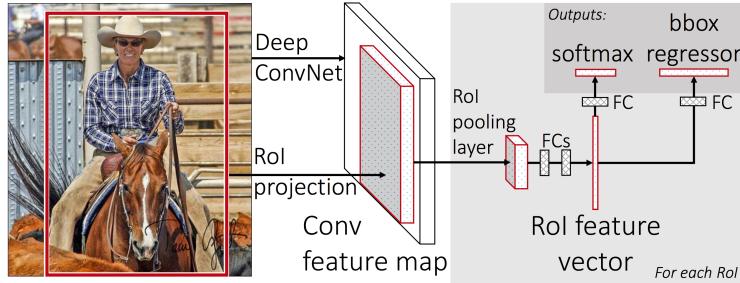


Figure 3.6: Fast-RCNN architecture. [43]

Algorithm Architecture

In Figure 3.6 one can observe the general Fast-RCNN architecture. The input for the Fast-RCNN network are the entire desired image that one wants to process and the set of object proposals. Algorithm stages are:

1. The first step in the network is to process the whole image with several convolutional and max pooling layers to produce a convolutional feature map.
2. Then, for every object proposal present in the image, a region of interest (ROI) pooling layer extracts a fixed-length feature vector from the feature map.
3. Each feature vector is after, introduced into a sequence of fully connected (FC) layers that finally diverge into to output layers.

The first one produces probabilities estimates over K object classes.

The second one outputs four real numbers for each of the K object classes. This set of 4 values encodes the final bounding box positions for one of the objects from the K classes.

In this case, both Fast-RCNN detector and MCG object proposal extract are implemented within external Matlab libraries and so, should be computed offline and introduced externally to the application.

3.2.5 PSP-Net

As we explained in Section 3.1 PSP-Net has been trained to detect people as a class. In our proposed system we also use this approach by encapsulating connected components and pixels labeled as pedestrians.

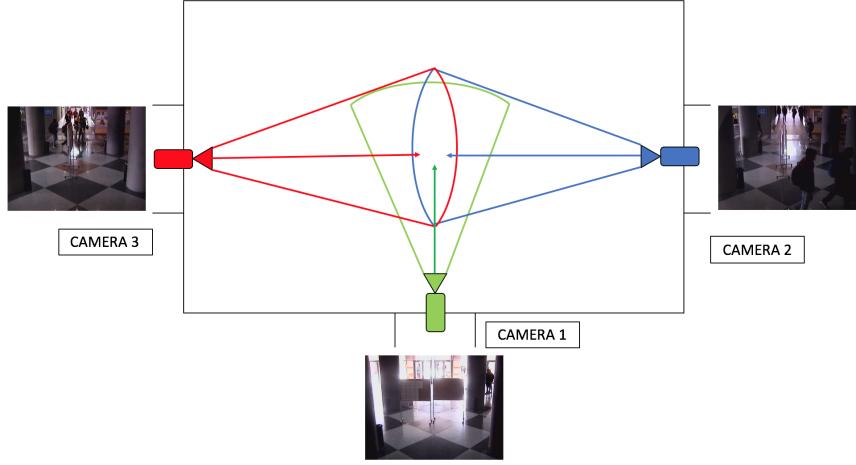


Figure 3.7: Multi-camera configuration

3.3 Fusion and Filtering

Once the semantic model and pedestrian detections are obtained the next step is to combine the information from different camera sources and project detections onto a common plane. We call this process fusion and filtering.

3.3.1 Multi-camera Scenario

Our system runs on a multi-camera scenario, hence, the scene can be observed from different point of views. In Figure 3.7 we can observe how the scene is configured with 3 static video-surveillance cameras. Two of them are placed at the sides of the scene, while one is at the bottom part. The starting views from the three static cameras are displayed in Figure 3.8.

Analyzing a multi-camera scenario has some advantages. One can see in Figure 3.8 we can observe the same scene area from three different points of views, which means that for instance, detections from one camera can be used to detect in other camera or even refine the available detections.

3.3.1.1 Cenital Plane Homography

Homography calculation is a pivotal task in our work and the main base for the fusion of all the different information coming from the three cameras.

The objective is to compute an homography matrix $\pi_{ref} H_F$ that relates one camera frame F_t at a time t , to a so called cenital plane π_{ref} , i.e. a bird-eye representation of the scene.

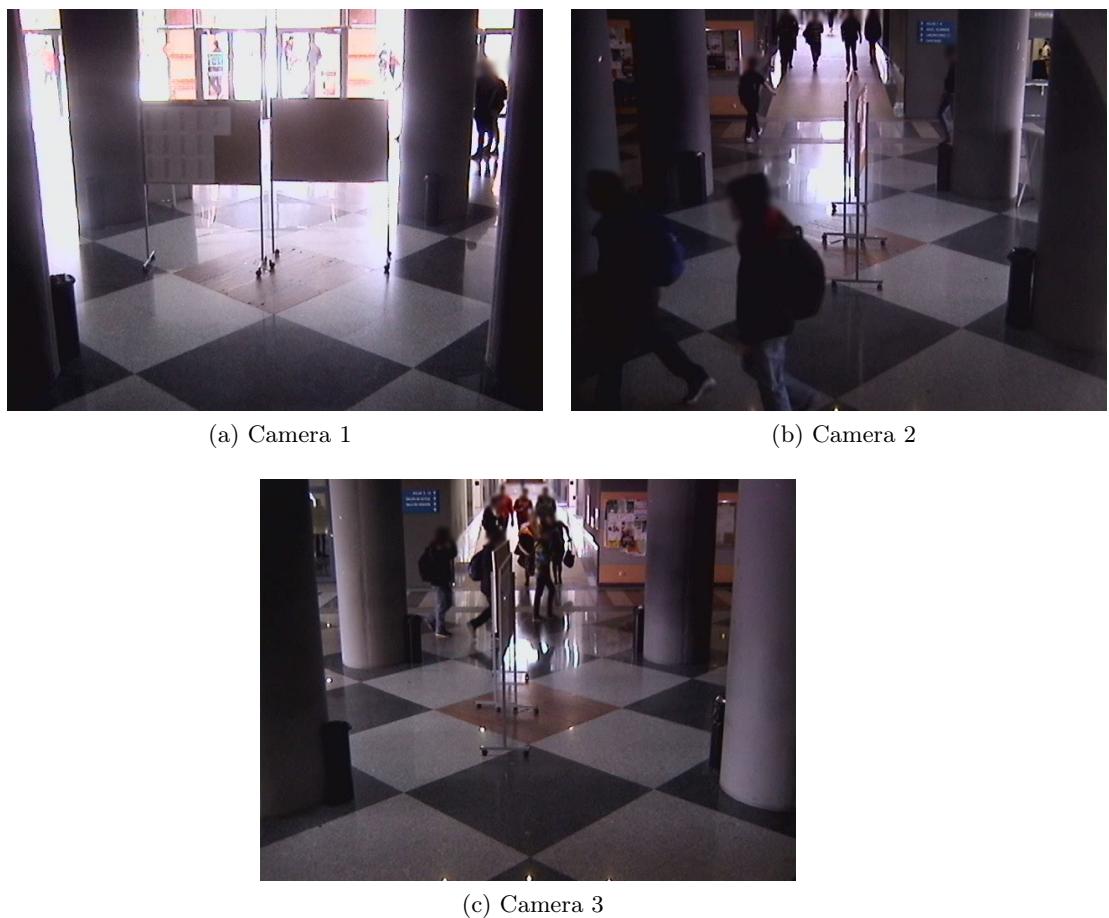


Figure 3.8: Initial Camera Views

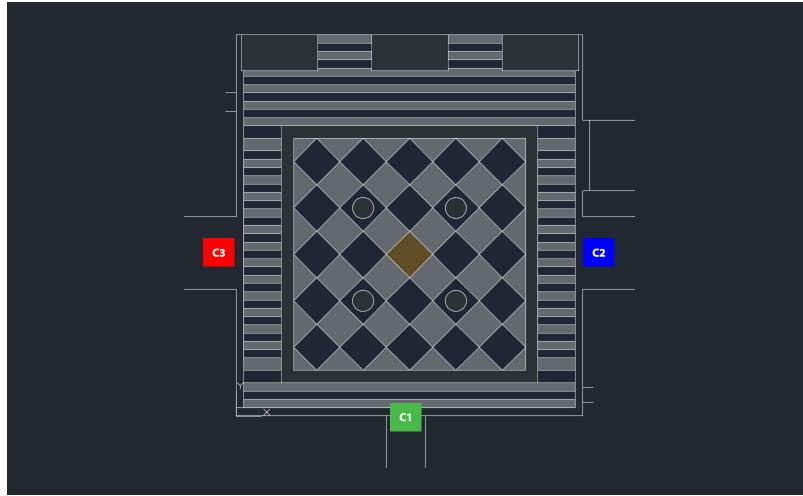


Figure 3.9: Cenital plane with camera positions.

With this homography matrix we are able to transform every frame pixel $F_t(x_1, y_1)$ to its correspondent position in the cenital plane $\pi_{ref}(x_2, y_2)$, i.e. from $2D \rightarrow 3D$. This process is done by Eq 3.1.

$$P_i = \lambda \cdot \pi_{ref} H_F \cdot p_i \quad (3.1)$$

where λ is a scale factor, P_i is the subset of points from π_{ref} and p_i is the subset of points of F .

Frame perspective is so, transformed as if it were being viewed from the top. The homography process is also important as it enables to project pedestrians detections to the cenital plane with the same proceeding.

In order to compute an homography, a relation between at least 4 points in each of the two images should be computed. In this case a relationship between a real camera frame and a cenital view plane created by computer is needed. This means that there is not a real correspondence for pixels and so it is not possible to use a point descriptor to extract common points in both images. User has to manually select the points that represent the same spatial place in both images using the Graphical User Interface and then the algorithm computes the transformation matrix.

3.3.1.2 Cenital Plan Design

In order to share detections between cameras a common plane is needed. For the proposed system the plane depict in Figure 3.9 has been used. The complete creation of the plane is detailed in Appendix A.

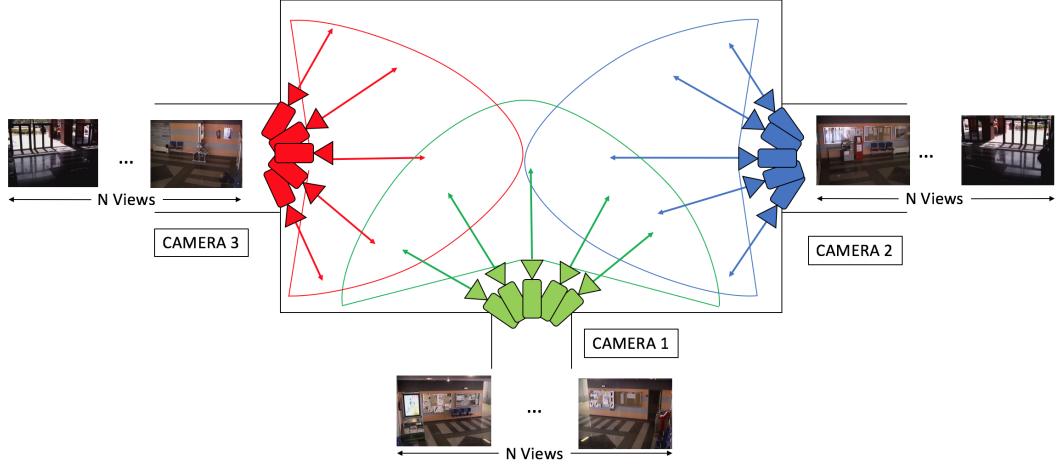


Figure 3.10: Multi-camera configuration with panning set up

This plane, due to its high amount of scene details allows the user to correctly select common points between it and a camera frame for the homography calculation.

3.3.1.3 Camera Panning

When cameras are static and pointing to the center of the scene as in Figure 3.8 and 3.8 one can easily observe that the vision range compared to the hole scenario representation is quite limited due to the low vision range of the cameras. The cameras are covering the middle part leaving completely unattended the lateral parts of the scenario. This approach is a limitation to the project objectives as high percentage of the scene semantic remains hidden.

The solution for this problem is to include panning movement thanks to the PTZ technology in all the cameras from left to right so cameras view range is increased and they are no longer focus only in the middle of the scene as before. This leads also to more common and not static areas which the cameras will share. This solution is presented in the scene diagram in Figure 3.10.

3.3.1.4 Camera Panning Discretization

Before including camera panning in the set up, homographies were calculated, at the beginning, for only one video frame per camera. This means that every camera only had one homography matrix as all its sequence was static and all its video frames were projected using the same matrix.

However, when including video panning this simple approach is no longer valid as the scene view is constantly changing. The ideal solution for this problem would be

having one homography matrix $\pi_{ref} H_{Frame}$ per video frame. As the homographies are calculated by the user selecting the points this solution is evidently impossible in terms of usability.

We propose discretization of the panning tour in which N views are selected from the video sequence. This leads to the obtention of an homography codebook in where the user computes a homography matrix $\pi_{ref} H_{View}$ for each of the N camera views. This codebook of views and homographies are responsible for projecting into the cenital plane the hole video.

3.3.1.5 View Selection

Due to the creation of the homography codebook we now have to choose between a set of N homography matrices to project detections into the cenital plane. The actual analyzed frame should be compared to each of the N views to obtain an spatial correspondence and so, use the correct homography matrix.

Inter-image Comparison

In order to compare two images we have done comparison between points of interest. AKAZE detector and descriptor [57] has been used for this task (AKAZE is explained in Appendix B).

AKAZE detects and describes points of interest in any image and then by a brute force comparator in terms of point distances we extract how many coincidence we have between a pair of images.

It is essential to say, that, as we have to compare the current frame with all of the N views there is a trade off. More views means a better space discretization and more overlapping between them, however, the computational time increases exponentially, this means that we have to choose the number of views N so it represents correctly the space and keeps the computational time relatively low.

In our proposed system we have choose $N = 9$ to maintain a trade off between projection quality due to the overlap between the views and time consumption.

Intermediate Homography

As we said before, the homographies are calculated for each of the views, this means, that if a frame is not positioned exactly as a view the homography matrix does not project the points correctly, there will be a small error.

Taking advantage of the calculation of common descriptors between the frame and its correspondent view this problem can be solved by calculating automatically the

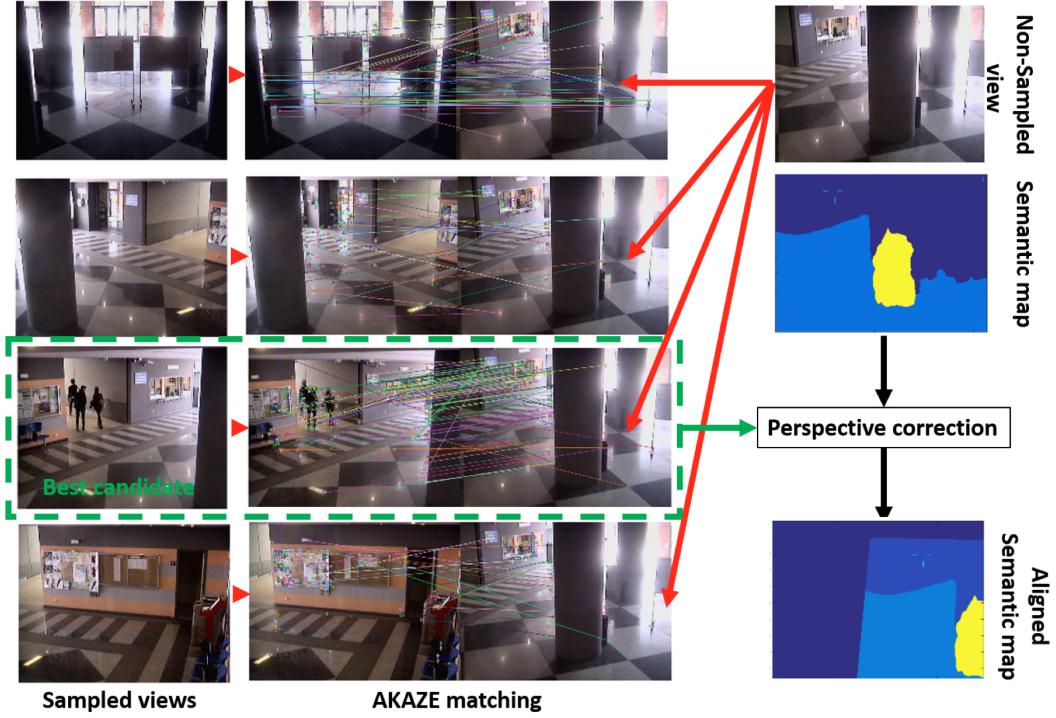


Figure 3.11: View selection process and homography between views computation.

homography between both images. This means that now, to project detections from the current frame to the cenital plane two homographies are used. First we change the perspective of the original frame to the perspective of the view with ${}^{View}H_{Frame}$. This process ensures that the used perspective is the same one as the one that was used previously to compute the homography.

Finally, the frame and its new perspective are projected to the cenital plane by means of ${}^{\pi_{ref}}H_{View}$. This process is graphically explained in Figure 3.11.

3.3.1.6 Video Sequence Synchronization

One of the main issues when dealing with multi camera systems, and specially, those that combine information between the cameras, is that video sequences coming from them should be correctly synchronized. This is really important because if the same frame number F is not representing the same exact moment in time in the three cameras, combination and fusion of the information is not possible.

We have synchronized the processed videos using the so called clapperboard technique. This technique aims to synchronize using concrete events that can be observed from all the cameras at the same time. By this, we are able to set a synchronize

starting point from where the following video frames are correctly representing the same moment in time.

3.3.2 Pedestrian Filtering and Constraining

During this Section we explain how the pedestrian detections are semantic filtered and also how the information between people from the three cameras is fusion in the system to increase the general algorithm performance.

3.3.2.1 Cylinder Estimation

In [54] a cylinder estimation technique is proposed. The detected bounding boxes on the camera frame do not correspond spatially with the exact position of the detected object due to the camera perspective. If this detection error is not corrected when the bounding boxes are projected either to another camera instance or to the common cenital plane there is a distance between the represented bounding box and the real object. Figure 3.12 shows the case for bounding box transference between cameras.

As we can observe in Figure 3.12 when the blue bounding boxes are projected from image (a) to image (c) they are not correctly on the pedestrian. The solution is to compute the cylinder that embraces the square whose side is the bounding box (b). Once the cylinder is estimated, the person will hypothetically be in the middle of the cylinder. In (c) we can observe that the estimation, if correctly computed has great accuracy.

In Figure 3.13 the same method is applied but in this case, the bounding box is not projected to another camera perspective but to the cenital plane. The projection of the bounding box, represented by the green line is not in the same position as the center of the cylinder, which is at the end of the purple line. The center of the cylinder so, corresponds to a more approximate position of the detected person.

3.3.2.2 Gaussian Representation of Bounding Boxes

Using also the cylinder representation one can change the way bounding boxes are represented in the cenital plane. As said before, in the simple cylinder representation pedestrian are represented with two perpendicular lines, however, in this representation one does have spatial information about the detection but none about the detection accuracy. This means that with the two perpendicular lines one can observed where the detection has been achieved but not the score associated with this detections. To solve this issue we propose a new representation method based on a Gaussian function placed at the middle of the cylinder.

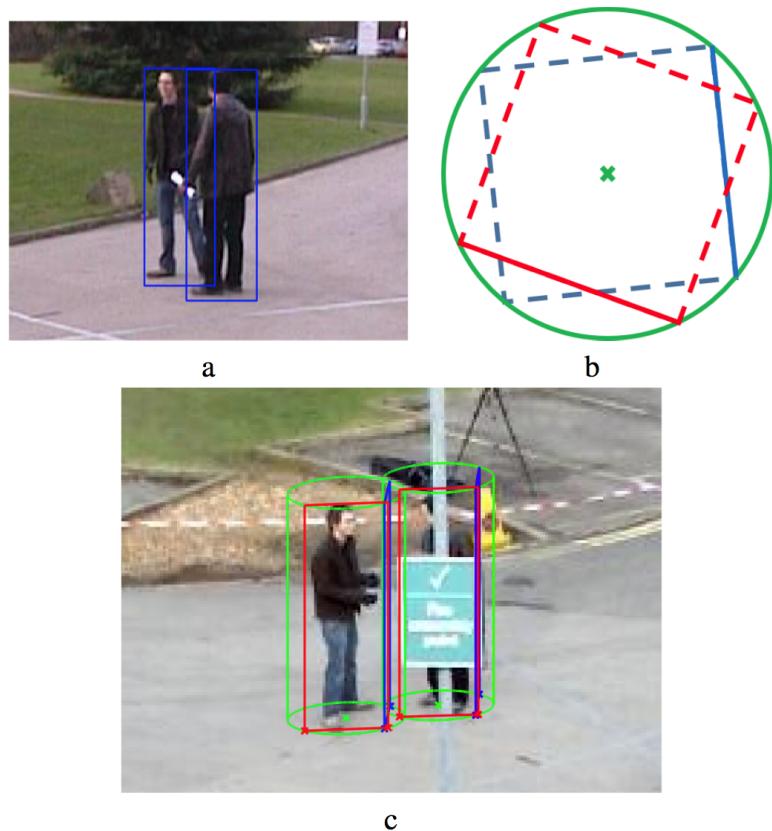


Figure 3.12: Cylinder estimation for camera instance projections. Extracted from [54].

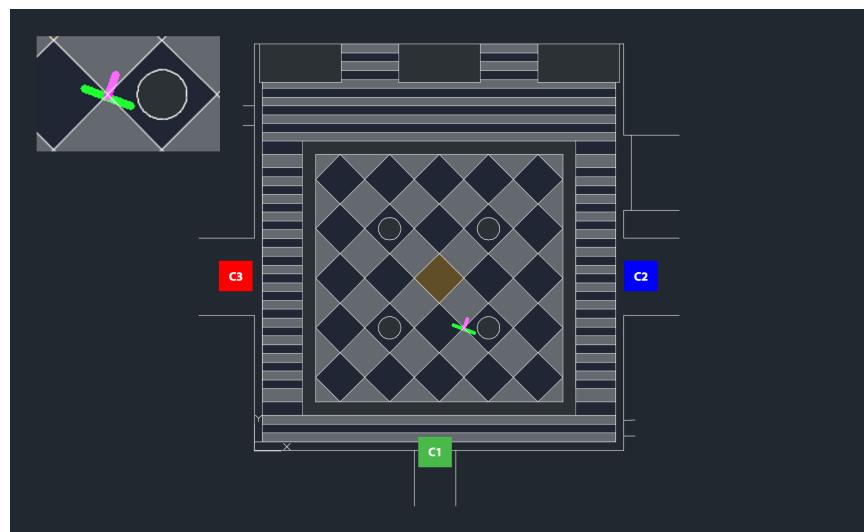


Figure 3.13: Cylinder estimation for cenital view projections.

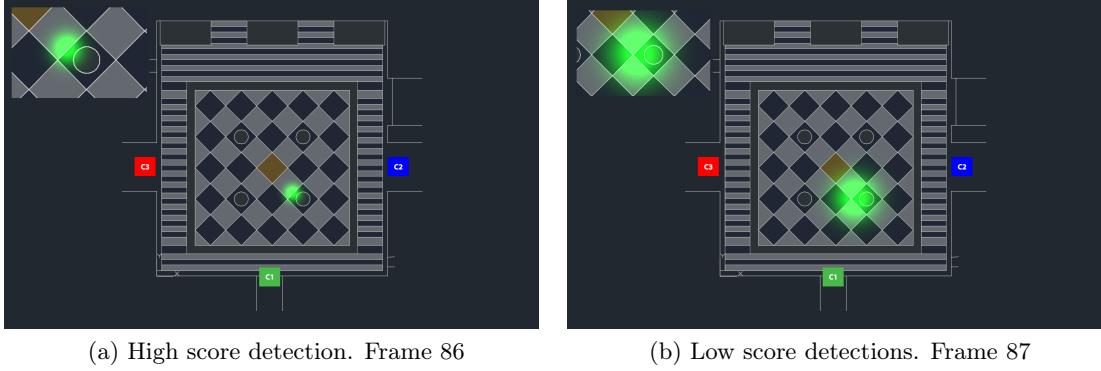


Figure 3.14: Gaussian representation examples.

Every pedestrian detection is so, represented as a Gaussian function of the form described in Eq. 3.2.

$$f(x, y) = A \exp \left(- \left(\frac{(x - x_0)^2}{2\sigma_x^2} + \frac{(y - y_0)^2}{2\sigma_y^2} \right) \right) \quad (3.2)$$

In this case A is the amplitude which typically is set to $A = 1$, x_0, y_0 represent the mean of the gaussian which, in our case, is the center of the cylinder associated to the bounding box and σ_x, σ_y which is the standard deviation and in our case represents the accuracy of the detection.

In order to relate a detection score S_n with a determinate standard deviation a simple rule explained in Eq 3.3 is proposed.

$$\sigma_x = \sigma_y = (S_n \cdot 10) + 5 \quad (3.3)$$

This Equation has been adapted to the size of the used cenital plane. This means that detections with higher scores are represented as a narrow gaussian function (Figure 3.14a) while lower scores lead to wide gaussians that represent a bigger area in where one can find that person (Figure 3.14b). Some Gaussians examples depending on the scores can be observed in Figure 3.14.

3.3.2.3 Pedestrian Reprojection Between Cameras

When all the pedestrian detections have been projected into the cenital plane one can start reprojecting those projections from the cenital plane back to the other camera frames.

This process has a fundamental importance as it is from where the pedestrian detector accuracy can be increase by the use of the multi-camera system. For instance,

detections from Cameras 2 and 3 are reprojected into Camera 1. Sometimes, those reprojections are not in the frame because the cameras are not aiming to the same spatial area, however, when they are seeing the same spatial space cameras share those detections. Ideally, if the three cameras detect the person reprojection could not be necessary, but, if one of them misses the detection, the other two detections, by the use of reprojection, lead to suppress that miss detection.

Projected frames from one camera to another are treated as another detector blob. This means that if the detection is accurate enough in the original camera, when projected to another frame it could be either joined by a Non-Maximum Suppression (NMS), if that camera has already one blob, or maintained as a detection on the person. However, if the detection is not accurate enough, the projection could be not on the person and so that blob leads to a false positive. This can be described as the main drawback of the method.

In addition, when blobs are projected from one frame to another blob height is lost due to the cenital projection. In order to achieve a complete blob once it has been reprojected, aspect ratio is used to, using the width of the blob obtain the height.

This process can be observed in Figure 3.15.

3.3.2.4 Semantic Constraining of Pedestrian Detections

Once all the detections have been reprojected into all the cameras one can use semantic information in order to filter false positives detections in not correct semantic scene areas.

Ideally people are always walking on the floor so the common floor areas between pairs of cameras, whose calculation is deeply explained Section 3.4.2, are used to constrain the pedestrian detections to this hypothesis. This process is done based on the cenital plane builded using the cylinder estimation explained before.

The constrain idea is that a bounding box coming from a camera C_1 is assumed correctly computed if the center of its related cylinder is on the common floor between C_1 and C_2 and also between C_1 and C_3 . An example of the constraining can be observed in Figure 3.16.

3.4 Reference RGB and Semantic Planes

During this Section our proposed method to create both RGB and semantic reference planes combining information from all the cameras positions is explained. The final objective is to have complete maps in each camera that represent the scene.

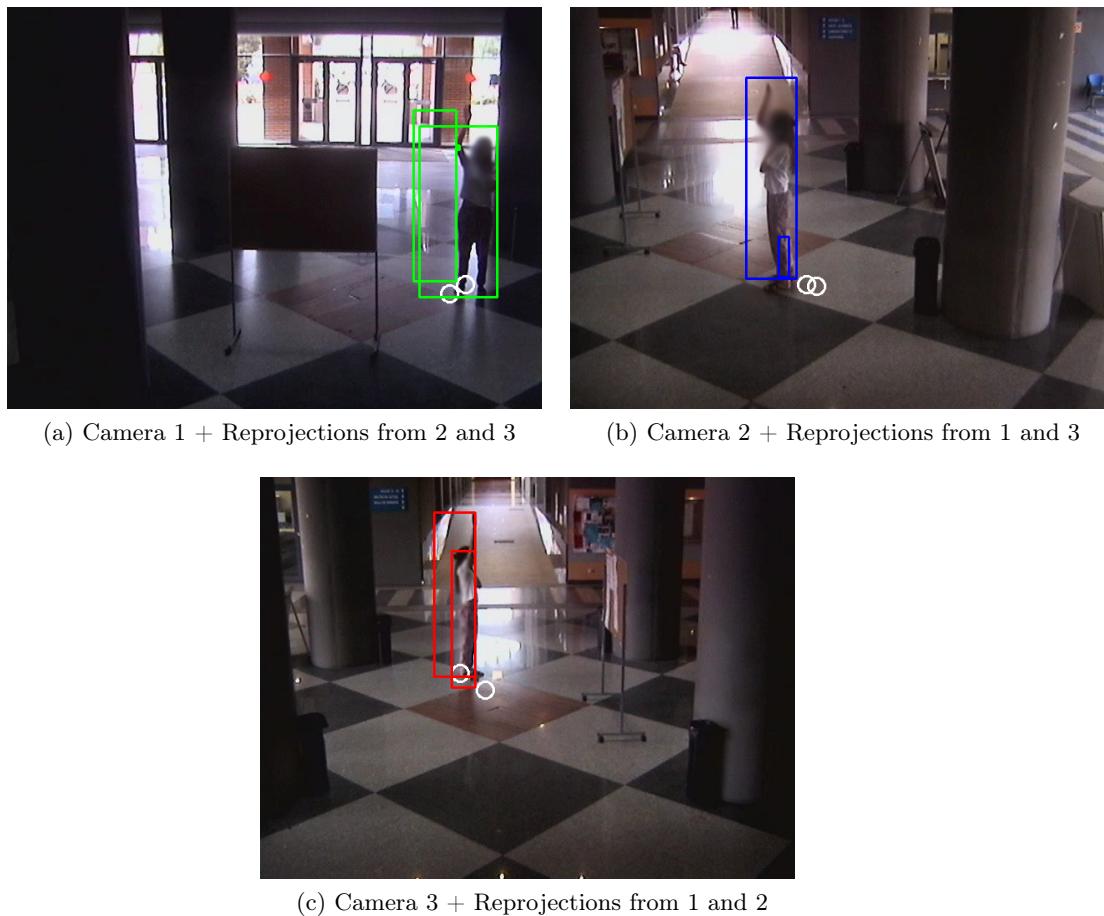


Figure 3.15: Pedestrian detection reprojection.



Figure 3.16: Pedestrian semantic constraining. Blobs that are not in the right semantic area are just displayed by a circle and text label.

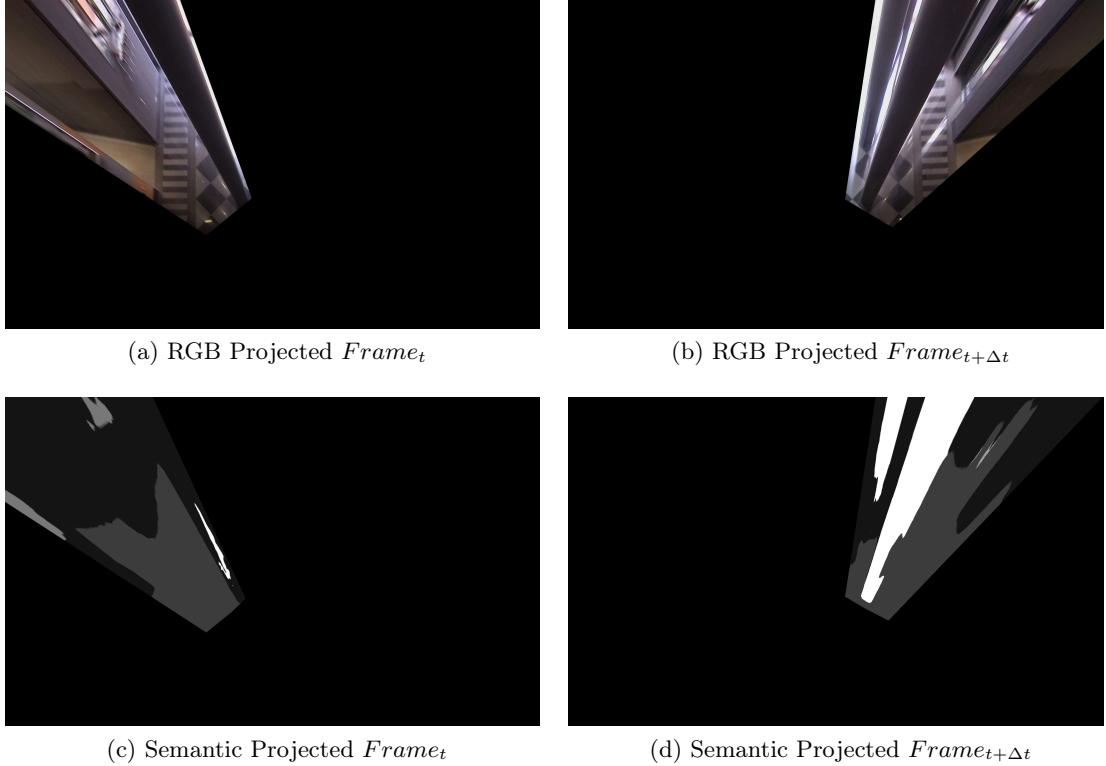


Figure 3.17: RGB and semantic frames projected.

3.4.1 Semantic and RGB Reference Plane Generation

Following with the process scheme of Figure 3.11 one can project every video frame into a cenital perspective. This way a set of projected RGB and semantic frames for each camera is obtained. Some examples of these images can be observed in Figure 3.17.

Once all the video frames are projected we propose to apply a temporal average through a median filter for the three cameras separately. By this method a reference plane π_{ref-C} for each camera contained in the ground plane is created by joining all the separated projected frames. Given a set of pixels $P(x_0, y_0, N)$ from the same spatial position (x_0, y_0) from N different projected images one can define median filter as in Eq 3.4

$$P(\tilde{x}_0, \tilde{y}_0) = \begin{cases} P_{(N+1)/2} & \text{if } N \text{ is odd} \\ \frac{1}{2}(P_{N/2} + P_{1+N/2}) & \text{if } N \text{ is even} \end{cases} \quad (3.4)$$

The result coming from the temporal averaging can be observed in Figure 3.18 for

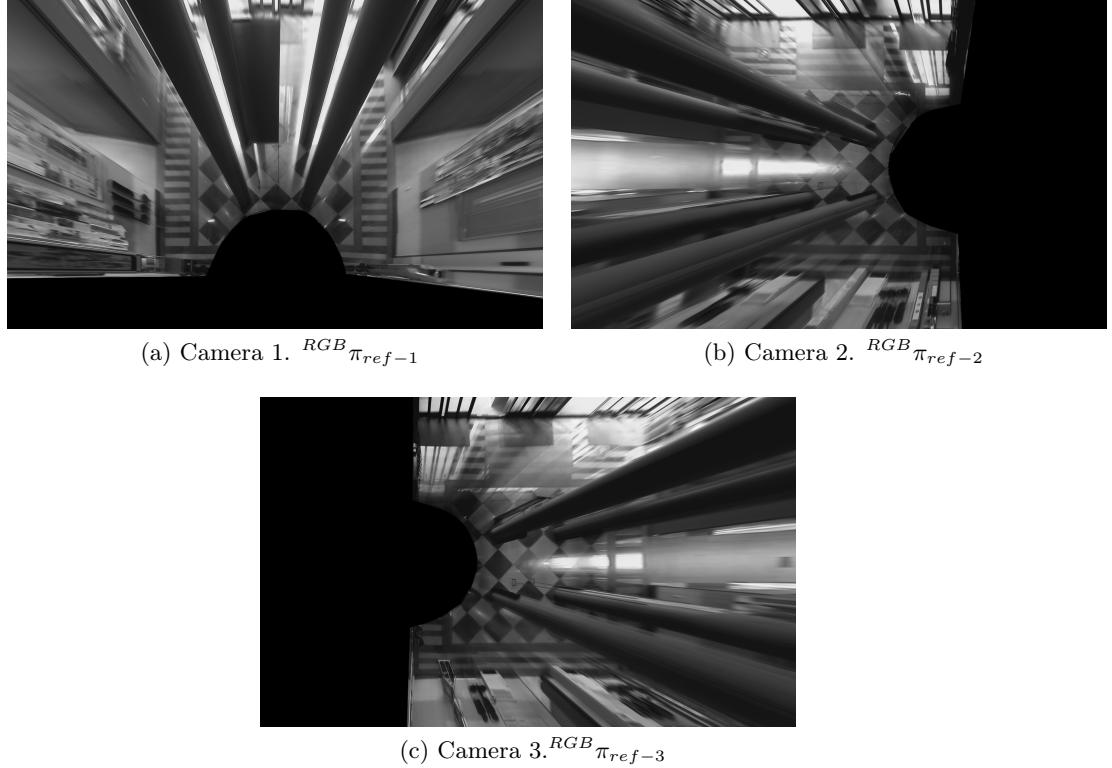


Figure 3.18: RGB Reference Planes $\text{RGB}_{\pi_{ref-C}}$.

RBG case and in Figure 3.19 for semantic information.

As one can observe smooth semantic and RGB maps from all the hall area are created. It is important to recall that only if the homographies and the ground floor are correctly calculated and projected, the base of for instance the columns from all the different cameras, should lay in the same spatial point in these images.

3.4.2 Semantic Fusion in the Reference Plane

Once the cenital maps have been extracted one can combine them to create semantic fusion in the reference plane. The main objective is to obtain common areas, i.e. pixels with the same label within two or more cameras. This method allows us to increment the confidence for those pixels and so, we can use this information for example, to constrain pedestrians detections. As the reference plane is obtained by the homography to the ground plane we have only take into account floor pixels to create the common areas, as the others have projection errors that lead to fusion errors.

First, pairs of cameras are combined to create three common semantic areas that

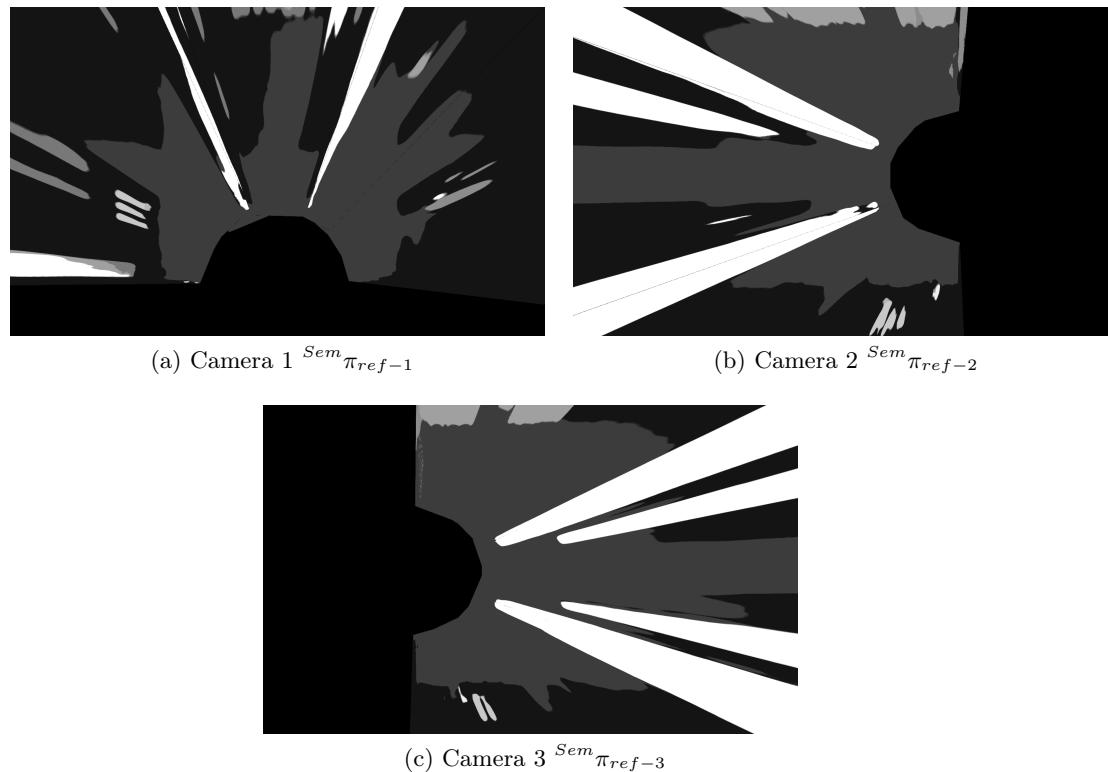


Figure 3.19: Semantic Median Average $S_{sem}^{\pi_{ref-C}}$.

are shown in Figure 3.20.

Once pairs of cameras have been combined one can go further and extract common areas for the three cameras at the same time. Those areas are pixels with the exact same label for the three cameras so its probability to be that class is really high. Figure 3.21 represents the common areas.

As we can observe, due to the multi camera setup and some scene occlusions such as scene panels or columns, the common area shared by all the cameras represents a small portion of the room. This is the reason why we have choose common areas between pairs of cameras to work along with pedestrian detection rather than common area between all of them.

3.5 Statistical Usage Data

Along this Section we explain the proposed system to extract statistical usage data frame by frame given a sequence. The aim is to obtain one graph per selected class that measures the amount of people at a moment t of time in the determined semantic area during the sequence.

3.5.1 Statistical Semantic Map Generation

In order to extract usage data the first step is to create a projected common semantic map. For this task, rather than creating common areas between cameras as done for the PD constraining, information for all the cameras have been joined without strictly being shared. This have been done to preserve as much semantic as possible and so, have bigger floor or door areas than if we were more restrictive.

The result from this step is a unique map that represents the hole spatial area (Figure 3.22).

3.5.2 Usage Curves and Paths

As one can observe we have decided to only take into account for this process floor, doors and chairs, which are the most accurate detections from the semantic segmentation. Following the same principle as in the previous Section, pedestrians are projected on to the statistical semantic map and so one can be able to know how many people are in each of the designed areas.

In addition, one can divide Figure A.3 into different subregions of fixed size and extract in which of them pedestrians are more likely to be. This process leads to a set of most used paths per subregions.

(a) Cameras 1 and 2. ${}^{Sem}\pi_{ref-1} \wedge {}^{Sem}\pi_{ref-2}$ (b) Cameras 2 and 3. ${}^{Sem}\pi_{ref-2} \wedge {}^{Sem}\pi_{ref-3}$ (c) Cameras 3 and 1. ${}^{Sem}\pi_{ref-3} \wedge {}^{Sem}\pi_{ref-1}$

Figure 3.20: Common semantic areas between pair of cameras.



Figure 3.21: Common semantic areas for all the cameras $Sem_{\pi_{ref-1}} \wedge Sem_{\pi_{ref-2}} \wedge Sem_{\pi_{ref-3}}$

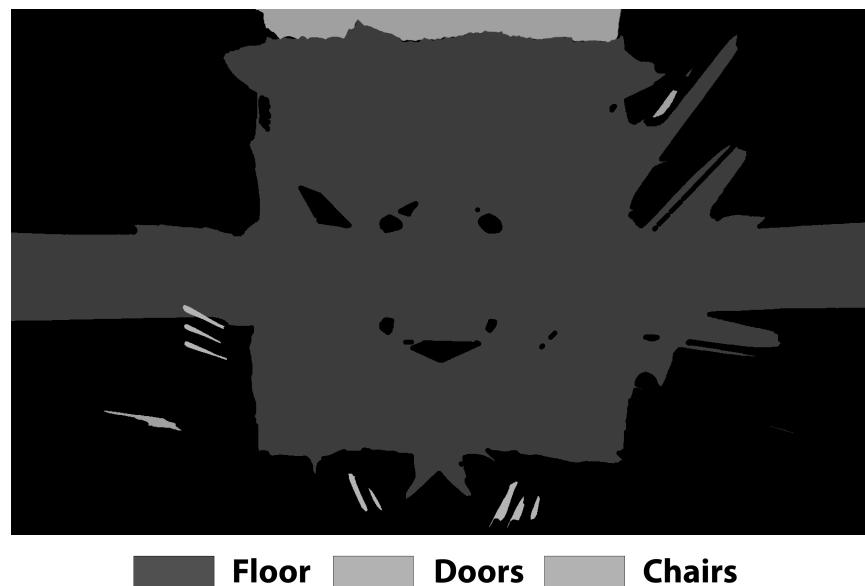


Figure 3.22: Statistical semantic map.

Chapter 4

Developed Application

Within this Chapter the developed application in terms of software development is described.

This application is the base for the integration of pedestrian detection algorithms as well as semantic segmentation. Visualization and arrangement of the usage statistics from the different areas of interest is also done by the software.

Application environment should be user-friendly to ensure a correct and easy usage by the end user. It has been developed completely from scratch for the purpose of this Master Thesis.

The application has been developed under [QT Creator](#)¹ coding environment in Mac OS Sierra. This decision has fundamentally been based on the following QT characteristics:

1. Its cross-platform characteristic which makes it easily portable from one operating system to another such as Windows or Linux distributions.
2. Its application window designer that allows the programmer to design software windows by using an interface instead of having to create windows by coding (Figure 4.1).
3. The possibility to add [OpenCV](#)² libraries to the project as well as independent external libraries.
4. Its multi-thread capabilities that enables to perform different code segments in various threads to increase computational speed.

¹<https://www.qt.io/>

²<http://opencv.org/>

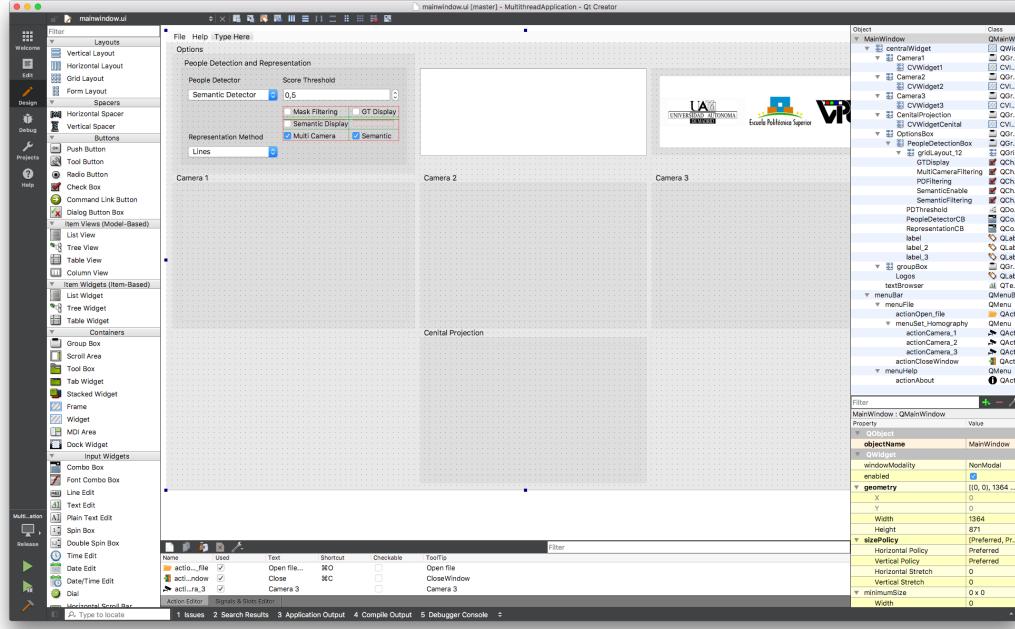


Figure 4.1: QT Main Window Designer

Due to the complexity of some of the algorithms used, in terms of parameter tuning and configuration, two separate applications have been developed: developer and user version.

The first version of the software corresponds to the developer application. It is designed so it can be used by programmers or engineers who generally understand concepts of the algorithms running at the backend application. This means that:

1. Variable parameters are available for tuning from the graphical interface.
2. Different pedestrian detectors can be selected.
3. Options and tuning for these algorithms can be done.
4. Results are displayed in different areas.

This version allows to change parameters and methods online. However, this assumes that the user has basic knowledge on how parameters affect the software performance.

For the developer application both single-thread and multi-thread versions have been developed and are discussed on the following lines. All the code for both approaches is available in a [GitHub Repository](https://github.com/alexlopezcifuentes/PCV-MasterThesis)³.

³<https://github.com/alexlopezcifuentes/PCV-MasterThesis>

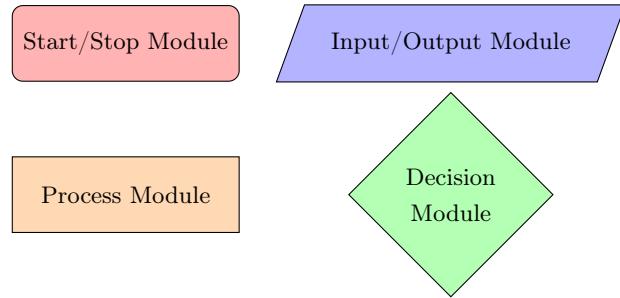


Figure 4.2: Flow-chart legend

During the analysis of the application through the Section some flow-charts are displayed. A common legend for all of them is included in Figure 4.2.

4.1 Single-thread Developer Application

During the first stages of the development and for the sake of simplicity the application has been designed and developed to run under a single thread. This means that all the processing has been done sequentially camera by camera. A simple flow-chart diagram that illustrates the execution path can be seen in Figure 4.3 .

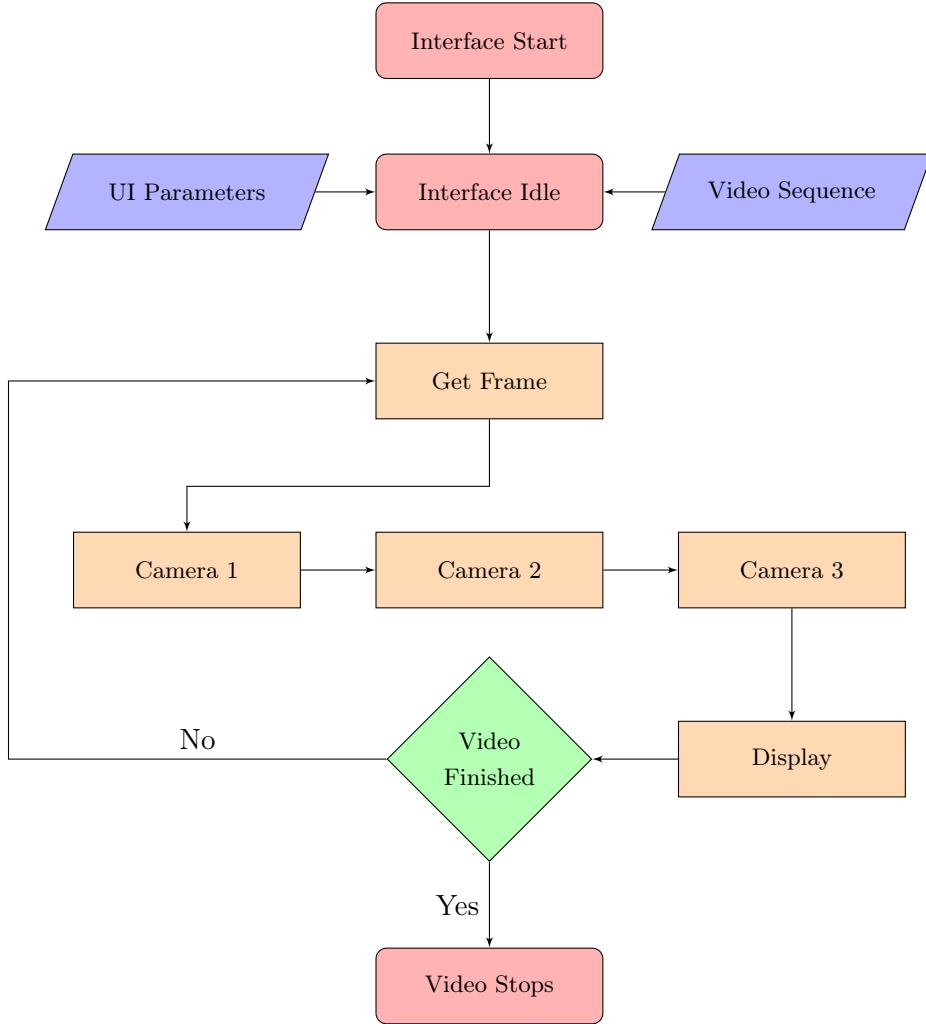


Figure 4.3: Flow-chart diagram for the single-thread application.

This approach has the advantage that all the code is executed in the same memory segment. This makes really easy, for instance, to share information between cameras. This design is however, only valid if the computational effort is minimum. All the process for the three cameras should be computed one after the other which means that when calculating detections for one camera, the others remain idle. When working with such a multi-camera system with heavy algorithms running –as in the proposed method– the computational time increases exponentially and this design is no longer worthwhile.

4.2 Multi-thread Developer Application

Multi-thread Developer approach can be observed in the flow-chart displayed in Figure 4.4. Now different threads are running in parallel, one for each camera, and so, all the process is no longer done sequentially and computing power of the CPU can be further exploited.

However, as threads are running separately a synchronization strategy should be included to keep consistency in the application.

One thread can process a frame faster than another one due to multiple external reasons, nevertheless, the application should display the same exact frame for all the cameras at the same time. This is specially relevant if results are going to be shared between threads. In our case the synchronization is performed by two barriers –see diagram 4.4–.

- The first one ensures that all the threads have perform PD before sharing these detections to the rest of the threads.
- Second barrier creates a meeting point at the end of the frame processing so a thread waits to the others before sending results to the main display.

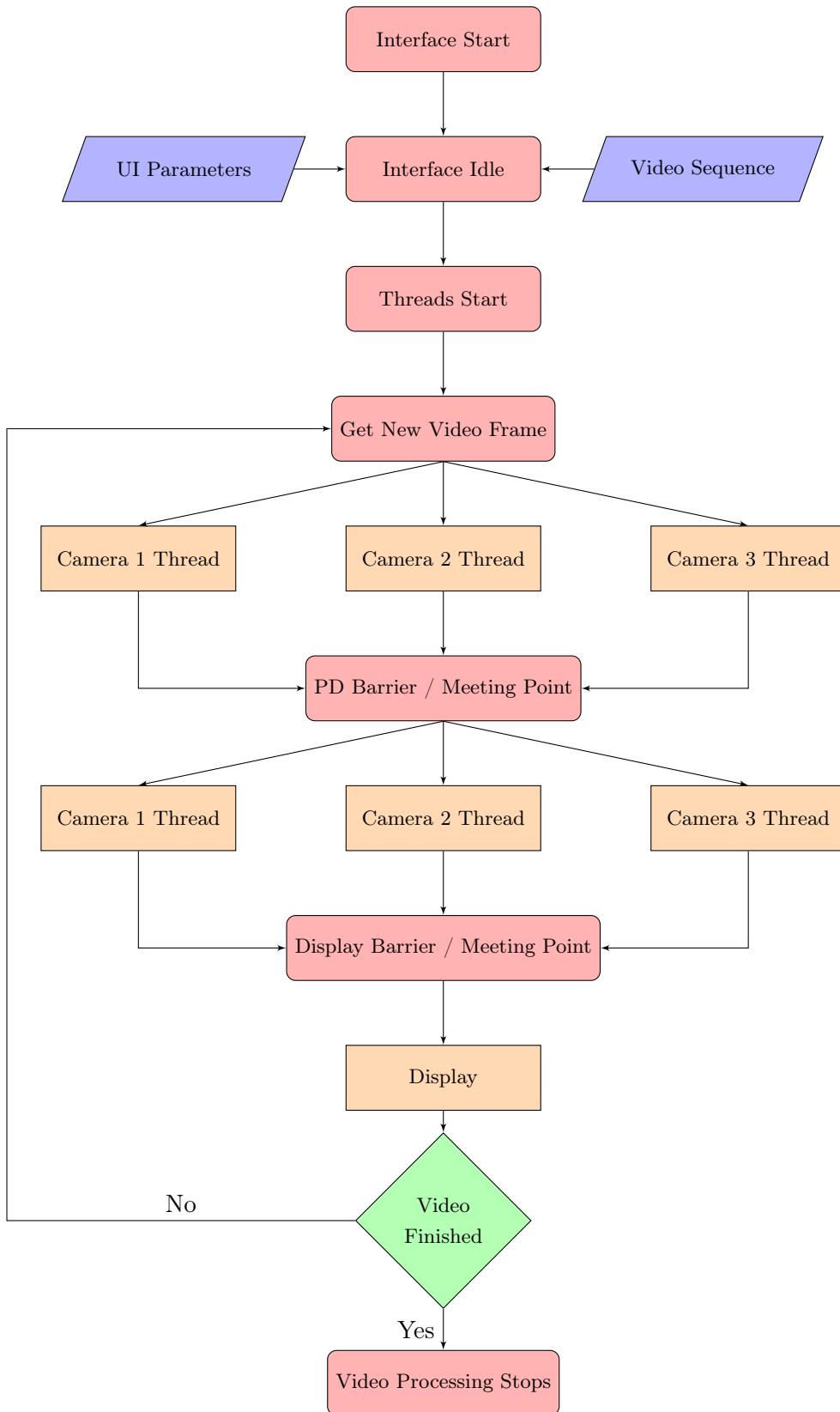


Figure 4.4: Flow-chart diagram for the multi-thread application.

4.2.1 Main Application Window

Main application window is shown in Figure 4.5. As one can observe it is composed of four separate areas:

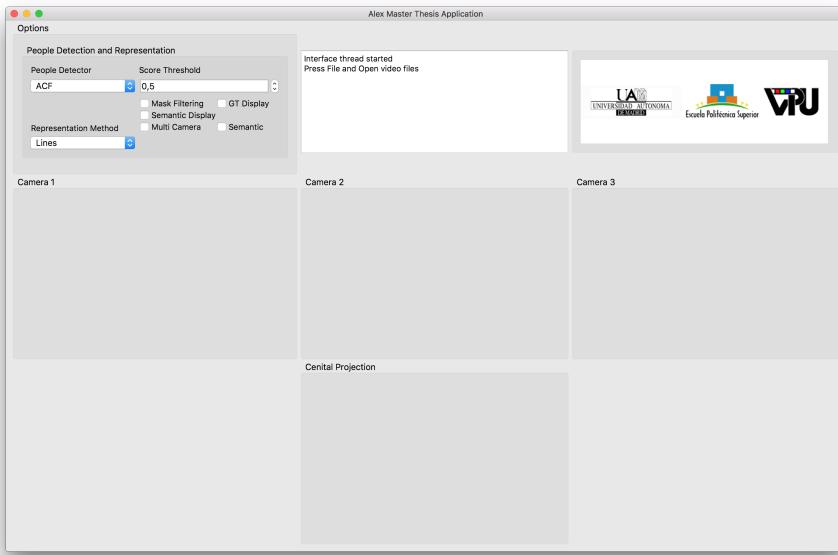


Figure 4.5: Main application window.

4.2.1.1 Application Menu Bar

In the menu depicted in Figure 4.6 the main application actions are contained. From here the user can:

1. Open a new video sequence.
2. Compute the set of needed homographies for the integrated algorithms.
3. Close the application.
4. Search through the help searcher.
5. Open the external information window.

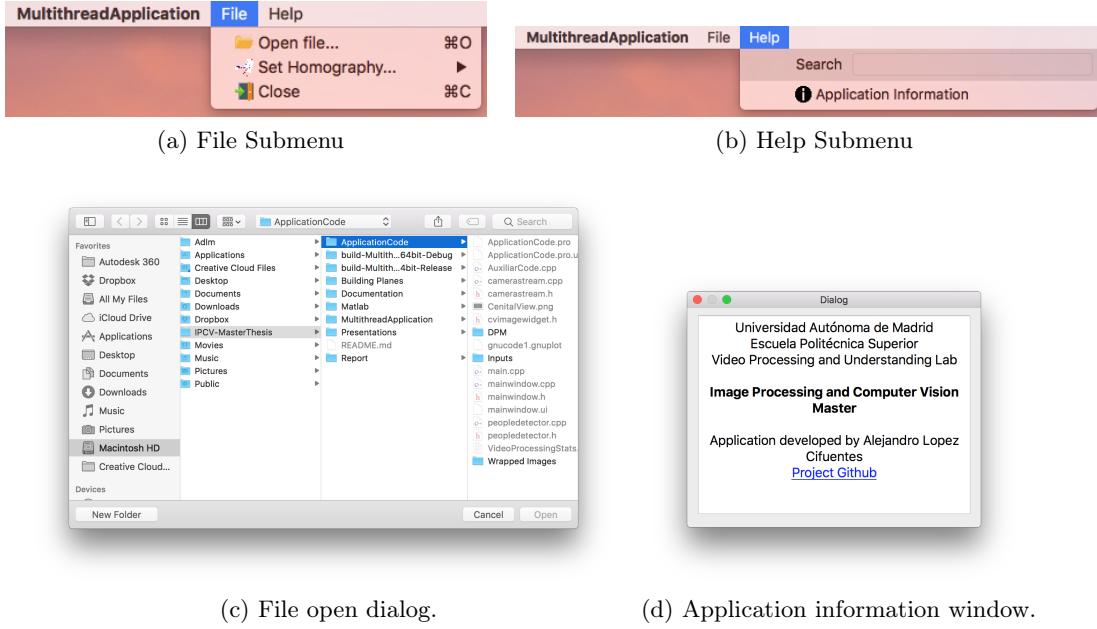


Figure 4.6: Application menu bar

4.2.1.2 Options Menu

The options box (see Figure 4.7) in the application contains all the possible parameters that can be tuned during the execution of the program. From here algorithms can be changed in real time so there is no need to restart the execution before changing some parameter. From here user can change:

- Pedestrian detectors. The user can select among the following ones:
 - ❖ PSP-Net detector
 - ❖ HOG
 - ❖ DPM
 - ❖ ACF
 - ❖ Fast-RCNN
- Different representation methods for PD detections as explained in Section 3.3.2:
 - ❖ Lines
 - ❖ Gaussians
- Enables the user to select the threshold for PD algorithms.

- PD Filtering or constraint as explained in Section 3.3.2 can also be changed.
The available options are:
 - ❖ Raw PD
 - ❖ PD with semantic constraining.
 - ❖ PD with multi camera reprojection.
- Mask filtering option to perform PD over a limited area.
- Ground truth check box to display or not ground truth information.

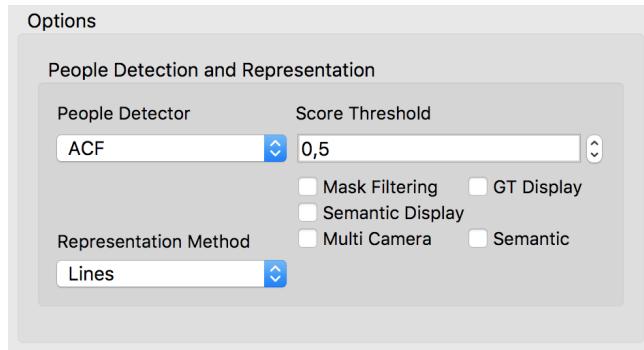


Figure 4.7: Options Menu.

4.2.1.3 Information Display

Along this text box status information is provided to the user. Messages such as “Open video files”, “Processing starts now” or “DPM Pedestrian Detector is now in use” appear during the execution of the application so the user can obtain some information about what to do, or what algorithm is in use. This can be observed in Figure 4.8.

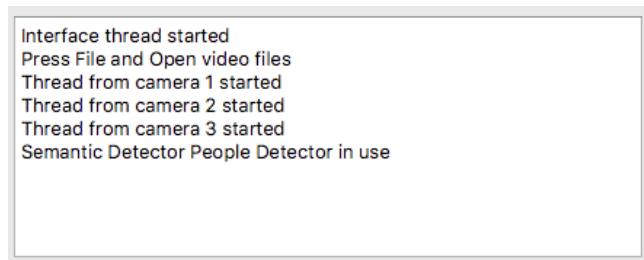


Figure 4.8: Information Display

4.2.1.4 Results Display

This is the main display area in the application in which all the visual results are presented.

We have three separate windows for each of the used cameras as well as one more display window for the cenital plane. Here the camera frames and associated detections and/or ground truth are shown: Besides all the projected semantic can be observed on the cenital frame. An example is depicted in Figure 4.9.

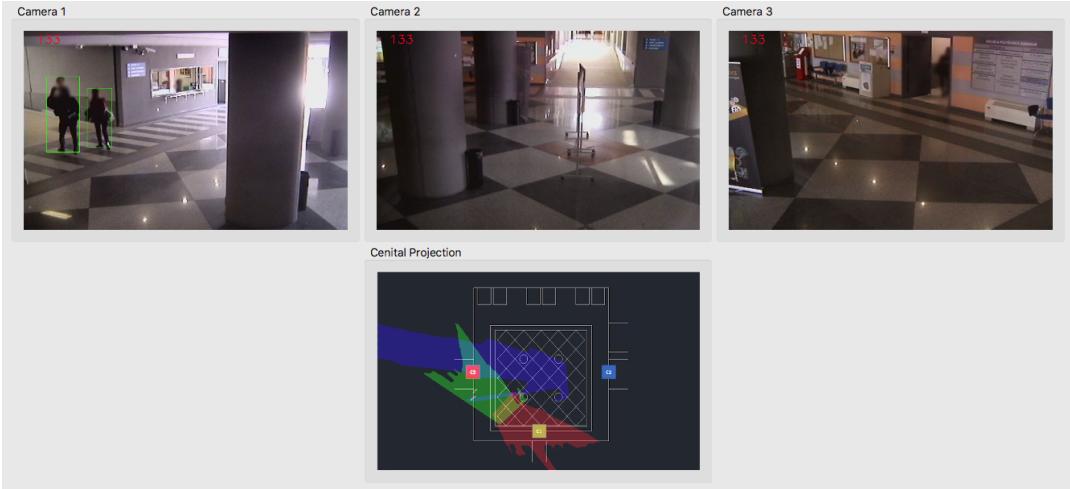


Figure 4.9: Results Display Area

4.2.2 Classes Distribution

In terms of C++ basic units the application has been divided into several classes for a better code design and to ease code comprehension / interpretation.

- **MainWindow:** This class corresponds to the main interface window and main application thread. It is the base for all the further processing as everything is inherited from this class. The reason for that is that **MainWindow** class is used to create threads and sharing procedures between them. Associated functions for this class are:

- ❖ Creating file open dialogs.
- ❖ Setting up and start all the camera threads.
- ❖ Update all the algorithms configurations from the UI.
- ❖ Displaying results through the CVImageWidget class.
- ❖ Sharing information between threads.

- **AboutWindow:** Class that executes the second available information window. This instance displays general application information.
- **CameraWorker:** Main class for all the execution in each of the cameras. **CameraWorker** class is linked with a unique thread that process all the algorithms inside. It has **CameraStream**, **PeopleDetector**, **Evaluation** and **Barrier** classes declared within it to distribute the processing.
- **CameraStream:** This class includes all the functions that are related to video processing except PD:
 - ❖ Main sequence reading loop.
 - ❖ Homographies calculations
 - ❖ Semantic projections.
- **PeopleDetector:** Main class to carry PD out. All the functions to detect, project and draw results either on the camera frame or on the cenital plane are in this class.
- **Evaluation:** Here ground truth is read and also the evaluation between system pedestrian detection and ground truth information is performed.
- **Barrier:** This class deals with thread synchronization. It is declared in **MainWindow** and passed by arguments to the thread so each of them has the same exact barrier object to perform the synchronization.
- **CVImageWidget:** Display representation class that deals with all the processes to draw OpenCV Mat images into the QT main window interface Widget.

In Figure 4.10 a hierarchical representation of how the different objects are arranged is presented. As one can observe everything is under the heritage of the **MainWindow** object. Here we have three **CameraWorker** threads that include the **barrier** object, and three **CameraStream**, **PeopleDetector** and **Evaluator** –one for each one–. In addition, **MainWindow** instantiates also **AboutWindow** and **CVImageWidget** objects.

4.3 Multi-thread User Application

On the contrary to the discussed version, the second developed application is focused on general users. This version only allows to load the video files and display results. All the parameters are set by default so the user does not have to fine tune any of

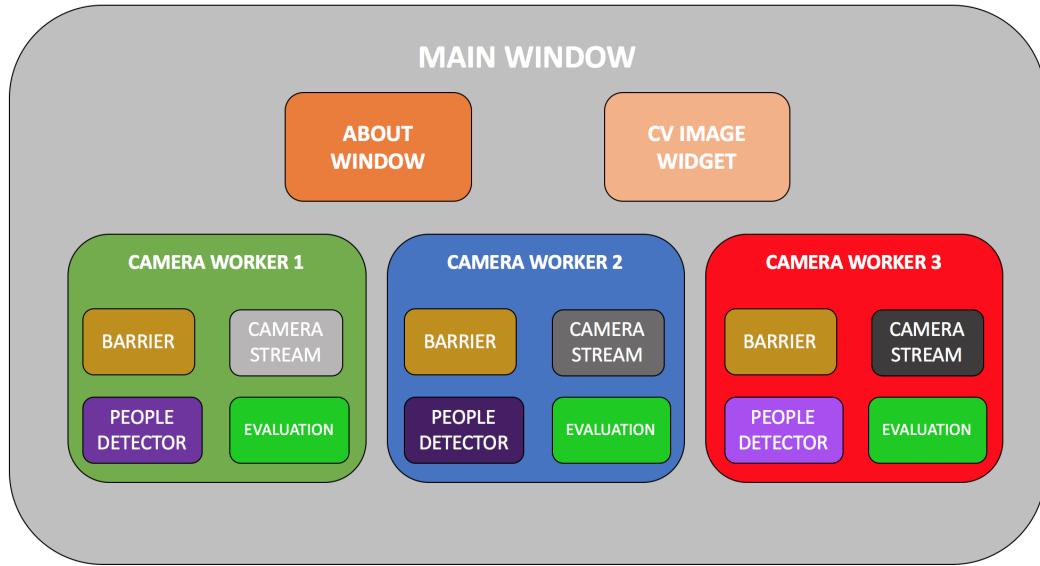


Figure 4.10: Hierarchical representation of the code.

them. This turns the usage and the general perception of the application much more simple and easy. Default setup is parametrized as the best observed configuration in the results (see Chapter 5).

Figure 4.11 displays the general application window where we can observe the difference with the developer version.

Conversely to the developer application, the only functionality of the main window is to display results and guide information. All the options presented in the developer application for parameter setup are no longer available.

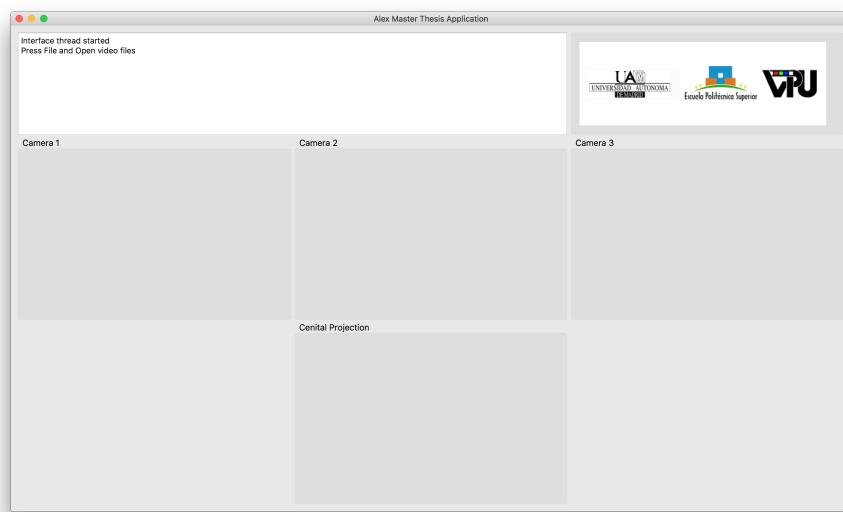


Figure 4.11: User version application main window.

Chapter 5

Results

During this Chapter all the achieved results are exposed and analyzed. Firstly, a brief analysis of the used hardware during the Thesis is done. Secondly, in detail, the experiment setup is explained, i.e. the generated data-set, ground truth, and the evaluation framework. Later, results concerning application performance are presented. Homography projections results are analyzed. Finally, the last two sections, aim is to present pedestrian detection results in term of performance and usage data extraction.

5.1 Used Hardware

The project has been developed in the Escuela Politécnica Superior (Universidad Autónoma de Madrid). Due to this fact, the testing environment has been the hall of the mentioned engineering school which has a setup of three Internet Protocol Cameras (IP Cameras). This type of cameras can send and receive data via a computer network and Internet which allows the user to set the configuration and receive frames from the cameras.

5.1.1 Camera Specifications

Specifically, the camera model used along the project has been the Sony SNC-RZ50P PTZ Camera. This is a PTZ camera which means that is able to Pan, Tilt and Zoom all over the scene. Precisely this camera has a pan range of 340 degrees and a tilt range of 115 degrees, enabling users to monitor a wide area over the scene if the camera is moved (Figure 5.1) .

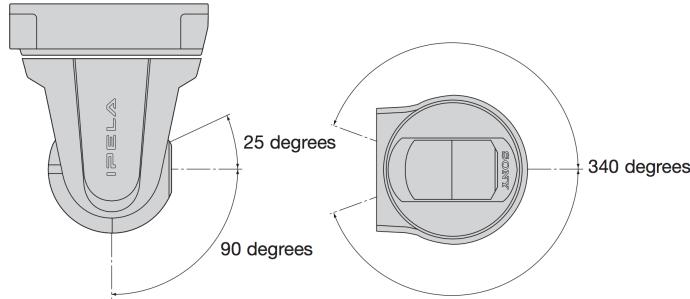


Figure 5.1: Camera Sony SNC-RZ50P Pan/Tilt Range diagram

The complete and relevant specifications from the cameras are detailed in Table 5.1. In the scope of our work the most important features are resolution and frames per second that enable us to have more quality in both image and video.

Camera	
Horizontal viewing angle	1.7 to 42.0 degrees
Focal length	f = 3.5 to 91.0 mm
F-number	F1.6 (wide), F3.8 (tele)
Minimum object distance	320 mm (wide), 1,500 mm (tele)
Pan angle	-170 to +170 degrees
Pan speed	300 degrees/s (max.)
Tilt angle	-90 to +25 degrees
Tilt speed	300 degrees/s (max.)

Image		
Image size (H x V)	640 x 480, 320 x 240, 160 x 120	
Compression format	JPEG, MPEG-4, H.264	
Maximum frame rate	JPEG/MPEG-4	25 fps (640 x 480)
	H.264	8 fps (640 x 480)

Table 5.1: Camera Sony SNC-RZ50P Specifications

5.1.2 Camera User Web Interface (GUI)

The camera comes with a built-in web interface that helps the user to visualize the visual range and set the different parameters that would change the camera behavior. The most important features that users are able to tune are described as follow:

- Camera control: Through this interface one can control and set the position of the camera in terms of pan, tilt and zoom (Figure 5.2). Changes in this three variables lead to different visualizations of the scene.

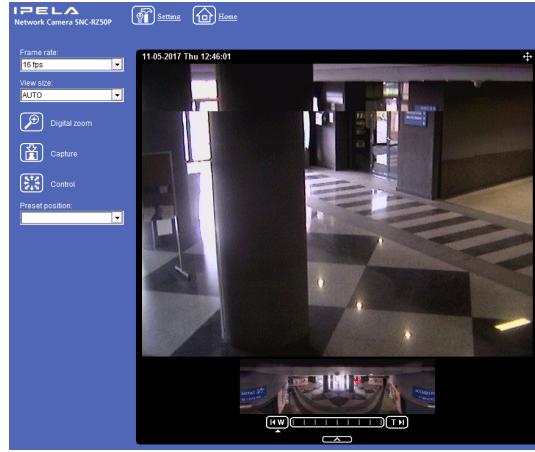


Figure 5.2: Visualization and control menu

- Preset position: In this menu (Figure 5.3) one could save the position that has been set in Figure 5.2 in order to recover the same position if the camera has been moved before in precise and easily manner.



Figure 5.3: Preset position setting menu

- Tour setting: One can set the camera to describe a tour over the scene. This process is done by setting at least two different preset positions from where the

camera is moving from one to the other at a set speed. The menu to configure this behavior is displayed in Figure 5.4.

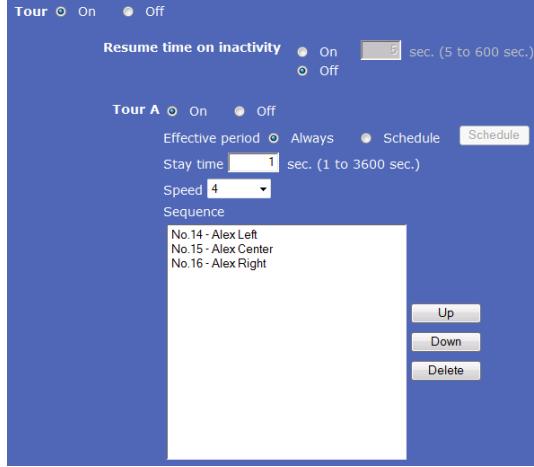


Figure 5.4: Tour setting menu

5.2 Experiment Setup

For all the experiments that are going to be analyzed along this Chapter the experiment setup should be explained. In this section, the generated data-set, ground truth and the proposed evaluation framework used to obtain results are detailed.

5.2.1 Data-set

Data-set has been generated in the university hall. A tour for each camera has been set so they move from left to right only changing position in the X/horizontal axis and not in the Y/vertical direction because of the reasons presented in Chapter 3.

Both settings are configured as previously displayed in the diagram of Figure 3.10.

The complete data-set consists on a set of three different recordings about 5 minute long from three different days. In each of the recordings all the camera sequences are available. Figure 5.5 shows some example of them. The technical characteristics of the videos are the detailed in the following lines:

- Number of videos per recording: 3
- Resolution: 640x480 pixels.
- Frame rate: 23.976 fps.
- Format: MPEG Video.



Figure 5.5: Data-set example frames

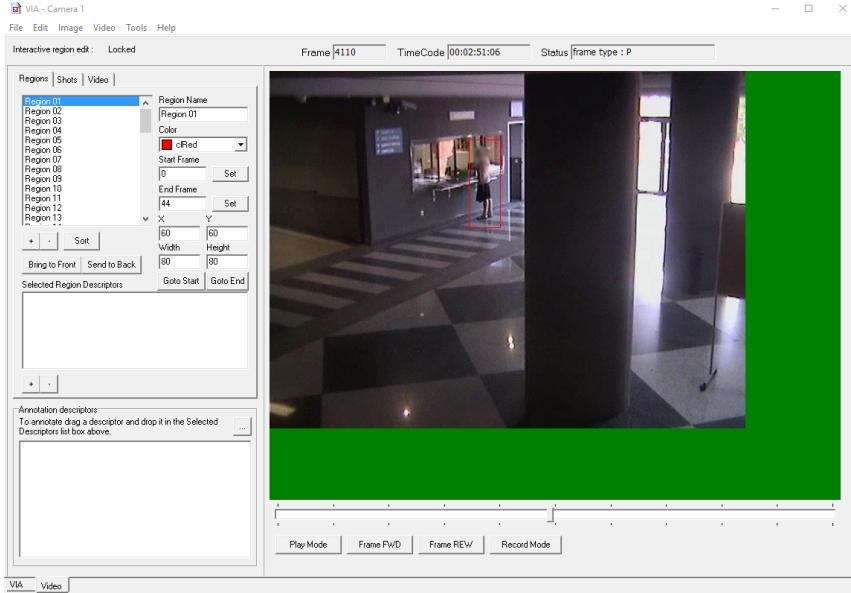


Figure 5.6: Via Annotation Tool software main window

5.2.2 Ground Truth Generation

We have selected one of the three recordings from the data-set to completely evaluate the system.

In order to do the evaluation ground truth information is needed. The video has been manually annotated with [Via Annotation Tool](#)¹ to generate real pedestrian bounding boxes. This software allows the user to generate XML files with bounding boxes positions through the whole sequences. The recording consists on three different videos of 8300 frames, which means that the total number of manually annotated frames has been 25.200.

In Figure 5.7 annotated frames for the recording in different video situations can be observed.

As one can observe in Figure 5.7. The selected sequence for evaluation purposes is a mixture between easy and complex situations for a pedestrian detector.

Situations go from people walking alone across the hall, to people pushing objects like a wheelchair, in addition to big groups of people both inside and outside the building. In general due to image quality in terms of resolution and illumination and pedestrian situations the selected sequence is in the medium-high range of difficulty for pedestrian detection.

¹<https://sourceforge.net/projects/via-tool/>

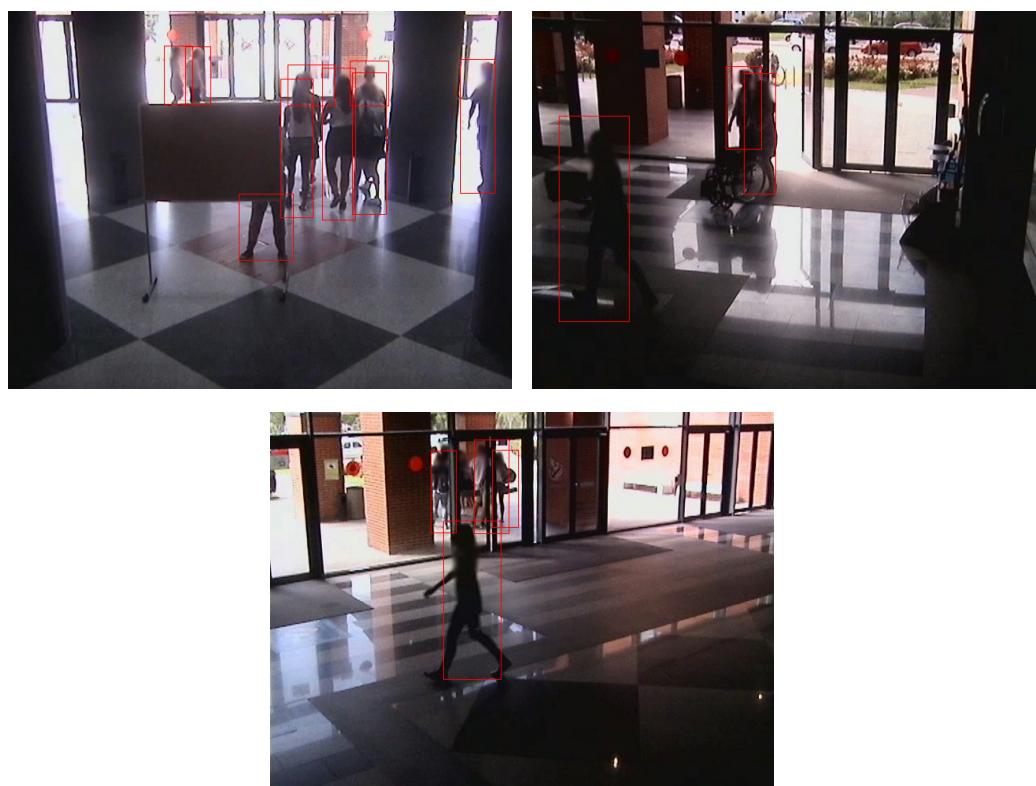


Figure 5.7: Annotated ground truth frames

5.2.3 Evaluation Framework

Once the data-set recording has been correctly annotated one can proceed to evaluate the performance of pedestrian detectors. In order to evaluate these algorithms precision and recall metrics are used. The calculation of these metrics is defined in Eq. 5.1 and Eq. 5.2, respectively.

$$\text{Precision} = \frac{\#\text{True Positives}}{\#\text{True Positives} + \#\text{False Positives}} \quad (5.1)$$

$$\text{Recall} = \frac{\#\text{True Positives}}{\#\text{True Positives} + \#\text{False Negatives}} \quad (5.2)$$

Precision gives us information about the classification stage, whereas recall provides overall system performance information.

In addition, to measure performance in a more general process F-Score is used. It combines both Precision and Recall in a unique metric measure. Eq 5.3 express the general metric expression.

$$F - \text{Score} = \frac{(2 \cdot \text{Precision} \cdot \text{Recall})}{(\text{Precision} + \text{Recall})} \quad (5.3)$$

To extract the mentioned metrics we should have a method to compute which bounding boxes match with the ground truth. In order to achieve this objective Intersection over Union (IoU) or Jaccard metric has been selected. IoU is computed as in Eq 5.4 and it measures the relation between bounding boxes overlapping and its union areas.

This means that even if the bounding boxes overlap perfectly, but the union area of them is high, IoU metric will have a small value. This effect can be observed in Figure 5.8 where one can see that a almost perfect overlapping between bounding boxes does not lead to a high IoU value.

For our evaluation framework a detection is consider correct if $\text{IoU}(\text{Detection}, \text{Ground Truth}) \leq 0.5$ which is the usually metric value selected in the State of the Art.

$$\text{IoU}(\text{Detection}, \text{Ground Truth}) = \frac{\text{OverlapArea}(\text{Detection}, \text{Ground Truth})}{\text{UnionArea}(\text{Detection}, \text{Ground Truth})} \quad (5.4)$$

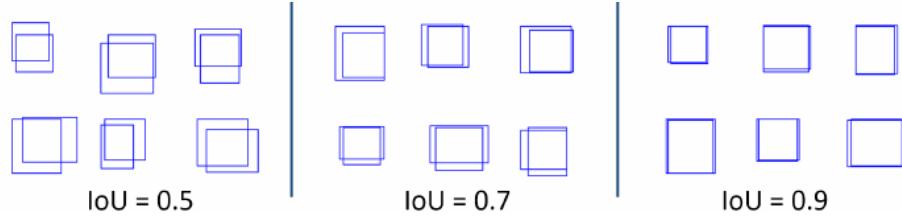


Figure 5.8: Intersection over Union or Jaccard Index

5.3 Application Performance

Along this section evaluation for the application in terms of general computational time is done.

Firstly, results concerning differences between single thread and multi thread applications are exposed. Measures have been obtained by using the same sequence and the same PD. In Table 5.2 one can observe the difference in term of seconds per frame and frame per second for the same video sequence and the same pedestrian detector approach (HOG).

	Seconds/Frame	FPS
Single-Thread	3.8 s	0.26 fps
Multi-Thread	0.7 s	1.42 fps

Table 5.2: Comparison between the use of single thread or multi threads

In addition, all the PD are compared when evaluated in the multi thread environment with the same video sequence. This results can be seen in Table 5.3

	Seconds/Frame	FPS
HOG	0.737 s	1.35 fps
ACF	1.642 s	0.60 fps
DPM	1.062 s	0.94 fps
PSP-Net (Offline)	0.624 s	1.60 fps

Table 5.3: Comparison between the use of single thread or multi threads

5.3.1 Application Performance Results Discussion

Looking at Table 5.2, one can easily derive that the speed up of using multithreads is clear. This difference implies that in the amount of time the single thread application obtains results for three frames, the multi-thread application has been executed almost 5.5 times. This is a huge improvement in the computational time.

In terms of pedestrian detection comparison, taking a look at the table one can easily observe that PSP-Net is the one which runs faster. However, its offline computation makes it evidently the fastest. Comparing those which are truly computed online, HOG is the one that performs faster followed by DPM and ACF.

Nevertheless, it is noticeable that even the slowest PD multi thread approach, performs faster than HOG in the single-thread application.

5.4 Homography Calculation

First results concerning with the proposed system deal with homographies. Correctly calculated homographies are essential for the right performance of further algorithms. In order to test that homographies have been correctly done by the manual selection of points, views have been projected with its correspondent matrix. In addition, selected points from both the frame (blue points) and the cenital plane (red points) have been overlapped in the same image. Some examples can be observed in Figure 5.9.

5.4.1 Homography Results Discussion

As one can observe in the mentioned Figure, homographies are correctly calculated in most of the cases. Blue points are not visible because of the overlapping with the red ones, which means that their correspondence is perfect. However, some small errors can be noticed in Figure 5.9a where there exists a displacement between some red points and the blue ones (top of the image).

It is important to notice in this results, how precision and image quality is reduced the further the points are from the camera. This is a standard problem when dealing with such type of homographies.

5.5 Semantic Segmentation

In this Section some results concerning the performance of PSP-Net in the proposed environment are presented. In order to test the algorithm and for the lack of semantic

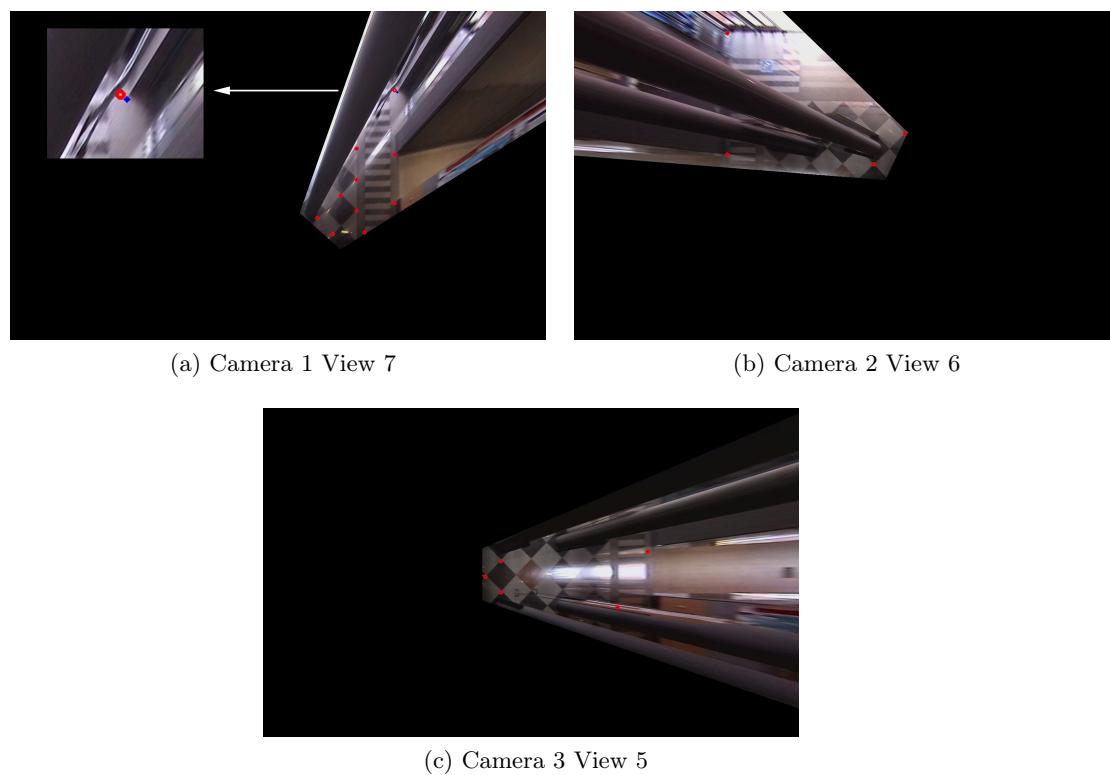


Figure 5.9: Projected views with overlapping selected points.

ground truth results should be visually evaluated. For this purpose Figure 5.10 is used.

5.5.1 Semantic Segmentation Results Discussion

Taking a look to results one can derive that PSP-Net extracts high quality results in terms of semantic. It is noticeable that most of the floor, if it is not occluded as in the top image from Figure 5.10, is correctly extracted.

However, it is important to point that segmentation is not perfect when dealing with difficult illumination cases or complex objects. It can be observed that columns are not correctly segmented. Some of them are misclassified as walls or floor. This is one of the main issues for our system as an accurate floor detection is needed. In addition, crystal doors such as the ones displayed in the top and bottom images, are not correctly distinguish from building and wall labels which could also be a main issue for further statistical data extraction algorithms.

5.6 Pedestrian Detection

During this section results in the context of pedestrian detection are presented. For this results HOG, DPM and PSP-Net algorithms have been selected to test. We propose a set of four different experiments to measure performance of our developed system.

The first one, tests pedestrian detections working in a mono-camera environment. This means that neither information is shared by the cameras nor semantic constrains are applied.

We follow testing the algorithms in a multi-camera setup. As explained in 3.3.2.3 detections from cameras are projected into the others to try to achieve better results than in a mono-camera scenario.

Thirdly, we analyze people detection again in a mono camera environment but this time, applying semantic constraining as proposed in Section 3.3.2.4.

Finally, the last test embraces all the proposed algorithms working within the multi-camera reprojection setup and applying semantic constraints.

It is important to remind that results are presented within Recall versus Precision curves. In order to create these graphs algorithms have been evaluated in all the four test using five different values from the score thresholding [0, 1] interval. Each threshold iteration (0, 0.2, 0.4, 0.6, 0.8) for an algorithm in a test, provides one point of the curve.

To proceed with this process all the algorithms scores have been normalized from its original values to the mentioned interval. Maximum and minimum scores have been

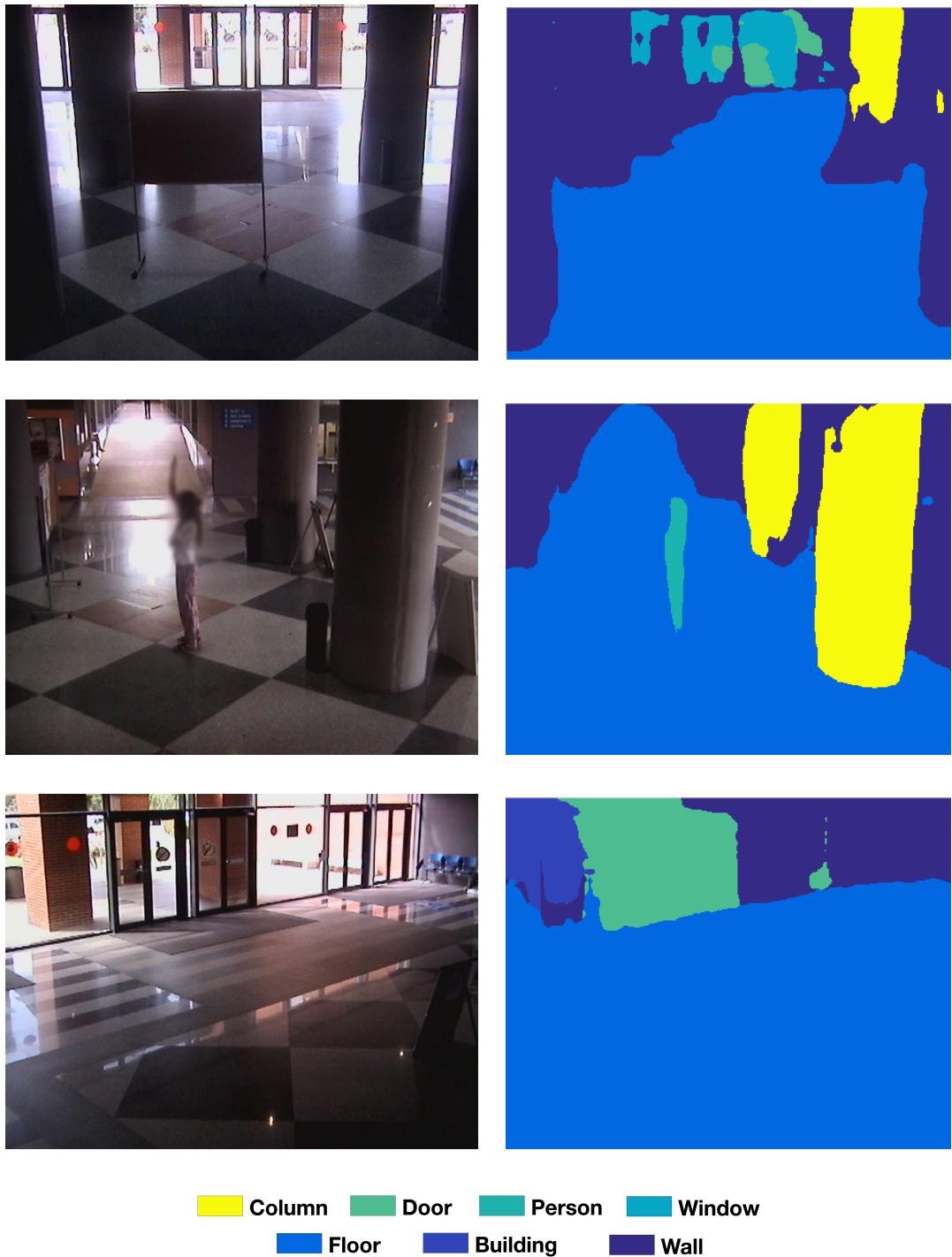


Figure 5.10: RGB frames and PSP-Net Semantic segmentation results. From top to bottom: Camera 1, Camera 2, Camera 3.

used for this purpose. In order to obtain scores distributions for all the detectors a train sequence has been used.

5.6.1 Mono-Camera

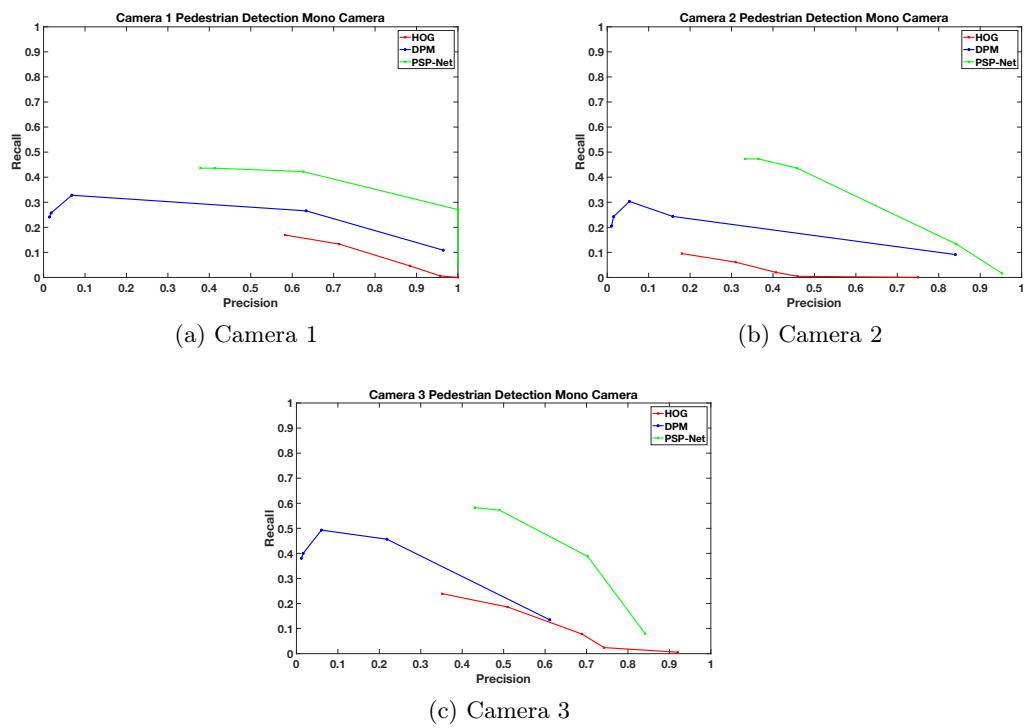


Figure 5.11: Recall / Precision graphs for mono camera pedestrian detection.

5.6.2 Multi-Camera

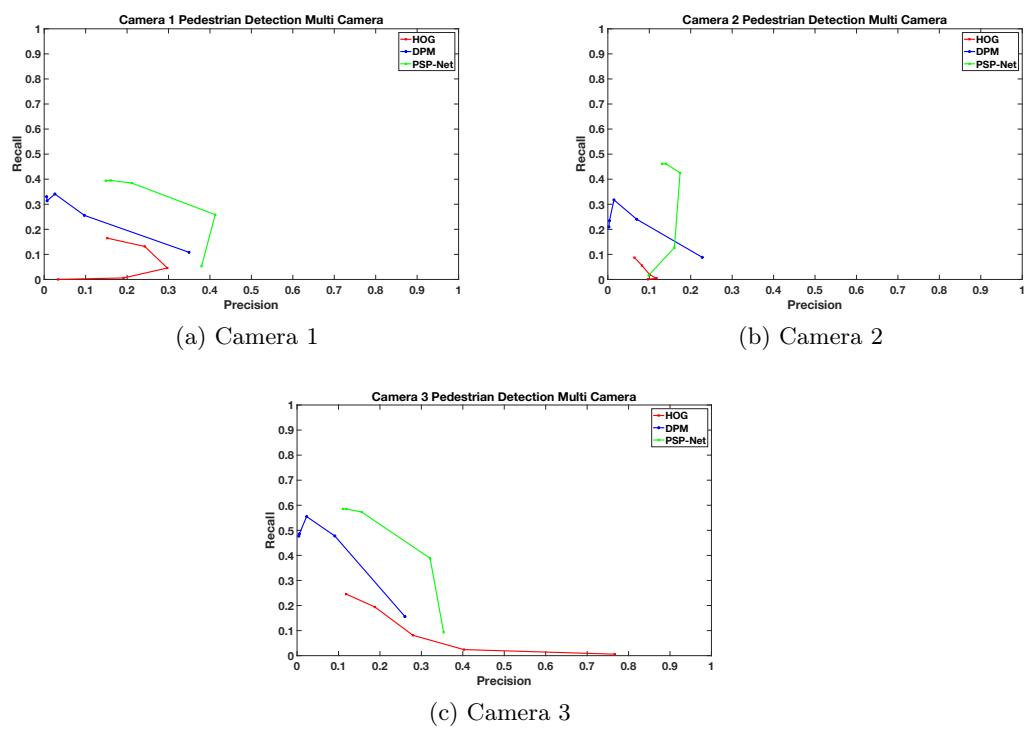


Figure 5.12: Recall / Precision graphs for multi camera pedestrian detection.

5.6.3 Mono-Camera with Semantic Constraining

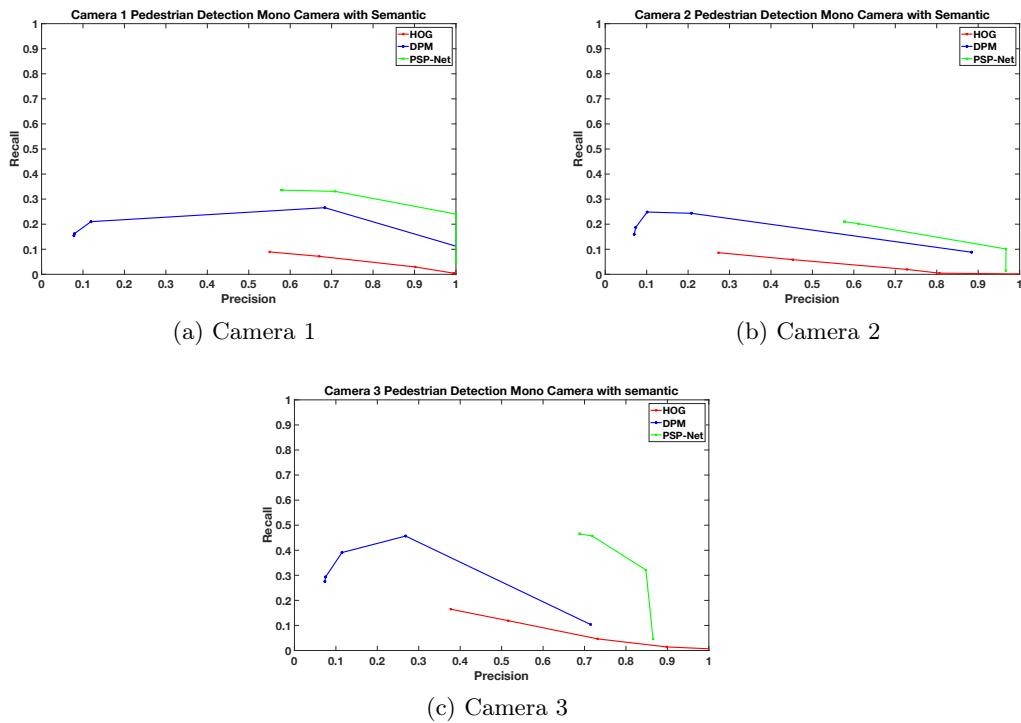


Figure 5.13: Recall / Precision graphs for mono camera pedestrian detection with semantic constraining.

5.6.4 Multi-Camera with Semantic Constraining

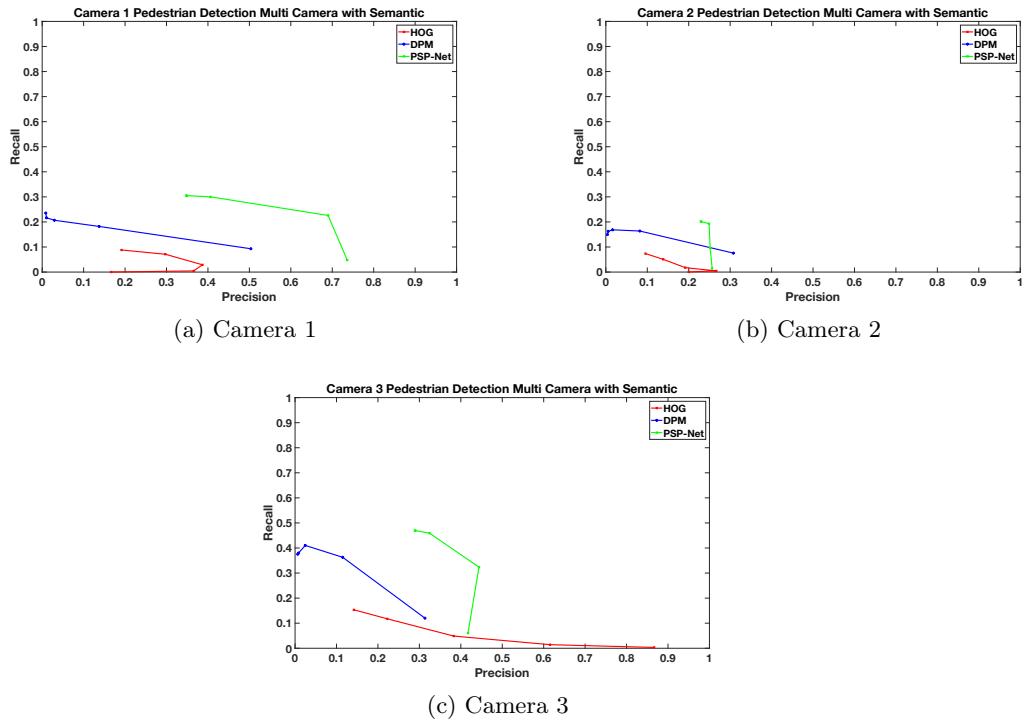


Figure 5.14: Recall / Precision graphs for multi camera pedestrian detection with semantic constraining.

In order to compare numerically all the experiments with the PD approaches F-Score results have been obtained and presented in Table 5.4. For this Table the best F-Score metric for a PD in a particular experiment has been chosen.

Pedestrian Detector		Mono Camera	Multi Camera	Mono Camera Semantic	Multi Camera Semantic
HOG	Camera 1	0.26	0.17	0.1537	0.12
	Camera 2	0.12	0.07	0.14	0.08
	Camera 3	0.28	0.19	0.2294	0.15
DPM	Camera 1	0.37	0.16	0.39	0.16
	Camera 2	0.19	0.12	0.23	0.13
	Camera 3	0.29	0.19	0.34	0.17
PSP-Net	Camera 1	0.50	0.31	0.45	0.34
	Camera 2	0.44	0.24	0.30	0.22
	Camera 3	0.52	0.35	0.56	0.38

Table 5.4: F-Score comparison between all the experiments and pedestrian detector approaches. Best result for an approach in a camera marked as bold. Best result of the table marked as red bold.

5.6.5 Pedestrian Detection Results Discussion

General PD Performance

If one analyzes Figures 5.11, 5.12, 5.13 and 5.14 it is easily observable that PSP-Net outperforms DPM and HOG in every experiment and in all the cameras. This is quite logical as algorithms based on CNNs usually perform better than classical PD approaches.

However, none of the detectors are getting excellent results in the testing sequence. Ideal results should lead to high recall and low precision values when low threshold are chosen and on the other hand, low recall and high precision values when the highest thresholds are used.

Semantic Constraining Performance

If one focus the analysis on the system general performance, taking into account both PD fusion and constraining we have to center the attention to Table 5.4. Here algorithms are compared in terms of F-Score which represents a global performance measure.

If we take a look at HOG descriptor two out of the three best results are obtained using raw HOG. Only Camera 2 with semantic constrain seems to outperform just by a small F-Score the raw detector. PSP-Net presents the same effect where only camera 3 gets a better result using semantic constraining. However, DPM using semantic filtering outperforms in every camera the scores obtained with raw DPM. This means that the inclusion of semantic information leads to an increase in DPM performance.

The reason why HOG and PSP-Net are not improved could be due to the number of extracted detections. Both approaches, although their scores have been normalized to the $[0, 1]$ interval, extract good results even with the lowest threshold. This is totally contrary to DPM detector, in which, using $Th = 0$ a high number of detections are obtained. (See Figure 5.15 for visual results).

This means, that if detections are already accurate and due to our common semantic maps, they lay on some of the floor holes, detections are suppressed. This leads to obtaining a worse performance than using the raw detector. In other words, false positives laying in areas different than floor should be higher than those wrongly suppressed detections if semantic constraining is to be considered worthwhile.

In addition, as depicted in Figure 5.7 ground truth has been made so every pedestrian is annotated in the sequence. This is a problem when using semantic constraining as every detection placed at the outside is suppressed for not being inside the hall. Furthermore, PSP-Net as is depicted in Figure 5.16 does not obtain pedestrian detections when analyzing people outside the building which could be one of the reasons of its poor performance.

Multi Camera Reprojection System Performance

If we now focus on the performance of the multi camera system when dealing with reprojections of pedestrian, Table 5.4 shows that it is not properly working. In both experiments in which reprojection is used (second and fourth experiments) worse results than even the regular detectors are obtained. This could be due to two reasons:

1. The projection from one frame to another has some error. This can be because of the homography calculation or pedestrian detection errors.

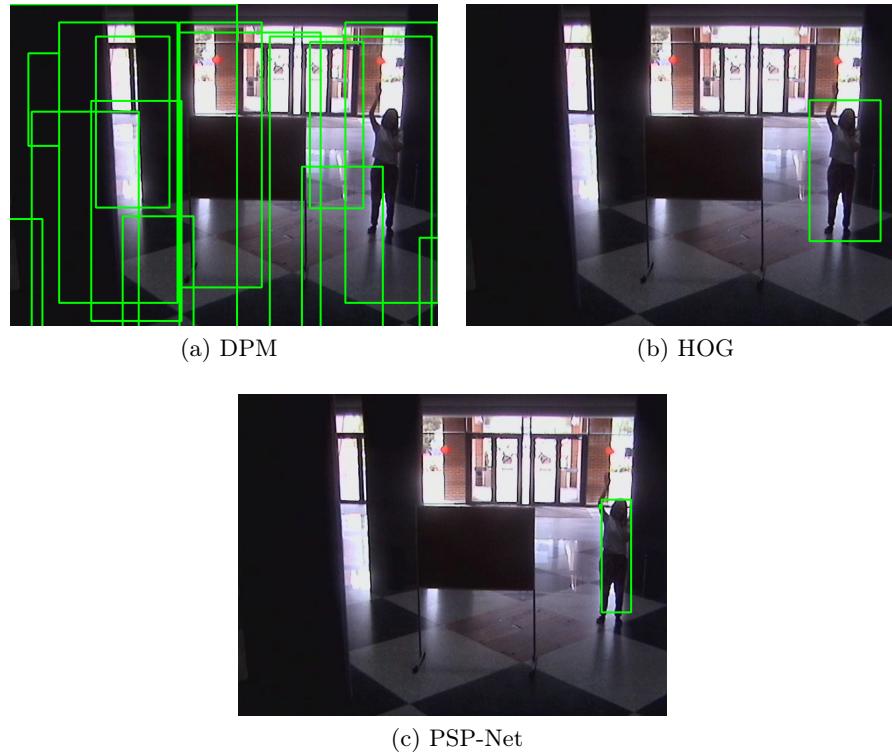


Figure 5.15: Differences in number of detections for $Th = 0$ and same t .

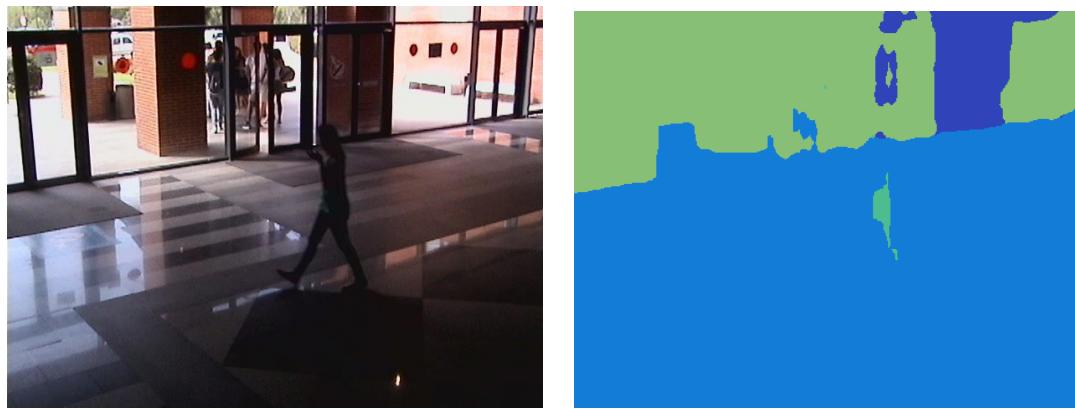


Figure 5.16: Misclassification of people in semantic segmentation. Camera 3.



(a) Camera 1. Red circle corresponds to ground truth and green one corresponds to the actual cylinder estimation.
 (b) Camera 2. White circle represents reprojection from Camera 1.

Figure 5.17: Pedestrian detection error leads to a reprojec-tion displacement.

In the first case if one homography is not accurate enough projections have small errors that translated to distances in the cenital plane could be a high displacement.

The second case is related to pedestrian detection errors. If a person is accurately detected but, the bounding box does not surround entirely the person, the reprojection is not accurate. This effect can be observed in Figure 5.17, in which the detection is correct but feet are not in the bounding box. This leads to a spatial difference between the cylinder center estimation and the actual person (see Figure 5.17b).

2. When reprojecting detections, blob height is lost and further reconstructed with respect to an aspect ratio. This process could lead to examples as the ones in Figure 5.18. Projections are finally reconstructed as small or big blobs if blob width is small or big respectively. This means that reprojections are not correctly combined with the original detection from the camera and they produce a false positive.

5.7 Statistical Usage Data

Finally, results concerning statistical usage data are presented. One graphic per selected semantic area is obtained. In our case and due to the mentioned limitations of the map in Figure 3.22 floor and doors semantic areas have been chosen. Curves represent pedestrian density per frame in the recording. The obtained curves are



Figure 5.18: Pedestrian reprojection error due to height.

displayed in Figures 5.19 and 5.20.

In addition, some results about the most used areas of the hall have been extracted. All the scene has been divided into regular sub regions of fixed size in which pedestrian flow has been measured. By using this method one could have information about the most used paths by people over the scene.

Results concerning this process are displayed in Figure 5.21 for a subregion size of 40×40 pixels and in Figure 5.22 for a subregion size of 70×70 pixels.

Yellow squares represent the amount of pedestrians found in that sub region. The more strong the color appears in the image it means that pedestrian density over that area has been higher than in the neighborhood. A temporal comparison has been made with the correspondent frame.

5.7.1 Statistical Usage Data Results Discussion

One can observe in Figures 5.19 and 5.20 that both semantic areas are measured in terms of number of people by frame. If one takes a look to the correspondent frame, displayed at the top of the graph, it can be observed that there is a correspondence between periods when most people are walking through the scene and the different peaks from the curve. It is important to notice, that also in Figure 5.19 at $t_1 = 1594$ and $t_2 = 4877$ we are not getting an accurate measure due to the small blobs detected, and also due to the group detection as a whole.

When talking about the doors usage one can observe that there is people passing through doors are in frames $t_1 = 1541$, $t_2 = 2808$, $t_3 = 3569$ and $t_4 = 5091$ approximately.

In both Figures 5.21 and 5.22 it is observable that most used paths increase accordingly pedestrian flow in the sequence. For that reason when a person enters

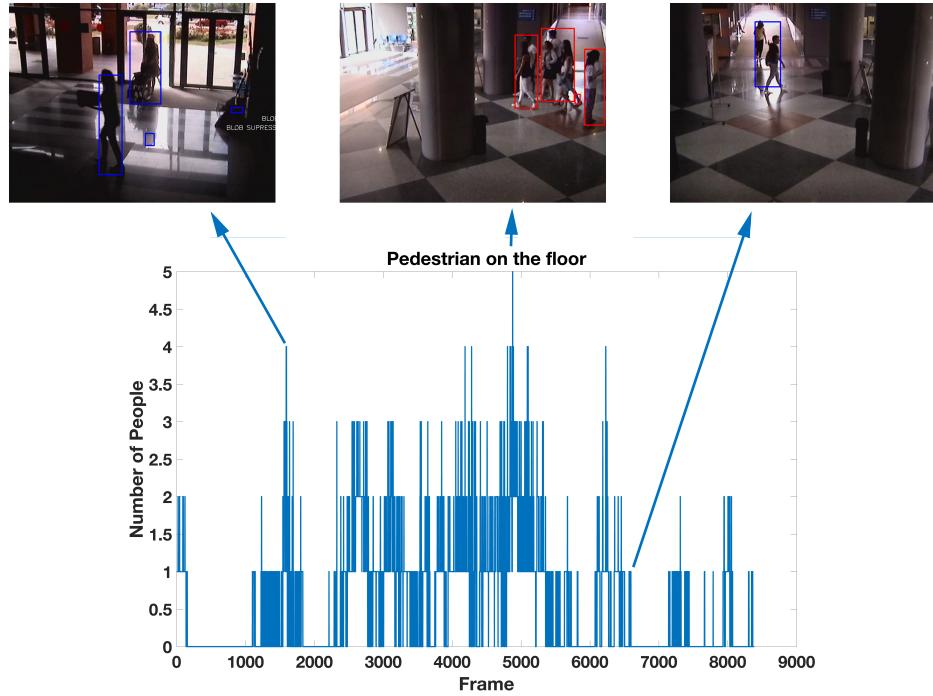


Figure 5.19: Floor usage graph. From left to right in camera frames: $t_1 = 1594$, $t_2 = 4877$ and $t_3 = 6228$.

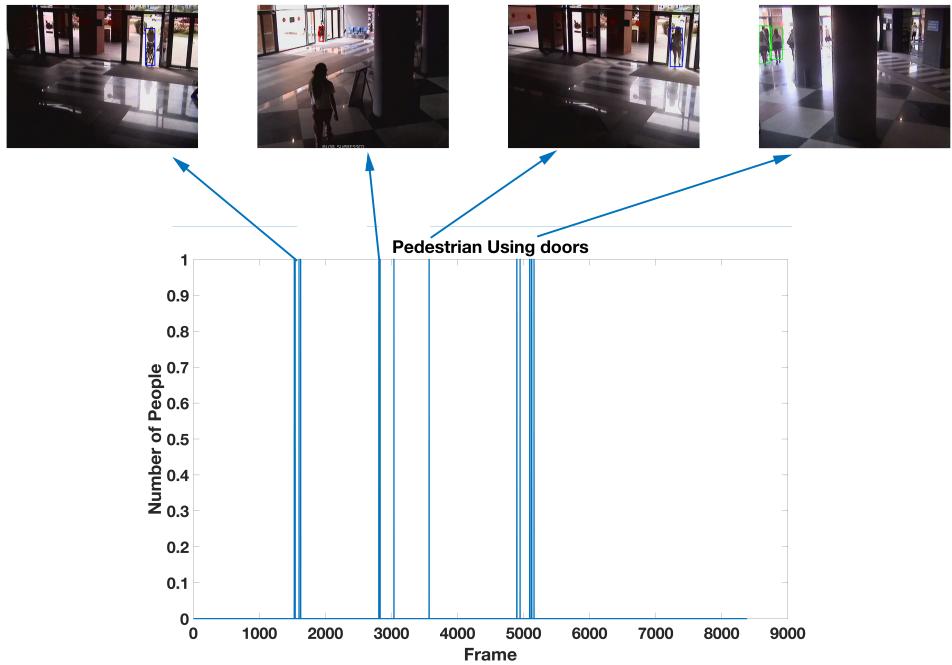


Figure 5.20: Doors usage graph. From left to right in camera frames: $t_1 = 1541$, $t_2 = 2808$, $t_3 = 3569$ and $t_4 = 5091$.

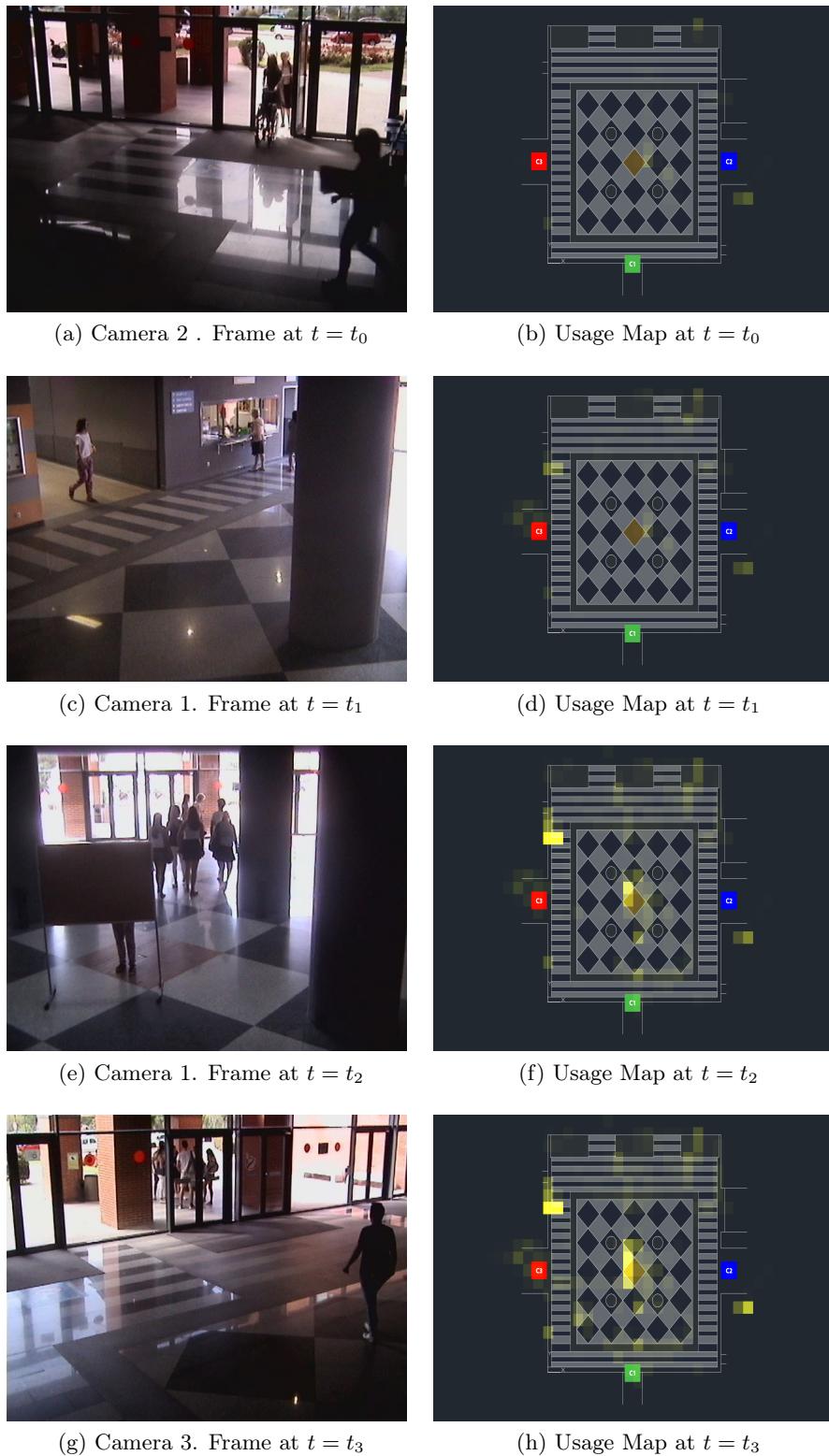


Figure 5.21: Pedestrians paths usage along the Hall. Scene has been divided in sub regions of 40 pixels.

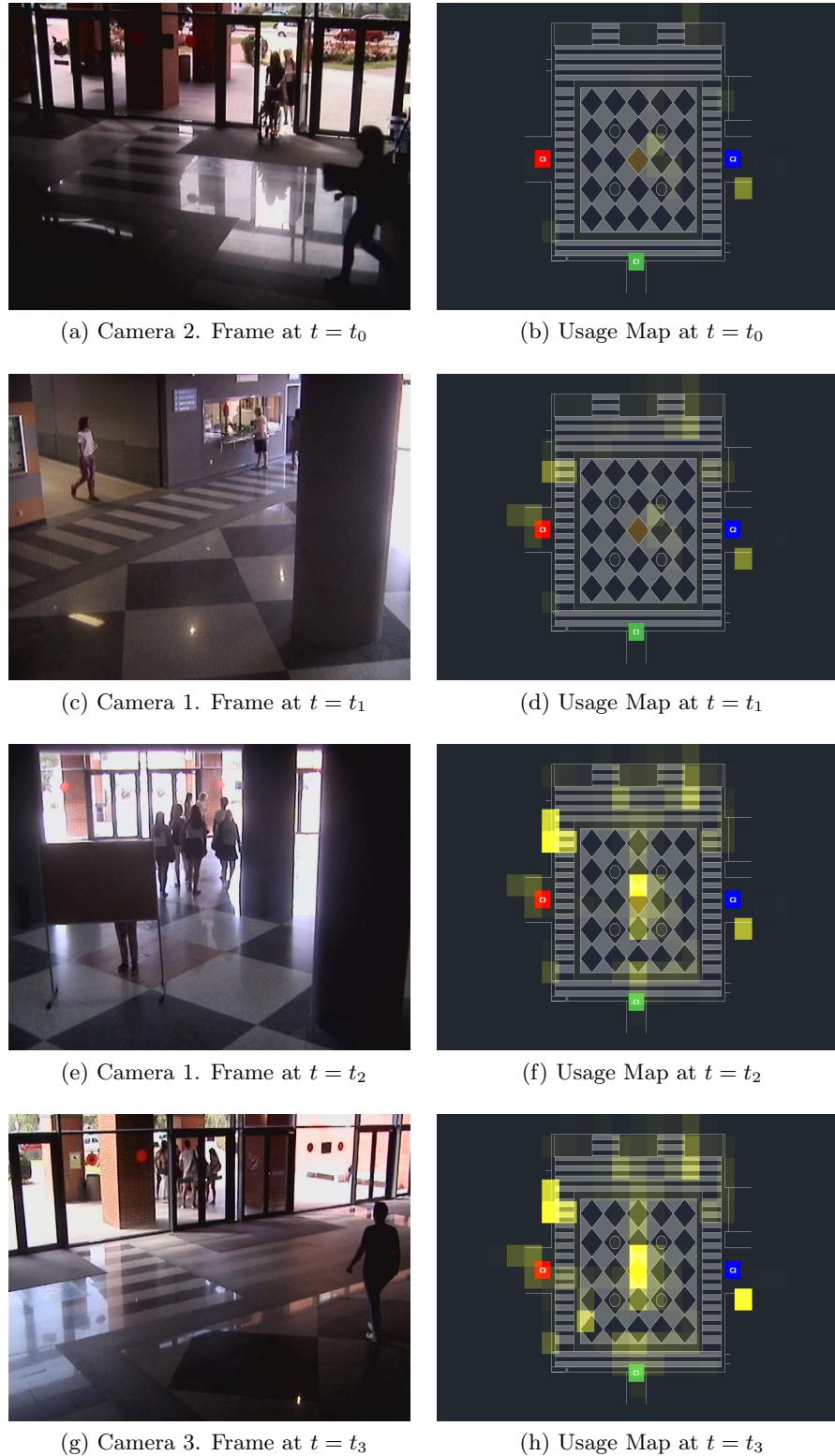


Figure 5.22: Pedestrians paths usage along the Hall. Scene has been divided in sub regions of 70 pixels.

the building areas surrounding the door get highlighted. In the same manner, when a big pedestrian group walks along the scene all their trajectory is being highlighted as pedestrian density over that areas is highly increased.

5.8 Overall Discussion

Taking into account all the presented results one can observe that some proposed algorithms from the system are performing correctly while others need to be improved.

Homography calculation is accurately enough for the purpose of the system. We observed that almost all the homographies were correctly computed and the small errors do not lead to further problems in algorithms.

Semantic segmentation using PSP-Net has some main drawbacks when dealing with columns or doors segmentation. In the first case, when misclassifying columns, sometimes, leads to a false floor label. This means that finally, when creating the semantic reference plane π_{ref} some floor areas are really columns. Doors errors in some specific frames could lead to not obtaining door areas, however, when combining all the sequence frames by the temporal median filter this problem is solved as depicted in Figure 3.22.

Moving to performance in terms of pedestrian detection above results show that multi cameras reprojection is not working properly in any of the tested approaches. However, semantic constraining performs promising when dealing with algorithms that need some filtering due to the high amount of extracted blobs.

Finally, semantic usage extraction performs good although errors from PD are trespassing to these results and so, they are influence by PD.

Chapter 6

Conclusions and Future Work

6.1 Conclusions

This master thesis has described a system capable of performing pedestrian detection and semantic segmentation over a multi camera setup at the Escuela Politecnica Superior, Universidad Autonoma de Madrid (Spain). Different pedestrian detection approaches such as HOG, DPM, ACF or Fast-RCNN and a semantic segmentation algorithm, PSP-Net, have been selected among others from a complete study of the state of the art.

A complete system to perform pedestrian fusion and filtering has been proposed. Information from three cameras are combined into a common cenital plane by the use of homographies. This process allows PD to be shared and reprojected from one camera to another. In addition, common multi camera semantic information has been used to constraint PD detections.

Statistical data usage from different semantic areas has been extracted. Also, pedestrian paths in terms of density have been computed based on some semantic information such as floor paths and doors.

In order to control and tune algorithms a multi thread application has been developed under QT Developing environment. It has a Graphical User Interface that sets the base for a user-friendly interaction between the user and the software. This application also allows the user to represent and visualize all the extracted results.

All the system has been tested in a manually annotated recording obtaining performance results for homography calculation, PD approaches under different environments and semantic areas statistical usage data. During these tests semantic constraining has improved DPM approach, however algorithms such as HOG or PSP-Net have obtained worse results than the raw detector. Similarly, multi-camera fusion has

not been able to increase performance for any of the algorithms. Problems related to pedestrian reprojection between cameras combined with PD errors lead to an increase in the number of false positive and so, in performance decreased.

Statistical data usage has been almost perfectly extracted, having exact frames in which semantic areas are used. In addition, semantic paths have been obtained with high precision in spatial and temporally terms.

6.2 Future Work

Considering current state of the art, obtained results and extracted conclusions one can set the stage for future work.

In terms of application and software development some improvements are proposed to be done. Nowadays, heavy computational work is achieved almost in real time by the use of graphical cards. GPU computation should be implemented in the scope of this work. Many of the used methods are also implemented with GPU computation functionalities and the inclusion of this kind of speed-up could lead to better and faster results.

In other term, the proposed system bases all its performance in the correct view selection as explained in Section 3.3.1.5. However, when high illumination changes occur this process fails. We propose to fix this problem by the use of camera spatial positions. With this information exact same position for each of the views can be obtained periodically and so, they can be updated during the video sequence in order to adapt to illumination changes.

One of the main problems discussed in the results section is that when pedestrians are reprojected the blob height is lost. Due to this, there are some frames that have small detections compared to the person size. To correctly reproject the blob we propose as future work to use real distances between camera and pedestrian detection to finally obtain an approximation to the real height for the blob.

Appendix A

Cenital Plane Design

One of the main objectives of the work is to project extracted detections from the three cameras, i.e. pedestrian and semantic, into a common plane for all of them. To achieve this goal a cenital plane that correctly represents the scene is needed.

First Approach

The first used approach to represent the cenital plane is depict in Figure A.1.



Figure A.1: First cenital plane approach

This cenital plane lacks of details from the scene. The information about the details of the scenario is minimum and also, the scene proportions are not correct. To compute a correct homography between the camera frame and the cenital plane one should be able to identify the same scene points in both images in the ground plane.

This means that the cenital plane should have enough details so the point selection is done correctly by the user and the homography is correctly computed.

Second Approach

For this reason, and driven by bad results in terms of projections, another cenital plane has been computed starting from zero. In this new approach the scene has correctly been measured by hand and the plane has been done with real measures and high floor detail.

For correctly drawing the plane [AutoCAD 2017](#) software has been used. The second plane approach with all the manual measures extracted from the real scene can be seen in Figure A.2. Figure A.3 represents the final cenital map with the correct camera positions.

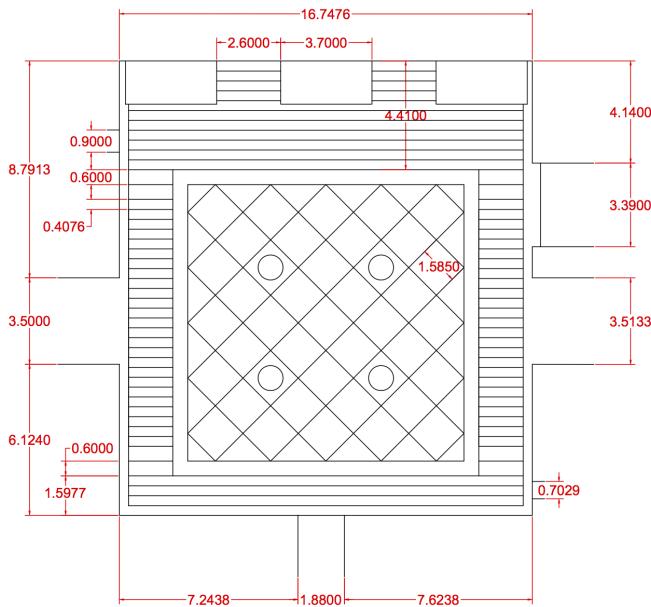


Figure A.2: Second cenital plane approach with real measures

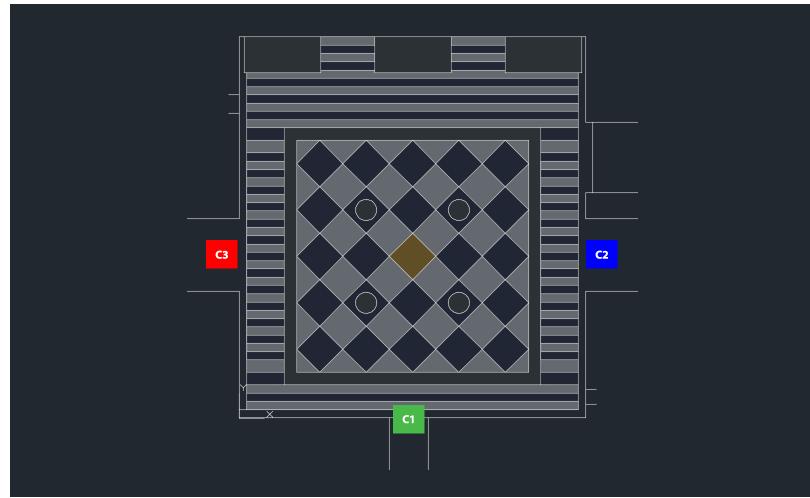


Figure A.3: Second cenital plane approach with camera positions

It is easily observable that differences between Figure A.1 and A.3 are outstanding both in floor details and in general construction proportions. This detail rise leads to a much more easy homography selection points by the user. This is due to the high amount of point options to choose from the new cenital map. Evidentially this means that the final homography matrix, and so, all the projections computed by it, have better accuracy.

Appendix B

AKAZE

AKAZE detector and descriptor [57] is a fast multi scale feature detection and description approach. It exploits the benefits of nonlinear scale spaces.

Previous approaches such as KAZE [58] or BFSIFT [59] have a main time consumption drawback in terms of nonlinear scale space creation.

Nevertheless, AKAZE uses recent numerical schemes called Fast Explicit Diffusion (FED) [60, 61] in order to build any kind of discretization scheme in a much more faster speed. These FED schemes are embedded in a pyramidal framework in order to achieve the speedup in terms of features detector.

In addition the use of the Modified-Local Difference Binary (M-LDB) descriptor which is described as highly efficient. It exploits gradient and intensity information from the nonlinear scale space. In addition, M-LDB is both scale and rotation invariant.

Appendix C

Parametric Homographies Between Inertial Planes

Homography matrices, as explained during the Thesis, aim to relate the floor plane present in a frame with the cenital view. However, everything that is not exactly in the same plane as the floor is not projected properly when the homography matrix is used. In our work, different semantic are needed to be projected, for instance a door. When using floor homography only its base is correctly projected, whereas the rest of it is disfigured.

In [62] a solution to this problem is proposed. The general idea is to create a multilayer reconstruction. Once the homography matrix ${}^{\pi_{ref-C}}H_{view}$ that relates the image view with the reference frame π_{ref-C} is calculated, one can obtain another matrix ${}^{\pi'}H_{view}$ that relates the same image frame with a parallel plane called inertial plane at a fixed height Δt .

${}^{\pi'}H_{view}$ can be expressed as a function of ${}^{\pi_{ref-C}}H_{view}$ and Δh as described in Eq C.1.

$${}^{\pi'}H_v^{-1}(\Delta h) = {}^{\pi}H_v^{-1} + \Delta h P \hat{\mathbf{k}}^T, \quad (\text{C.1})$$

where $P = [u_0 \ v_0 \ 1]^T$ is the principal point of camera C and $\hat{\mathbf{k}}$ is the unit vector of the Z axis.

All this process is described in Figure C.1.

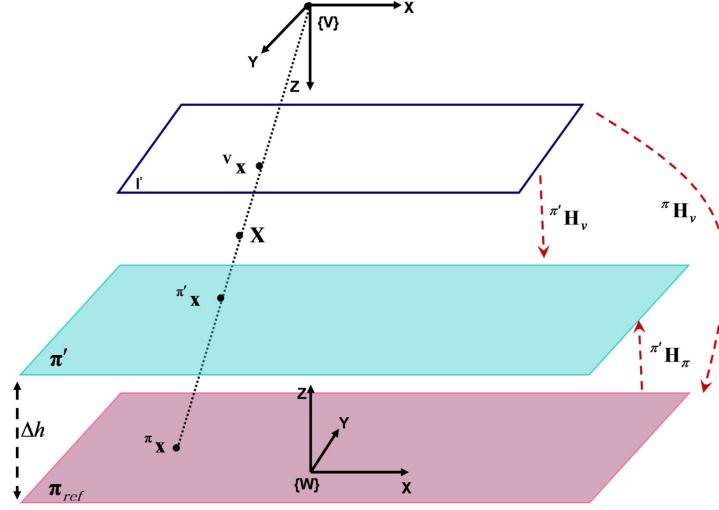


Figure C.1: Extending homography for planes parallel to π_{ref} . $\pi_{ref} H_V$ is the available homography between camera view and reference plane π_{ref} .

By this process, ideally a number k of planes could be generated (Figure C.2) in which different object sections are correctly projected. It could lead to a complete semantic map in which all the pixels represent semantic areas that have been correctly projected.

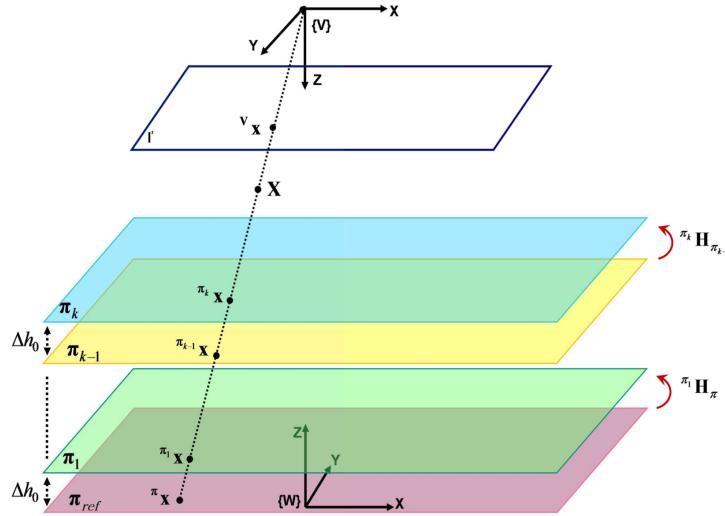


Figure C.2: Set of k inertial planes π_k . Each inertial plane is separated from the other by the same Δh height

Bibliography

- [1] A. K. Jain, L. Hong, and Y. Kulkarni, “A multimodal biometric system using fingerprint, face and speech,” in *2nd Int’l Conf. AVBPA*, vol. 10, 1999.
- [2] X. Wang, “Intelligent multi-camera video surveillance: A review,” *Pattern recognition letters*, vol. 34, no. 1, pp. 3–19, 2013.
- [3] P. Dollar, C. Wojek, B. Schiele, and P. Perona, “Pedestrian detection: An evaluation of the state of the art,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 4, pp. 743–761, 2012.
- [4] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, “Object detection with discriminatively trained part-based models,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 32, no. 9, pp. 1627–1645, 2010.
- [5] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, pp. 91–99, 2015.
- [6] J. Hosang, R. Benenson, P. Dollár, and B. Schiele, “What makes for effective detection proposals?,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 4, pp. 814–830, 2016.
- [7] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, “Pyramid scene parsing network,” *CoRR*, 2016.
- [8] R. Mazzon and A. Cavallaro, “Multi-camera tracking using a multi-goal social force model,” *Neurocomputing*, vol. 100, pp. 41–50, 2013.
- [9] A. Turner and A. Penn, “Encoding natural movement as an agent-based system: an investigation into human pedestrian behaviour in the built environment,” *Environment and planning B: Planning and Design*, vol. 29, no. 4, pp. 473–490, 2002.

- [10] P. Scovanner and M. F. Tappen, “Learning pedestrian dynamics from the real world,” in *Computer Vision, 2009 IEEE 12th International Conference on*, pp. 381–388, IEEE, 2009.
- [11] F. Jiang, J. Yuan, S. A. Tsaftaris, and A. K. Katsaggelos, “Anomalous video event detection using spatiotemporal context,” *Computer Vision and Image Understanding*, vol. 115, no. 3, pp. 323–333, 2011.
- [12] P. Dollár, C. Wojek, B. Schiele, and P. Perona, “Pedestrian detection: A benchmark,” in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pp. 304–311, IEEE, 2009.
- [13] A. Ess, B. Leibe, and L. Van Gool, “Depth and appearance for mobile scene analysis,” in *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*, pp. 1–8, IEEE, 2007.
- [14] C. Wojek, S. Walk, and B. Schiele, “Multi-cue onboard pedestrian detection,” in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pp. 794–801, IEEE, 2009.
- [15] N. Dalal and B. Triggs, “Histograms of oriented gradients for human detection,” in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, vol. 1, pp. 886–893, IEEE, 2005.
- [16] P. Felzenszwalb, D. McAllester, and D. Ramanan, “A discriminatively trained, multiscale, deformable part model,” in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pp. 1–8, IEEE, 2008.
- [17] P. Dollár, R. Appel, S. Belongie, and P. Perona, “Fast feature pyramids for object detection,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 8, pp. 1532–1545, 2014.
- [18] Á. García-Martín and J. M. Martínez, “People detection in surveillance: classification and evaluation,” *IET Computer Vision*, vol. 9, no. 5, pp. 779–788, 2015.
- [19] R. Cutler and L. S. Davis, “Robust real-time periodic motion detection, analysis, and applications,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 781–796, 2000.
- [20] J. Giebel, D. Gavrila, and C. Schnörr, “A bayesian framework for multi-cue 3d object tracking,” *Computer Vision-ECCV 2004*, pp. 241–252, 2004.

- [21] K. Okuma, A. Taleghani, N. d. Freitas, J. J. Little, and D. G. Lowe, “A boosted particle filter: Multitarget detection and tracking,” *Computer Vision-ECCV 2004*, pp. 28–39, 2004.
- [22] A. García-Martín, B. Alcedo, and J. M. Martínez, “Pdbm: people detection benchmark repository,” *Electronics Letters*, vol. 51, no. 7, pp. 559–560, 2015.
- [23] K. E. Van de Sande, J. R. Uijlings, T. Gevers, and A. W. Smeulders, “Segmentation as selective search for object recognition,” in *Computer Vision (ICCV), 2011 IEEE International Conference on*, pp. 1879–1886, IEEE, 2011.
- [24] S. Manen, M. Guillaumin, and L. Van Gool, “Prime object proposals with randomized prim’s algorithm,” in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2536–2543, 2013.
- [25] P. Rantalankila, J. Kannala, and E. Rahtu, “Generating object segmentation proposals using global and local search,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2417–2424, 2014.
- [26] K.-Y. Chang, T.-L. Liu, H.-T. Chen, and S.-H. Lai, “Fusing generic objectness and visual saliency for salient object detection,” in *Computer Vision (ICCV), 2011 IEEE International Conference on*, pp. 914–921, IEEE, 2011.
- [27] J. Carreira and C. Sminchisescu, “Constrained parametric min-cuts for automatic object segmentation,” in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pp. 3241–3248, IEEE, 2010.
- [28] I. Endres and D. Hoiem, “Category independent object proposals,” *Computer Vision-ECCV 2010*, pp. 575–588, 2010.
- [29] A. Humayun, F. Li, and J. M. Rehg, “Rigor: Reusing inference in graph cuts for generating object regions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 336–343, 2014.
- [30] P. Krähenbühl and V. Koltun, “Geodesic object proposals,” in *European Conference on Computer Vision*, pp. 725–739, Springer, 2014.
- [31] P. Arbeláez, J. Pont-Tuset, J. T. Barron, F. Marques, and J. Malik, “Multiscale combinatorial grouping,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 328–335, 2014.

- [32] B. Alexe, T. Deselaers, and V. Ferrari, “What is an object?,” in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pp. 73–80, IEEE, 2010.
- [33] E. Rahtu, J. Kannala, and M. Blaschko, “Learning a category independent object detection cascade,” in *Computer Vision (ICCV), 2011 IEEE International Conference on*, pp. 1052–1059, IEEE, 2011.
- [34] M.-M. Cheng, Z. Zhang, W.-Y. Lin, and P. Torr, “Bing: Binarized normed gradients for objectness estimation at 300fps,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3286–3293, 2014.
- [35] C. L. Zitnick and P. Dollár, “Edge boxes: Locating object proposals from edges,” in *European Conference on Computer Vision*, pp. 391–405, Springer, 2014.
- [36] J. Feng, Y. Wei, L. Tao, C. Zhang, and J. Sun, “Salient object detection by composition,” in *Computer Vision (ICCV), 2011 IEEE International Conference on*, pp. 1028–1035, IEEE, 2011.
- [37] Z. Zhang, J. Warrell, and P. H. Torr, “Proposal generation for object detection using cascaded ranking svms,” in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pp. 1497–1504, IEEE, 2011.
- [38] M. Van den Bergh, G. Roig, X. Boix, S. Manen, and L. Van Gool, “Online video seeds for temporal window objectness,” in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 377–384, 2013.
- [39] J. Kim and K. Grauman, “Shape sharing for object segmentation,” *Computer Vision–ECCV 2012*, pp. 444–458, 2012.
- [40] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov, “Scalable object detection using deep neural networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2147–2154, 2014.
- [41] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- [42] Z. Cai, M. J. Saberian, and N. Vasconcelos, “Learning complexity-aware cascades for deep pedestrian detection,” *CoRR*, vol. abs/1507.05348, 2015.
- [43] R. Girshick, “Fast r-cnn,” in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1440–1448, 2015.

- [44] D. Navon, “Forest before trees: The precedence of global features in visual perception,” *Cognitive psychology*, vol. 9, no. 3, pp. 353–383, 1977.
- [45] R. A. Rensink, J. K. O’Regan, and J. J. Clark, “To see or not to see: The need for attention to perceive changes in scenes,” *Psychological science*, vol. 8, no. 5, pp. 368–373, 1997.
- [46] A. Torralba, A. Oliva, M. S. Castelhano, and J. M. Henderson, “Contextual guidance of eye movements and attention in real-world scenes: the role of global features in object search.,” *Psychological review*, vol. 113, no. 4, p. 766, 2006.
- [47] A. Torralba and P. Sinha, “Detecting faces in impoverished images,” tech. rep., Massachusetts Inst of Tech Cambridge Artificial Intelligence Lab, 2001.
- [48] P. Salembier and T. Sikora, *Introduction to MPEG-7: Multimedia Content Description Interface*. New York, NY, USA: John Wiley & Sons, Inc., 2002.
- [49] S. Oh, A. Hoogs, M. Turek, and R. Collins, “Content-based retrieval of functional objects in video using scene context,” in *European Conference on Computer Vision*, pp. 549–562, Springer, 2010.
- [50] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, “Semantic understanding of scenes through the ade20k dataset,” *CoRR*, vol. 1608, 2016.
- [51] 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,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3213–3223, 2016.
- [52] Z. Wu, C. Shen, and A. v. d. Hengel, “Wider or deeper: Revisiting the resnet model for visual recognition,” *CoRR*, 2016.
- [53] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. W. Cottrell, “Understanding convolution for semantic segmentation,” *CoRR*, vol. abs/1702.08502, 2017.
- [54] A. Miguélez and R. Nieto, “Detcción de personas en entornos multicámara utilizando informacion contextual,” Master’s thesis, Universidad Autónoma de Madrid. Escuela Politécnica Superior, 2016.
- [55] J. Xiao and L. Quan, “Multiple view semantic segmentation for street view images,” in *Computer Vision, 2009 IEEE 12th International Conference on*, pp. 686–693, IEEE, 2009.

- [56] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- [57] P. F. Alcantarilla and T. Solutions, “Fast explicit diffusion for accelerated features in nonlinear scale spaces,” *IEEE Trans. Patt. Anal. Mach. Intell.*, vol. 34, no. 7, pp. 1281–1298, 2011.
- [58] P. Alcantarilla, A. Bartoli, and A. Davison, “Kaze features,” *Computer Vision–ECCV 2012*, pp. 214–227, 2012.
- [59] S. Wang, H. You, and K. Fu, “Bfsift: A novel method to find feature matches for sar image registration,” *IEEE Geoscience and remote sensing letters*, vol. 9, no. 4, pp. 649–653, 2012.
- [60] J. Weickert, S. Grewenig, C. Schroers, and A. Bruhn, “Cyclic schemes for pde-based image analysis,” *International Journal of Computer Vision*, vol. 118, no. 3, pp. 275–299, 2016.
- [61] S. Grewenig, J. Weickert, and A. Bruhn, “From box filtering to fast explicit diffusion,” in *Joint Pattern Recognition Symposium*, pp. 533–542, Springer, 2010.
- [62] H. Aliakbarpour, V. S. Prasath, K. Palaniappan, G. Seetharaman, and J. Dias, “Heterogeneous multi-view information fusion: Review of 3-d reconstruction methods and a new registration with uncertainty modeling,” *IEEE Access*, vol. 4, pp. 8264–8285, 2016.