



Pós-Graduação em Ciência da Computação

Rafael Alves Roberto

Incremental Semantic Tracking on Mobile Devices



Universidade Federal de Pernambuco

posgraduacao@cin.ufpe.br

<http://cin.ufpe.br/~posgraduacao>

Recife
2018

Rafael Alves Roberto

Incremental Semantic Tracking on Mobile Devices

Ph.D. Thesis presented to the Graduation Program in Computer Science of the Informatics Center of Federal University of Pernambuco in partial fulfillment of the requirements for the degree of Philosophy Doctor in Computer Science

Advisor: Veronica Teichrieb

Co-Advisors: João Paulo Silva do Monte Lima
and Hideaki Uchiyama

Recife
2018

Catálogo na fonte
Bibliotecária Monick Raquel Silvestre da S. Portes, CRB4-1217

R642i Roberto, Rafael Alves
Incremental semantic tracking on mobile devices / Rafael Alves Roberto. –
2018.
121 f.: il., fig.

Orientadora: Veronica Teichrieb.
Tese (Doutorado) – Universidade Federal de Pernambuco. CIn, Ciência da
Computação, Recife, 2018.
Inclui referências.

1. Visão computacional. 2. Rastreamento semântico. I. Teichrieb, Veronica
(orientadora). II. Título.

006.37 CDD (23. ed.) UFPE- MEI 2018-103

Rafael Alves Roberto

“Incremental Semantic Tracking on Mobile Devices”

Tese de Doutorado apresentada ao Programa de Pós-Graduação em Ciência da Computação da Universidade Federal de Pernambuco, como requisito parcial para a obtenção do título de Doutor em Ciência da Computação.

Aprovado em: 27/06/2018.

Orientadora: Profa. Veronica Teichrieb

BANCA EXAMINADORA

Prof. Silvio de Barros Melo
Centro de Informática / UFPE

Prof. João Marcelo Xavier Natário Teixeira
Departamento de Eletrônica e Sistemas / UFPE

Prof. Douglas Guimarães Macharet
Departamento de Ciência da Computação / UFMG

Prof. Dieter Schmalstieg
Faculty of Computer Science and Biomedical Engineering /
Graz University of Technology

Prof. Takafumi Taketomi,
Graduate School of Information Science /
Nara Institute of Science and Technology

ACKNOWLEDGEMENTS

O doutorado é caminho longo, tanto que quando este aqui começou não existia nenhuma piada com os números 7 e 1 e a gente não sabia que o placar da FIFA tinha barra de rolagem. Falo com tranquilidade que eu não teria conseguido chegar aqui sem o suporte de muita gente. Primeiramente, meus pais, irmão e irmã, que me acompanham numa jornada ainda mais longa que essa. Vocês moram de pantufas no meu coração! Também à minha esposa, por estar presente durante os sacrifícios e caos absoluto e eterno que é fazer um doutorado. Por me acompanhar até o outro lado do mundo (literalmente). Você é minha guerreirinha! Aos meus orientadores por me darem oportunidade de fazer o meu nome, por servirem de inspiração e por todas as discussões e contribuições. Me vejo obrigado a concordar que isso foi fundamental para a qualidade final deste trabalho. Só conselho top! Não posso esquecer os companheiros e companheiras do Voxar Labs por todos os momentos de discussão, descontração e mão na brasa que me fizeram crescer bastante. Vocês tem talento pra isso! Por fim, mas não menos importante, aos amigos e amigas que, apesar do ritmo ragatanga insano, a gente conseguiu manter um convívio social honestíssimo. Vocês me ajudaram a não ficar mais louco que o padre baloeiro.

I would like to say a few words in English to my friends overseas. First, I would like to thank my advisors in Japan and Austria for having me there and for all the discussions. Your contributions were very important to this Ph.D. I also would like to thank people at LIMU and ICG for teaching me so much. And last but not least, I want to leave my gratitude to the people at ICG and Kroisbach for the good moments they helped us to feel at home when we were far away from our comfort zone.

ABSTRACT

Tracking is an important task that is used for several applications. The improvement and popularization of mobile devices in recent years allowed these applications to be executed on such devices, which provides a mobility that is not possible on desktop computers. However, there are still several challenges in this field. Thus, the goal of this Ph.D. is to investigate methods to perform tracking on mobile devices considering the characteristics of such platform. To achieve this goal, it was conducted a systematic mapping on tracking for mobile devices. This study collected 2,602 papers, from which 444 were selected to be classified. The results indicated a growing interest in this field and a preference for works that use the device's sensors to perform tracking locally on the device. The mapping was used to elaborate preliminary experimental scenarios. First, the Google Tango platform was evaluated to establish a ground base of the state-of-the-art trackers. It was observed that the precision in indoor spaces is suitable to provide a good user experience, including for augmented reality applications. Another experiment evaluated the use of parallelism, distributed approach and native implementation. This test showed that, on average, native development was the most efficient. Besides that, experiments were designed intending to test different tracking techniques. One is a face tracking technique using machine learning that was adapted to consider the characteristics of mobile devices and it runs in approximately eight milliseconds on such equipments. The other one is a SLAM technique that was developed in desktop and was ported to a Tango tablet device. There were several lessons learned from the experiments. One of them was the importance of finding high-level semantic information from a scene, which can improve tracking and provide more realistic rendering. In this Ph.D., it was developed a technique that incrementally detects and tracks primitives using the generating process of point clouds of visual SLAM systems, called Geometric and Statistical Incremental Semantic Tracker (GS-IST). The experiment indicates that GS-IST was able to improve both precision and stability of existing methods. However, since it focuses on precision, it compromises the recall to ensure the detection and tracking of correct shapes. In order to evaluate how GS-IST would perform running on mobile devices, it was ported to the Android platform. The evaluation showed that the mobile version is 8.5 to 9.9 times slower in comparison with the desktop implementation. Moreover, it uses up to 30.5% of the CPU load, which allows this implementation to run on a separate thread of the main tracking technique. Additionally, the energy consumption was not a concern because GS-IST can run for more than 4 hours in the worst case. Finally, the memory usage was less than 8% of the total RAM memory of the test devices, which did not have an impact on the execution time.

Key-words: Semantic tracking. Visual SLAM. Mobile devices. Tracking. Android.

RESUMO

Rastreamento é uma atividade importante usada em várias aplicações. Atualmente, estas aplicações podem ser executadas em dispositivos móveis graças à popularização e à melhoria desses aparelhos, dando uma mobilidade que não é encontrada nos computadores. Porém, ainda existem problemas em aberto nessa área. Assim, o objetivo deste doutorado é o de investigar métodos para realizar rastreamento em dispositivos móveis considerando as características desta plataforma. Para atingir esse objetivo, foi conduzido um mapeamento sistemático sobre rastreamento para dispositivos móveis. Este estudo coletou 2.602 artigos, dos quais 444 foram selecionados para classificação. Há um interesse crescente na área e uma preferência por artigos que usam os sensores para rastrear localmente no aparelho. O mapeamento foi usado na criação de experimentos preliminares. Inicialmente, o Google Tango foi avaliado para encontrar um referencial de precisão para os rastreadores atuais. Foi observado que sua precisão em espaços fechados é adequada para uma boa experiência do usuário. Outro experimento avaliou o uso de paralelismo, execução distribuída e implementação nativa. Estes testes mostraram que, na média, a implementação nativa foi a mais eficiente. Além disso, foram criados experimentos para testar diferentes técnicas de rastreamento. Uma delas rastreava faces usando aprendizagem de máquina, foi adaptada considerando as limitações dos celulares e é executada em aproximadamente oito milissegundos. O último experimento está relacionado com uma técnica de SLAM chamada STAM. Ela foi desenvolvida para desktop e portada para um aparelho com suporte ao Google Tango. Várias lições foram aprendidas a partir dos experimentos. Uma delas foi a importância de encontrar informações de semântica de alto-nível de uma cena, que podem ser usadas para melhorar o rastreamento ou criar renderizações mais realísticas. Neste doutorado foi desenvolvida uma técnica que detecta e rastreia primitivas geométricas de maneira incremental usando o processo de geração de nuvens de pontos dos sistemas de SLAM, chamada GS-IST (sigla para Rastreador Semântico Incremental Geométrico e Estatístico). A avaliação do sistema indicou que o GS-IST foi capaz de melhorar a precisão e a estabilidade dos métodos existentes. Porém, a técnica peca na revocação para garantir a detecção e rastreamento das primitivas corretas. Para avaliar como o GS-IST se comportaria em dispositivos móveis, ele foi portado para a plataforma Android. A avaliação mostrou que a versão móvel é de 8,5 a 9,9 vezes mais lenta quando comparada com a versão desktop. Mais do que isso, ele usa até 30,5% da capacidade da CPU, permitindo a execução em uma thread separada da técnica de rastreamento. Além disso, o consumo de energia não foi uma preocupação, uma vez que o GS-IST pode ser executado por mais de 4 horas. Finalmente, o uso de memória RAM foi inferior a 8% do total disponível nos aparelhos testados, o que não apresentou nenhum impacto no tempo de execução.

Palavras-chaves: Rastreamento semântico. SLAM visual. Dispositivos móveis. Rastreamento. Android.

LIST OF FIGURES

Figure 1 – Example of mobile applications that need tracking to deliver experience to the users.	15
Figure 2 – Systematic mapping process. The research question guides the definition of the search strategy, which is used to collect the works. Some criteria are defined to select the relevant studies that are classified in order to provide the systematic mapping.	19
Figure 3 – Tracking type classification diagram.	22
Figure 4 – Selection process shows the number of papers included and excluded and the reasons for exclusions.	25
Figure 5 – Publications over time. Annual trend of included papers.	25
Figure 6 – Tracking type distribution over the database.	27
Figure 7 – Degrees of freedom distribution over the database.	28
Figure 8 – Tracking platform distribution over the database.	29
Figure 9 – Research type distribution over the database.	29
Figure 10 – Annual trend per tracking type. Trends of all tracking types (top); yearly evolution of tracking types combining all vision-based and sensor-based techniques (bottom).	31
Figure 11 – Annual trend per degree of freedom.	31
Figure 12 – Annual trend per tracking platform.	32
Figure 13 – Two dimensional bubble chart: left side presents the tracking type by tracking platform and the right side presents the tracking type by degree of freedom.	32
Figure 14 – Two dimensional bubble chart: left side presents the combined tracking type by tracking platform and the right side presents the combined tracking type by degree of freedom with location service systems combined.	33
Figure 15 – One dimensional bubble chart of degree of freedom by tracking platform.	34
Figure 16 – Evaluation setup consists of a graph paper with precision of one millimeter and measuring 1.5 x 0.55 meters. Red circle highlights the needle used to get the exact position on the paper.	40
Figure 17 – Error dispersion for the small workspace experiment.	42
Figure 18 – Distribution of the device positions on the graph paper (green) and their correspondent positions calculated by the Yellowstone tablet (red).	42
Figure 19 – Error dispersion for the large indoor environment experiment.	43

Figure 20 – Screenshot of one of the paths computed using the Yellowstone tablet. The green arrow points to the initial place and the red one to the final position calculated after a free walk. The error is the average Euclidean distance between them.	43
Figure 21 – Screenshot of the depth estimation of two chessboards printed on a paper. On the right, the one printed with a mix of cyan, magenta and yellow inks. On the left, the same pattern printed with a black ink. Note that the sensor is not able to estimate depth on the black squares of the left paper.	44
Figure 22 – Mean depth estimation error with respect to distance between Tango device and chessboard pattern using different depth interpolation methods.	45
Figure 23 – Door width estimation using the Yellowstone tablet. Left side shows the initial measurement and the right side shows the ruler position after moving the device. Door actual width is 0.7 meters.	45
Figure 24 – Multi-Thread Partial Native flow diagram.	48
Figure 25 – Client/Server flow diagram.	49
Figure 26 – Full Native flow diagram.	49
Figure 27 – Screen capture of the system tracking the template. The blue square is the virtual content that is placed on top of the model.	50
Figure 28 – Execution time in milliseconds on every tracking stage for each implemented architecture. For Client/Server approach, feature matching, matching filtering and pose calculation also includes the time to transfer the data to the server and back to the device.	51
Figure 29 – RAM memory consumption in MB during the execution of each architecture implementation.	51
Figure 30 – ROM memory in MB required to storage each architecture implementation.	52
Figure 31 – For every image on the training dataset, each local feature is learned individually from the landmarks (white circles). The intensity difference between two random pixels (white crosses) is used as decision function to split the training images. Each node has a distinct pair of features.	55
Figure 32 – Local region around the landmark on every stage.	56
Figure 33 – Traverse of three different landmarks of one of the training images on the generated forest. The sequence of zeros and ones of every landmark is the local binary feature.	56
Figure 34 – Cascade shape regressor, in which the landmark position is incremented every stage.	57
Figure 35 – Result from C++ implementation of the LBF technique.	58
Figure 36 – Standard landmark configuration with 68 points (left) and the 31 landmarks selected (right).	59

Figure 37 – Initial guess chosen as if the image is in the upright position (a). If it is used on the rotated image, it will lead to an incorrect result (b). Therefore, if the initial guess is rotated using the device’s orientation, it will be on the correct position (c).	60
Figure 38 – Tracking results with different device orientations.	61
Figure 39 – STAM flow diagram.	64
Figure 40 – Frame from which the patches were extracted. Six samples of these patches are on the right. The 2D position of each patch should be found in the same image.	65
Figure 41 – Six samples of patches (top right) were extracted from the first frame (top left). These patches should be tracked along the image sequence (bottom row).	65
Figure 42 – Six samples of patches (top row) were extracted from the first frame and nine sample images from the sequence to be tracked. Note that the patches are disappearing along the image sequence.	66
Figure 43 – Schematic of the on-site tracking area. Each rectangle represents a table with different objects. The numbers X/Y indicate the poster used to assign the four competition points in which X represents the point order and Y the competition session.	67
Figure 44 – Images from the tracking area. The chessboard used to calibrate the tracking system is seen on the top left image. The other images show the tables with the trackable objects and the posters to mark the 3D points.	67
Figure 45 – Evaluation environment with the calibration chessboard in the middle and objects with different types of texture around it.	70
Figure 46 – Execution time in milliseconds to run the main steps of STAM.	70
Figure 47 – Tracking procedure running on Tango device. Small points are the keypoints extracted from the mapping while the large green dot shows the tracking point, which is a known point in the real world based on the chessboard template.	71
Figure 48 – Test scene (a) and (d) and their reconstructed point cloud (b) and (e). Eight shapes were detected on the first keyframe (c) and twelve primitives were found on the second one (f). Even with the point cloud being very similar, the red points represent the six primitives that were differently detected between keyframes.	75
Figure 49 – Flow of the Geometric and Statistical Incremental Semantic Tracking (GS-IST) approach.	77
Figure 50 – Difference between the input points to their correspondent projected points on a shape estimated correctly (left) and incorrectly (right). . .	79

Figure 51 – Fusion of parameters for different classes of shapes.	80
Figure 52 – Comparison of precision, recall and $F_{0.5}$ -Score between Efficient RANSAC (SCHNABEL; WAHL; KLEIN, 2007) and GS-IST.	85
Figure 53 – Precision over time of Efficient RANSAC (SCHNABEL; WAHL; KLEIN, 2007) and GS-IST on <i>Case 1</i>	85
Figure 54 – The first and third rows show the result of GS-IST on each test case. Blue labels represent planes, green ones are for spheres and red for cylinders. The second and fourth rows show one particular view of the input point cloud (in red) and some of the estimated primitives (in blue).	87
Figure 55 – Average distance in millimeters of each input point to its projection on the estimated primitive over time for <i>Case 1</i>	87
Figure 56 – Error in millimeters over time of the cylinder and the sphere radii detected in <i>Case 1</i>	88
Figure 57 – Percentage of points that were labeled to each primitive.	89
Figure 58 – Memory (in KB) required to describe a scene using the point cloud and the data structure of the detected primitives for <i>Case 1</i>	89
Figure 59 – The application indicates the shape the user has to find (left image). When the correspondent shape is centered (top-right), it gives a positive feedback (bottom-right) and moves to the next primitive.	90
Figure 60 – Results of GS-IST running on a Samsung Galaxy S8 and an ASUS ZenFone 3. Blue labels represent planes, green ones are for spheres and red for cylinders.	96
Figure 61 – GS-IST execution time in milliseconds divided by stages on desktop and two different mobile devices.	97
Figure 62 – CPU load and Normalized CPU load over time of all five test cases on a ZenFone 3 (top) and Galaxy S8 (bottom).	97
Figure 63 – Energy consumption (in W) over time of all five test cases on ZenFone 3 (top) and Galaxy S8 (bottom).	98
Figure 64 – Memory usage over time of GS-IST running on ZenFone 3. Each test case has a different time scale.	99
Figure 65 – Memory usage over time of GS-IST running on Galaxy S8. Each test case has a different time scale.	100

LIST OF TABLES

Table 1	– List of the most popular publication forums.	26
Table 2	– List of sensors and their most used combinations.	28
Table 3	– Relevant papers for each classification.	30
Table 4	– Evaluation of main aspects of the most promising works. Cells with asterisk mean that there is not a clear value for the feature.	54
Table 5	– Mean reprojection error in millimeter of all points on the last frame of off-site category Level 3.	68
Table 6	– Mean reprojection error in millimeters of on-site category.	69
Table 7	– History of the estimated shape from a primitive. For each sample that represents a keyframe K_i , it was classified as plane (P), sphere (S) or cylinder (C).	83
Table 8	– Details of the dataset used for the evaluation, in which P, S and C stands for Plane, Sphere and Cylinder, respectively.	84
Table 9	– The influence of modifications in GS-IST on the final precision and recall in <i>Case 1</i>	86
Table 10	– Summary of the technical specifications of the devices used for evaluation.	94
Table 11	– Time that different applications take to fully drain a battery fully charged on ZenFone 3 and Galaxy S8 devices.	101

CONTENTS

1	INTRODUCTION	15
1.1	OBJECTIVES	16
2	SYSTEMATIC MAPPING	18
2.1	METHODS	19
2.1.1	Research Questions	19
2.1.2	Scientific Databases and Search Strategy	20
2.1.3	Screening of Papers	21
2.1.4	Classification	21
2.1.4.1	Tracking Type	21
2.1.4.2	Degrees of Freedom	22
2.1.4.3	Tracking Platform	23
2.1.4.4	Research Type	23
2.1.5	Threats to Validity	24
2.2	RESULTS	25
2.3	MAPPING	27
2.3.1	Classification Distribution	27
2.3.2	Classification Trends	29
2.3.3	Classification Relationship	30
2.4	DISCUSSION	34
2.4.1	Implications for Future Studies	36
3	PRELIMINARY EXPERIMENTS	38
3.1	EVALUATION OF TANGO PLATFORM	38
3.1.1	Tango Platform	38
3.1.2	Evaluation Methodology	39
3.1.2.1	Motion Tracking	40
3.1.2.2	Depth Sensing	41
3.1.3	Results	41
3.1.3.1	Motion Tracking	41
3.1.3.2	Depth Sensing	44
3.1.4	Discussion	44
3.2	EFFICIENT TRACKING ON MOBILE DEVICES	46
3.2.1	Android Architectures for Computer Vision Tracking	46
3.2.1.1	Multi-Thread Partial Native Implementation	47
3.2.1.2	Client/Server Implementation	48

3.2.1.3	Full Native Implementation	48
3.2.2	Architectures Evaluation	50
3.2.3	Discussion	52
3.3	MACHINE LEARNING TRACKING ON MOBILE DEVICES	53
3.3.1	Research Methodology	53
3.3.2	Local Binary Features Technique	54
3.3.2.1	Learning The Feature Mapping Function ϕ_t	55
3.3.2.2	Learning The Global Linear Regression Matrix W_t	56
3.3.2.3	Tracking a New Face	56
3.3.3	Implementation on Android	57
3.3.3.1	Improvements	59
3.3.4	Discussion	60
3.4	SIMPLE SLAM SYSTEM ON MOBILE DEVICES	61
3.4.1	SLAM Systems Works on Mobile Devices	62
3.4.2	Simple Tracking and Mapping	62
3.4.3	STAM Evaluation	63
3.4.3.1	Off-Site Competition	64
3.4.3.2	On-Site Competition	66
3.4.3.3	Results	67
3.4.4	Implementation on Tango Device	68
3.4.4.1	Android Programming	69
3.4.4.2	Results	69
4	INCREMENTAL SEMANTIC TRACKING	72
4.1	SEMANTIC MODELING	72
4.2	RANSAC-BASED METHOD ON SPARSE POINT CLOUD	74
4.2.1	Method Overview	74
4.2.2	Evaluation	75
4.3	GEOMETRIC AND STATISTICAL INCREMENTAL SEMANTIC TRACKING	76
4.3.1	Efficient RANSAC	76
4.3.2	Shape Fusion	78
4.3.2.1	Parameter Computation	78
4.3.2.2	Inclusion Criteria	80
4.3.3	Shape Matching	81
4.3.4	Shape Update and Recovery	81
4.3.5	Reliability Computation	82
4.3.5.1	Geometric Analysis	82
4.3.5.2	Statistical Analysis	83
4.4	EVALUATION	83
4.4.1	Metric Evaluation	86

4.4.2	Runtime Evaluation	88
4.4.3	Segmentation Evaluation	88
4.4.4	Point Cloud Representation Evaluation	88
4.4.5	Proof of Concept	90
4.4.6	Evaluation on Dense Point Cloud	90
5	MOBILE EVALUATION OF INCREMENTAL SEMANTIC TRACKING	92
5.1	SEMANTIC MODELING ON MOBILE DEVICES	92
5.2	MOBILE IMPLEMENTATION	93
5.2.1	Evaluation	94
5.2.1.1	Execution Time	95
5.2.1.2	Energy Consumption	96
5.2.1.3	Memory Usage	98
5.3	DISCUSSION	101
6	FINAL CONSIDERATIONS	103
6.1	FUTURE WORK	104
6.2	CONTRIBUTIONS	105
6.3	PUBLICATIONS	105
	REFERENCES	107

1 INTRODUCTION

Several applications require tracking, which is the computation of an object placement relative to a real world element or location over a time period. For instance, some augmented reality software use the camera position related to a marker to display a virtual content registered with the pattern (ROBERTO et al., 2013). Another example is a GPS-based navigation system that calculates its location relative to the road in order to show the driver directions to a destination (LESHED et al., 2008), seen in Figure 1 (a). However, this is still a challenging task. Moreover, determining this placement can demand a lot of computational power and memory depending on the approach and the required information.



Figure 1 – Example of mobile applications that need tracking to deliver experience to the users.

Mobile devices, such as phones and tablets, are becoming increasingly popular. Research shows that a median of 43% of the world's population owns a smartphone (Pew Research Center, 2016). Moreover, these devices are continuously improving regarding processing power and memory space available (HALPERN; ZHU; REDDI, 2016), which makes them powerful enough to perform complex tasks, such as tracking. In fact, the processing power on mobile devices is increasing rapidly enough to reduce the gap with desktop computers. While in 2009 the desktop CPU clock frequency was 5.7 times the mobile CPU clock, in 2015 this distance dropped to 2.1 times. This scenario favors the creation of numerous types of applications since such devices create several opportunities that are only possible when the user can be mobile. One example is to annotate relevant information precisely on the facade of buildings in an outdoor environment (YOVCHEVA; BUHALIS; GATZIDIS, 2012), as illustrated in Figure 1 (b).

Besides all the improvement on the mobile devices itself, several tracking techniques that are capable to run on such devices were released in the last couple of years. There are examples in the academy and in the industry. The most distinguished were the ARCore¹ and ARKit², by Google and Apple and shown in Figure 1 (c) and (d) respectively. However, there are still several gaps in the field that can be illustrated with a challenging scenario, such as creating a 3D model of an entire house in order to redecorate it. Tracking an environment that is so large demands manipulation of a large amount of data, which impacts both the processing power, especially because the architecture of mobile processors compromise on speed to be more efficient on energy consumption and on temperature. Although memory is not so limited on these devices nowadays, using algorithms such as bundle adjustment can be an issue in a scenario like this. Additionally, the walls must be aligned in order to correctly recreate all the rooms. And in each room there are several objects that also need to be recognized and modeled, such as tables, sofas and wardrobe.

This particular scenario presents several research opportunities and one is extracting and tracking high-level semantic information. Several types of semantic can be collected, from the geometric primitives of the objects to the model of a piece of furniture and this process is referred to as semantic modeling. Shape parameters of geometric primitives and the relationship between them are valuable knowledge to be estimated especially in human-made environments such as a house. These semantics are useful for replacing redundant point clouds with more efficient data structures or aligning the walls based on their parameters. This data can also be gathered in different ways: from the input image, the scene map or a combination of both.

1.1 OBJECTIVES

In this sense, the objective of this Ph.D. thesis is to **investigate methods to perform tracking on mobile devices considering the characteristics of such platform.**

In order to achieve this goal, four specific objectives should be met. The first one is to **perform a strong literature review of tracking for mobile devices to find relevant studies and gaps.** One effective way to analyze the related works of a research area that has so many studies is through a systematic mapping, which is a research method to review, classify and provide an overview of a wide range of papers on a particular topic. *Chapter 2* describes how the systematic mapping was performed in this work and presents its major findings.

The second specific objective of this Ph.D. is to determine experimental scenarios aiming to **acquire a practical knowledge on the specificities of developing tracking techniques for mobile devices**, which are based on the results of the systematic mapping. The goal is to evaluate different tracking approaches for mobile device in order to

¹ <<https://developers.google.com/ar/>>

² <<https://developer.apple.com/arkit/>>

assess a state-of-the-art tracking method and find an efficient architecture for a computer vision system on mobile devices. Additionally, test different tracking techniques to decide which ones are suitable for such devices and which ones are not, as explained in *Chapter 3*.

The lessons learned from both the systematic mapping and the experiments were the foundation to **create a semantic tracking technique for sparse visual SLAM**, which is the third specific objective of this Ph.D. This approach incrementally detects and tracks primitive shapes using geometric and statistical analyses. *Chapter 4* details this method and presents the evaluation results.

And last but not least, the fourth specific objective of this study is to **port the semantic tracking approach to mobile devices and evaluate how it performs running on such platforms**. *Chapter 5* shows the different criteria used to evaluate this technique and discusses the implication of the results.

2 SYSTEMATIC MAPPING

During the past years, researchers have proposed different techniques to perform tracking on mobile devices. A preliminary search for related works revealed a significant amount of papers in this area. Therefore, it becomes important to summarize the current state of the art and provide an overview of the trends in this specialized field. In order to address this issue, it was performed a systematic mapping of the literature in this area. The main goal of this mapping is to analyze, classify and map existing papers about tracking for mobile devices, providing a primary study and an inclusive overview of this topic. The result of the systematic mapping was published in (ROBERTO; LIMA; TEICHRIEB, 2016), and an update of this study is detailed in sequence.

Systematic mapping is a method to review, classify and structure papers related to a specific research field (PETERSEN et al., 2008). It is frequently used in medical research and lately has been applied to software engineering. Unlike systematic reviews, the goal of this research method is not to perform a deep analysis of works in order to identify the best practices of a field, which usually includes a quality evaluation. The aim of a systematic mapping is to provide an overview of a wide range of papers. This broader analysis enables to observe more studies, which allows more general conclusions (PETERSEN et al., 2008). Nevertheless, both methods use a well-defined methodology, which reduces bias (KEELE, 2007). Moreover, systematic mapping papers have an educational value to provide valuable information for students and young researchers, being a useful first step for Ph.D. candidates (KITCHENHAM; BRERETON; BUDGEN, 2010).

To the best of the author's knowledge, there is currently no study that synthesizes or systematically analyzes, classifies and maps existing papers about tracking for mobile devices. However, some surveys were found on the field or one of its specific subareas. For instance, (LIU et al., 2007) evaluated wireless indoor localization techniques and (LANE et al., 2010) listed tracking algorithms for mobile phones that use only their sensors, as well as related applications. There are also surveys regarding mobile augmented reality, in which tracking is an important step. Examples are (OLSSON; SALO, 2011) that studied the overall acceptance and user experience of mobile augmented reality consumer applications, (HUANG et al., 2013) that presented the technologies and methods to perform augmented reality on mobile devices and introduces some applications, and (GRUBERT; LANGLOTZ; GRASSET, 2011) that conducted a survey about augmented reality browsers and performed a quantitative and qualitative analysis regarding the usability aspects of these tools.

As a definition, tracking for mobile devices means that an off-the-shelf cell phone or tablet extracts information from the environment and then processes it locally or remotely in order to compute the device's pose related to the world, which will be used by an

application or a service on the device itself.

2.1 METHODS

The systematic mapping was conducted based on the process proposed by (PETERSEN et al., 2008) and illustrated in Figure 2 in which a list of research questions is proposed, which guides the search strategy, the definition of inclusion and exclusion criteria for relevant studies and the classification schema of all the selected studies. The process steps performed in this study are described in the following subsections.

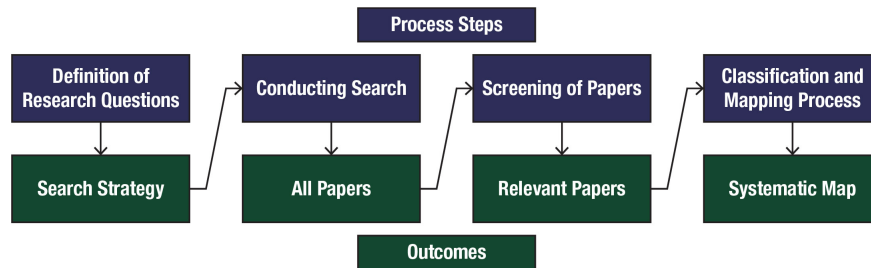


Figure 2 – Systematic mapping process. The research question guides the definition of the search strategy, which is used to collect the works. Some criteria are defined to select the relevant studies that are classified in order to provide the systematic mapping.

2.1.1 Research Questions

The goal of this systematic mapping is to provide an overview of the current research on the topic of tracking for mobile devices. The overall objective was defined in the following four research questions:

RQ1 How has the frequency of research on tracking for mobile devices changed recently?

RQ2 What are the most common approaches of tracking for mobile devices?

RQ3 In which platforms has tracking for mobile devices been executed?

RQ4 In which forums has research on tracking for mobile devices been published?

The first question aims to use the number of publications to investigate trends of the field in the past few years. The second and third questions explore the approaches and platforms researched in the field. The objective of the fourth question is to identify where tracking for mobile devices research can be found, which could be targets for the publication of future studies.

2.1.2 Scientific Databases and Search Strategy

Three online academic search engines were used to find the relevant papers:

- ACM Digital Library;
- IEEE Xplore Digital Library;
- ScienceDirect.

In order to perform an automatic search on the selected libraries, the search string consisted of two parts. The former regards the tracking domain and the later covers the device used. Thus, the search string was the following:

$$\begin{aligned} & (“tracking” OR “registration” OR “localization”) \\ & \quad AND \\ & (“phone” OR “tablet” OR “handheld” OR “smartphone”) \end{aligned}$$

Tracking is the key term of the first segment and the other ones are its most used synonyms. Other terms were not used because a preliminary analysis showed that the majority of the papers found would not be selected for classification. An example is “positioning”, which appears mostly in studies in which the device’s pose is used only by an external agent and not on the device itself, such as the phone’s position that is employed by the carrier to determine in which GSM antenna it will connect to. Moreover, the analysis revealed that the relevant papers were already found using the chosen terms.

Regarding the second segment, it was chosen to search for each device instead of using the terms “mobile” or “mobile device”. The reason is that these keywords returned too many papers and a preliminary analysis revealed that the vast majority of them use a broader concept of mobile devices than the desired in this mapping. For instance, some works use the term “mobile objects” for objects with sensors embedded that are tracked by computers. There are also several references to mobile device as a large object that is used for tracking-related activities, such as airplane radar or medical scanner, which was shrunk to become mobile. Therefore, using the types of mobile device as search terms showed to be more efficient.

An automatic search was performed in the aforementioned databases using an open-source paper crawler¹ software and applying the search in the title, abstract and keywords. The crawler was developed during this Ph.D. and aims to automate the process of retrieving papers. Hence, the crawler accesses the digital libraries, performs the search using the search strings, collects the papers, eliminates duplicate versions and creates a worksheet containing all the works with their title, year, source, abstract and web address.

¹ Available at <<https://goo.gl/7t8kG8>>

2.1.3 Screening of Papers

After collecting the papers, the crawler automatically remove duplicate works. Whenever a work had multiple publications, only the most complete version was selected and the other ones were removed as duplicates. Later, relevant papers were manually selected using the following inclusion and exclusion criteria.

- Inclusion criteria:
 - Papers about tracking techniques implemented on mobile devices;
 - Papers about mobile applications that use existing tracking techniques, even if they do not explain how tracking was implemented.
- Exclusion criteria:
 - Papers published before 2009, which is one year after the release of phone models that allowed 3rd party development;
 - Papers not written in English language;
 - Papers published on non-peer reviewed vehicles, such as books and magazines;
 - Papers not related to tracking techniques on mobile devices;
 - Papers about tracking techniques that were implemented only on desktop platform and that have no indication of how they can be developed for mobile devices.

2.1.4 Classification

Following, all included papers were classified according to four properties in order to answer the research questions. They are detailed next.

2.1.4.1 Tracking Type

Each paper was classified regarding its tracking type. The classification was adapted from (ZHOU; DUH; BILLINGHURST, 2008), which is shortly explained as follows and illustrated in Figure 3.

- **Sensor-Based Tracking:** techniques that calculate device's pose relative to real world using exclusively non-vision based sensors. This approach can be divided in two categories: *single sensor*, which uses only one sensor for tracking, and *sensor fusion*, which uses different sensors to perform the same task;
- **Vision-Based Tracking:** techniques that use images captured by the device cameras to calculate pose relative to the real world. This approach can also be divided in two categories: *marker-based* and *natural feature-based*. The former method calculates

device's pose from artificial markers placed in the scene and the latter performs the same task using natural characteristics from the environment, such as points and edges. The natural feature-based approach was also split into two subcategories: *static model* and *dynamic model*. The first one uses prior knowledge of the scene that does not change during tracking to compute device's pose and in the second one the tracker can use an initial model if it is available or build it entirely from scratch and this environment information is updated during computation of the device's pose;

- **Hybrid Tracking:** techniques that combine sensor-based and vision-based methods to calculate device's pose;
- **Severall:** papers that present techniques from several categories, such as surveys.

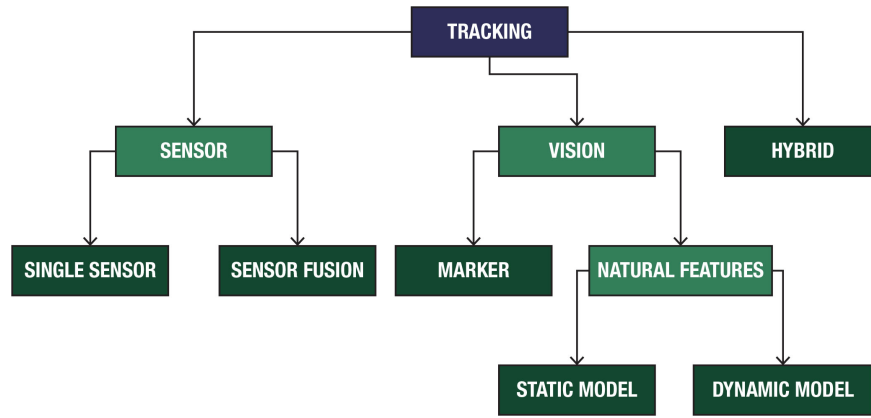


Figure 3 – Tracking type classification diagram.

2.1.4.2 Degrees of Freedom

This property details the degree of freedom required to compute the information desired. This classification was based on (NORMAND; MOREAU, 2012). One modification was the addition of the 3D degree of freedom, which was not mentioned in the original work. Thus, the complete degree of freedom classification used in this work is detailed in the following list.

- **0D:** techniques that detect a pattern and display an information about it without any relationship with its position and orientation;
- **2D:** techniques that provide information about the position, being indoor, outdoor or in the screen. It can also be called “2D Location”;
- **2D + θ :** techniques that extend the position information with orientation, providing the location with direction. It can also be named “2D Location + Orientation”;
- **3D:** techniques that compute the device's rotation in all three axis;

- **6D:** techniques that calculate the device's pose with rotation and translation. Systems that also compute scale were considered 6D as well;
- **Several:** papers that present techniques from several categories.

2.1.4.3 Tracking Platform

Two tracking platforms were considered to classify papers regarding this property, as detailed below.

- **Local Tracking:** techniques that compute all the required information at the mobile device;
- **Distributed Tracking:** techniques in which part or all the information is calculated on a server and the result is transmitted to the device and used to display the content;
- **Several:** papers that present techniques from both categories.

2.1.4.4 Research Type

The research type feature concerns the research approach used in the papers. This classification was adapted from (WIERINGA et al., 2005) and is summarized in the list below.

- **Evaluation Research:** papers that present implementation and extensive evaluation of existing techniques in order to determine their benefits and drawbacks;
- **Opinion Papers:** publications in which the author expresses a personal opinion whether a certain topic is good or bad without relying on related work;
- **Philosophical Papers:** papers that present new ways of looking at existing things, such as structuring the field in form of a new taxonomy;
- **Proposal of Solution:** works that propose solutions for problems, which can be based on novel or existing techniques;
- **Survey Papers:** papers that summarize and organize a research field based on other publications;
- **Technique Research:** publications in which the authors propose and implement a novel technique.

2.1.5 Threats to Validity

It is important to consider threats to validity in order to judge the systematic mapping strengths and limitations. The main issues are related to incomplete sets of relevant papers and researcher bias with regards to inclusion/exclusion criteria and classification.

Limitations with search string, scientific databases and search strategy can result in an incomplete set of relevant papers. As a way to mitigate that risk, three strategies were used. In order to validate the search string, the terms were discussed with five other experienced researchers in the field of tracking. The scientific databases that publish works from the most important conferences and journals in the area were selected. As for the search strategy, a different approach was used to maximize the number of papers found. Instead of using the complete search string, twelve different searches were performed using a two-by-two combination of every term in both parts of the search string. Using this strategy, it was possible to retrieve almost 34 times more papers than when the complete string was employed.

The Ph.D. candidate conducted the analysis to include/exclude and classify a paper. Since this may lead to a researcher bias, 11.74% of the studies were randomly selected before the subjective part of the screening phase to compose a set of control papers, and one of the co-advisors, which is a experienced researcher in the field of tracking, analyzed them. The results were compared using the Cohen's Kappa coefficient, which measures the agreement between the two classifications taking into account how much agreement would be expected to be present by chance (COHEN, 1960). The coefficient lies between -1.0 and 1.0 in which 1.0 denotes perfect agreement, 0.0 indicates that any agreement is due to chance and negative values present agreement less than chance. Cohen's Kappa was used to measure the reliability regarding inclusion and exclusion of papers and the classification of the included papers in common according to the classification schema. There is no consensus on what are good levels of agreement. Nevertheless, a common scale (ALTMAN, 1990) indicates that there is no agreement for negative values, poor agreement between 0.00 and 0.20, fair agreement between 0.21 and 0.40, moderate agreement between 0.41 and 0.60, good agreement between 0.61 and 0.80 and very good agreement for values higher than 0.80. At first, the classification ratio was in the range of good agreement. The main reason for that was the fact that the first classification schema was leading to dubious interpretations. For instance, natural feature tracking was divided into model-based and model-less approaches, in which it was not clear if information used could be considered a model or not. The classification schema was refined to the one previously presented, which uses a more straightforward classification, and all papers were then reclassified. Thus, the included/excluded papers Cohen's Kappa coefficient was 0.8062 ± 0.0495 and the Cohen's Kappa classification was 0.8345 ± 0.0303 .

2.2 RESULTS

The search was made on March 3rd, 2016 and resulted in 2,602 papers found. As can be seen in Figure 4, 611 papers were removed for being duplicated and 1,991 studies were available for the subjective steps of screening. Only 444 works remained for trends analysis and classification.

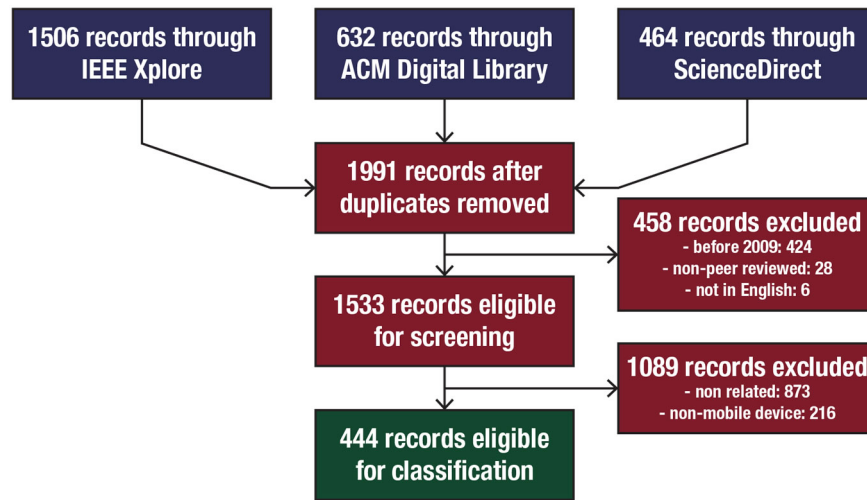


Figure 4 – Selection process shows the number of papers included and excluded and the reasons for exclusions.

Figure 5 shows the annual trend of papers. It is possible to see that the number of studies is growing between 2009 and 2014. Even with a reduction in the number of publications in 2015, there is an indication of increasing interest in tracking for mobile devices in recent years.

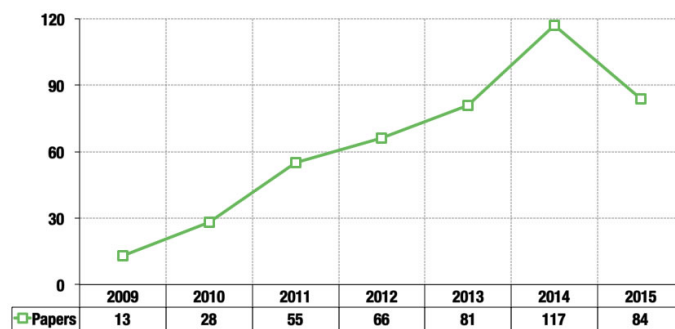


Figure 5 – Publications over time. Annual trend of included papers.

The 444 papers were published in 260 forums. As seen in Table 1, almost a quarter of all works came from the top 12 venues. ISMAR is the flagship event in the field with 30 studies. Preferable targets for such papers are conferences and symposiums in which 196 papers were published. They were followed by 52 journal works and 12 workshop studies.

Each paper was classified according to the scheme presented in the previous section. The full list of works can be accessed through an open-source web application (FIGUEIREDO

et al., 2015). Using this system, it is possible to filter the papers according to the year when the works were published and the classification criteria, as well as search for terms and words in the abstract. Moreover, collaborators can send new entries of studies about tracking for mobile devices, which will be revised and then added to the online data set. The web application can be accessed at http://cin.ufpe.br/~rar3/tracking_sm/.

Table 1 – List of the most popular publication forums.

Forum	Acronym	Type of Forum	Number of Papers	Percentage of the Total
International Symposium on Mixed and Augmented Reality	ISMAR	Symposium	30	6.76%
International Conference on Indoor Positioning and Indoor Navigation	IPIN	Conference	19	4.28%
Conference on Embedded Networked Sensor Systems	SenSys	Conference	7	1.58%
International Conference on Mobile Computing and Networking	MobiCom	Conference	6	1.35%
Pervasive and Mobile Computing	-	Journal	6	1.35%
IEEE Transactions on Mobile Computing	-	Journal	5	1.13%
IEEE Virtual Reality Conference	IEEE-VR	Conference	5	1.13%
International Conference on Computer Vision	ICCV	Conference	5	1.13%
Conference on Multimedia and Expo	ICME	Conference	5	1.13%
International Conference on Pervasive and Ubiquitous Computing	UbiComp	Conference	5	1.13%
International Conference on Pervasive Computing and Communications	PerCom	Conference	5	1.13%
Symposium on 3D User Interfaces	3DUI	Symposium	5	1.13%
Other 248 Forums	-	-	341	76.80%

2.3 MAPPING

From the classification of the studies it is possible to establish a mapping that aims to provide an overview of tracking for mobile devices and can help to identify potential research gaps. This map gives the distribution of works for each classification criteria, their annual trends and the relation between them.

2.3.1 Classification Distribution

It is possible to see in Figure 6 that most of the works, such as (ANDO et al., 2014), rely on a combination of the devices' sensors to calculate pose, and that marker-based tracking, which is used for example in (OUI; NG; KHAN, 2011), is the least used method for the same task. Moreover, it can be noted that vision-based methods like (WAGNER et al., 2010b) are present in three out of ten papers.

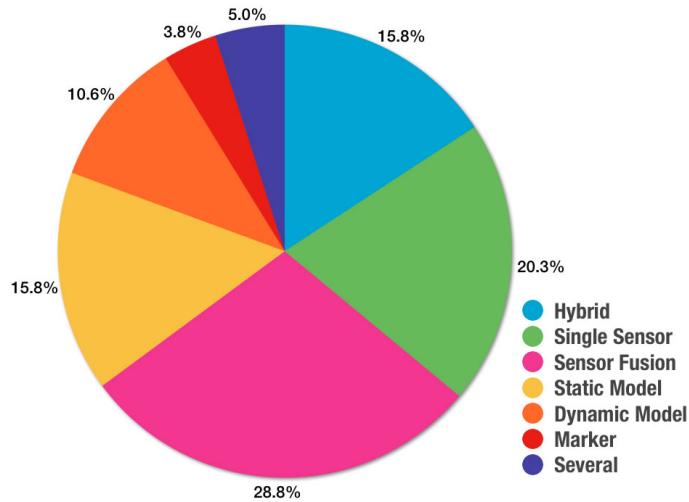


Figure 6 – Tracking type distribution over the database.

Table 2 lists the sensors used on each paper in which the tracking type is hybrid or based on sensors as well as the number of studies that uses them. Several works fuse different sensors and Table 2 also lists the 13 most common combinations.

Regarding the degree of freedom found in the works, several of them calculate a 6D pose (HAGBI et al., 2011), as shown in Figure 7. However, in 58.5% of the studies, a 2D position is computed. In some papers, a 2D position on the screen is found (SONG; LIU; CHEN, 2011), but the majority discovers this 2D position on the environment (HU et al., 2010). In the latter case, more works also find the orientation θ (SHIN et al., 2012), and the least common papers are the ones aiming 0D systems (TERAURA; SAKURAI, 2012).

In the majority of the works all the processing needed to calculate a pose is done at the device (ENGELKE et al., 2013). Only 14.0% use a remote server to assist in this task (VENTURA; HOLLERER, 2012) or completely perform pose calculation (HA et al., 2011), as illustrated in Figure 8.

Table 2 – List of sensors and their most used combinations.

Sensor	Number of Papers	Combination of Sensors	Number of Papers
Accelerometer	137	Wi-Fi	40
Magnetometer	98	GPS	36
Wi-Fi	97	Accelerometer, Gyroscope and Magnetometer	26
GPS	96	Accelerometer and Gyroscope	22
Gyroscope	90	Accelerometer	11
Cellular Network (GSM, CDMA)	32	Accelerometer, Gyroscope, Magnetometer and Wi-Fi	11
Acoustic	19	Accelerometer and Magnetometer	9
Barometer	10	GPS and Wi-Fi	8
Bluetooth	10	Cellular Network	8
Depth	6	Cellular Network and GPS	7
Illuminance	3	Accelerometer, GPS and Magnetometer	7
Thermal	1	Accelerometer and GPS	6
Radio	1	Other 50 combinations	97

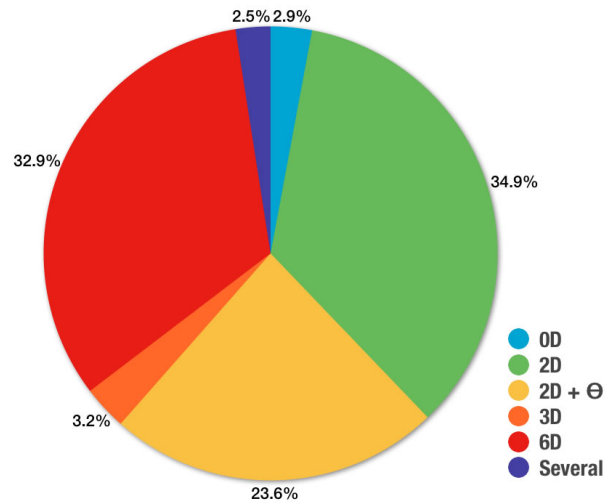


Figure 7 – Degrees of freedom distribution over the database.

Table 3 lists relevant studies for each classification, those with more citations per year.

As for research type, approximately two thirds of the papers propose a new technique to perform tracking (SHANKLIN; LOULIER; MATSON, 2011) and 28.6% use an existing method to develop a mobile solution that requires tracking (POLO et al., 2014), as can be seen in Figure 9.

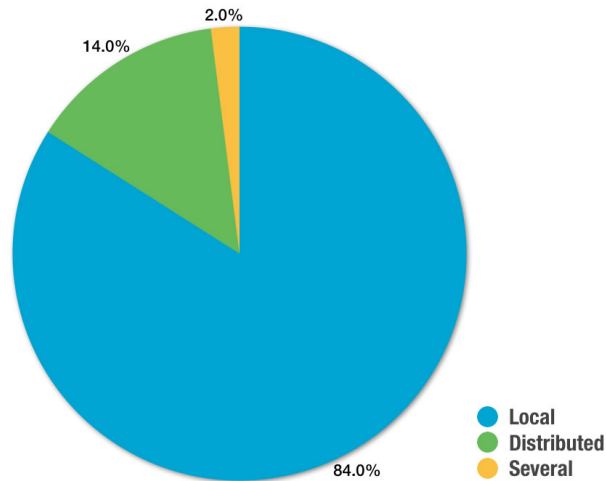


Figure 8 – Tracking platform distribution over the database.

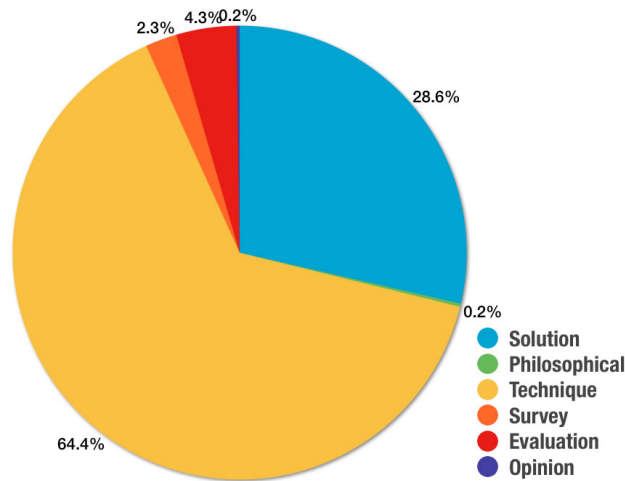


Figure 9 – Research type distribution over the database.

2.3.2 Classification Trends

The annual trend per tracking type shows that until 2015 the number of papers about sensor-based systems is increasing. Moreover, the number of works that use a single sensor in 2015 is almost 5 times higher than in 2009 and 13 times larger for sensor fusion techniques, as can be seen in Figure 10. From Figure 10 (top) it is also possible to conclude that the other tracking types have an overall growing tendency. From 2009 to 2015 natural feature solutions went from 3 works to 8 with static model studies and 3 to 8 with dynamic model papers. On the same period, hybrid solutions went from 0 to 20 studies and it was the only tracking type that had more publications in 2015 than in 2014. The growth of marker studies occurred in the last three years.

Regarding the annual trend per degree of freedom, it is possible to see in Figure 11 an increasing number of publications about 0D, 2D, 2D + θ and 6D trackers. Respectively, they went from 0 to 5, 4 to 27, 1 to 22 and 7 to 27 between 2009 and 2015. The image

Table 3 – Relevant papers for each classification.

Tracking Type	Relevant Studies
Hybrid	(KURZ; BENHIMANE, 2011)
Single Sensor	and (VENTURA; HOLLERER, 2012)
Sensor Fusion	(SHIN et al., 2010)
Static Model	and (GOZICK et al., 2011)
Dynamic Model	(CHON; TALPOV; CHA, 2012)
Marker	and (ZHANG et al., 2013)
	(WAGNER et al., 2010b)
	and (HU et al., 2013)
	(KLEIN; MURRAY, 2009)
	and (WAGNER et al., 2010a)
	(OUI; NG; KHAN, 2011)
	and (GHERGHINA; OLTEANU; TAPUS, 2013)
Degree of Freedom	Relevant Studies
0D	(RAI et al., 2012)
2D	and (XU et al., 2014)
2D + θ	(SHIN et al., 2010)
3D	and (LV, 2013)
6D	(SCHALL; MULLONI; REITMAYR, 2010)
	and (SHIN; CHON; CHA, 2012)
	(LI; KIM; MOURIKIS, 2013)
	and (ELLOUMI et al., 2013)
	(TAKACS et al., 2010)
	and (TANSKANEN et al., 2013)
Tracking Platform	Relevant Studies
Local	(ARTH et al., 2009)
Distributed	and (SCHÖPS; ENGEL; CREMERS, 2014)
	(CHEN et al., 2009)
	and (VENTURA et al., 2014)

also shows that the community did not demonstrate the same interest in systems with a 3D approach.

Most of the works use a local approach to calculate the device's pose and this fact is reflected in the annual trends per tracking platform, as shown in Figure 12.

2.3.3 Classification Relationship

The relationship between the classifications can provide a powerful and quick overview of tendencies on tracking for mobile devices. A bubble chart was used because it offers

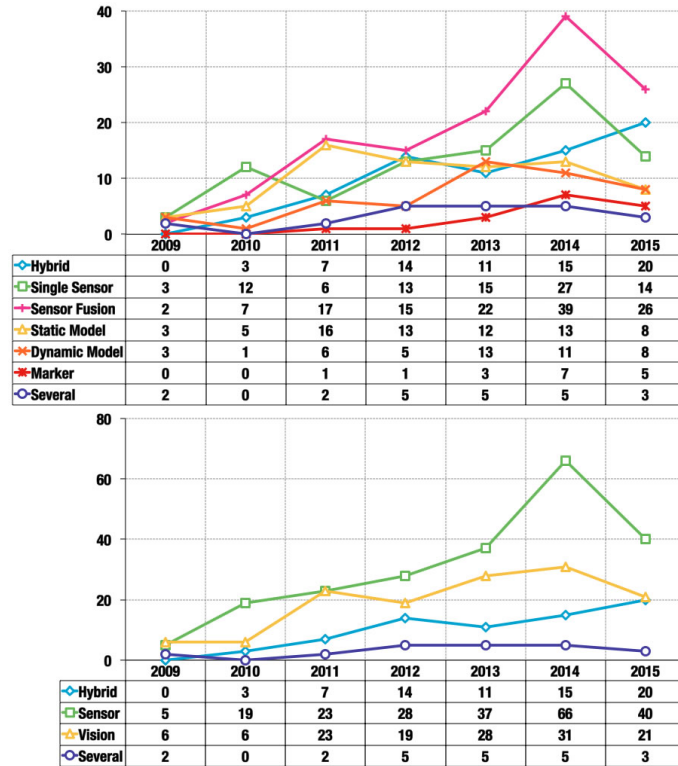


Figure 10 – Annual trend per tracking type. Trends of all tracking types (top); yearly evolution of tracking types combining all vision-based and sensor-based techniques (bottom).

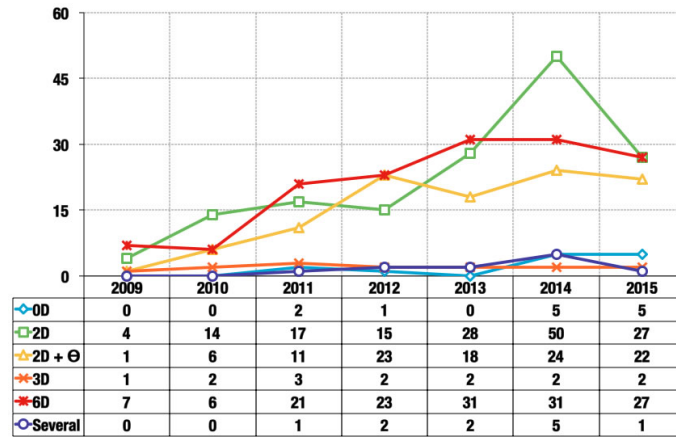


Figure 11 – Annual trend per degree of freedom.

a more visual result than tables. Figure 13 presents a bubble plot in two dimensions in which the leftmost represents the tracking type by tracking platform and the rightmost displays the tracking type by degree of freedom. It should be noted in the first dimension that the ratio of publications of local systems is at least four times larger than the ratio of distributed approaches.

The same balance cannot be seen in the second dimension, in which the majority of the sensor works are location-based solutions, such as (MICHAEL; CLARKE, 2013) that compute a 2D pose and (JUNG; LEE; JEONG, 2011) for 2D + θ papers. Only three publications

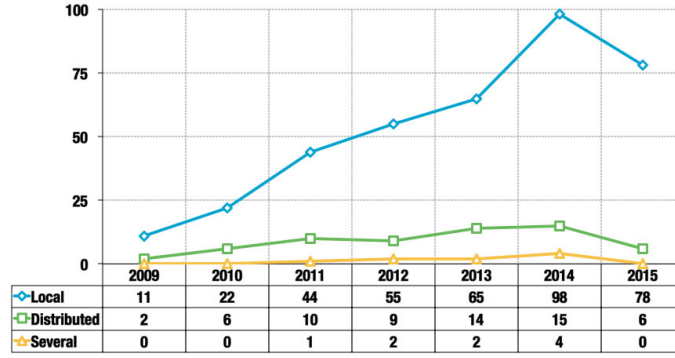


Figure 12 – Annual trend per tracking platform.

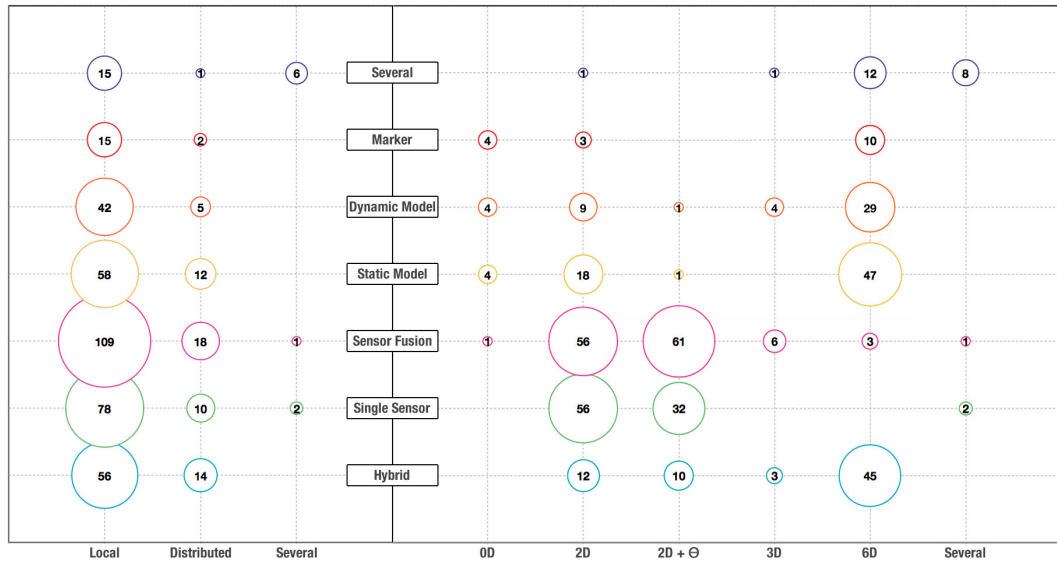


Figure 13 – Two dimensional bubble chart: left side presents the tracking type by tracking platform and the right side presents the tracking type by degree of freedom.

present a system that computes a 6D pose using only a combination of the device's sensors. One example is (ROBINSON et al., 2012), in which the authors append a pico projector to a mobile device in order to make projective drawings on the wall. The approximated position is computed in a calibration step using the sensors, in which the user has to move the device according to a projected guide. All 112 2D works that use sensors to compute the pose are location-based systems, as well as all three marker papers, such as (RAJ; TOLETY; IMMACULATE, 2013), seven of the hybrid works, like (SANTOS; RODRIGUES; OLIVEIRA, 2013), and four of the static model studies, such as (NGUYEN; ANDERSEN; HOILUND, 2010). All the other 29 works compute a 2D position on the screen, as in (AN; HONG, 2011).

Single sensor and sensor fusion systems are the only two tracking types in which a 6D pose is not the most common information required. For all other tracking types at least 58.8% of the papers are about a system that calculates a full rotation and translation pose, as exemplified in (ISSARTEL; GUENIAT; AMMI, 2014; HERLING; BROLL, 2012; PIRCHHEIM; SCHMALSTIEG; REITMAYR, 2013; KURZ; BENHIMANE, 2012). Additionally, the dynamic model and sensor fusion techniques are the only tracking type that has at least one paper

for every degree of freedom.

The bubble chart in Figure 14 presents the same dimensions of Figure 13. The difference is that it combines the vision-based and sensor-based techniques. Additionally, it also combines both location-based solutions. These combinations make more evident that most of the vision-based techniques calculate a 6D pose and that the majority of the sensor-based approaches are location-based services. Regarding the first dimension, it is possible to see that the ratio between the number of local and distributed solutions for each tracking type stays almost the same.

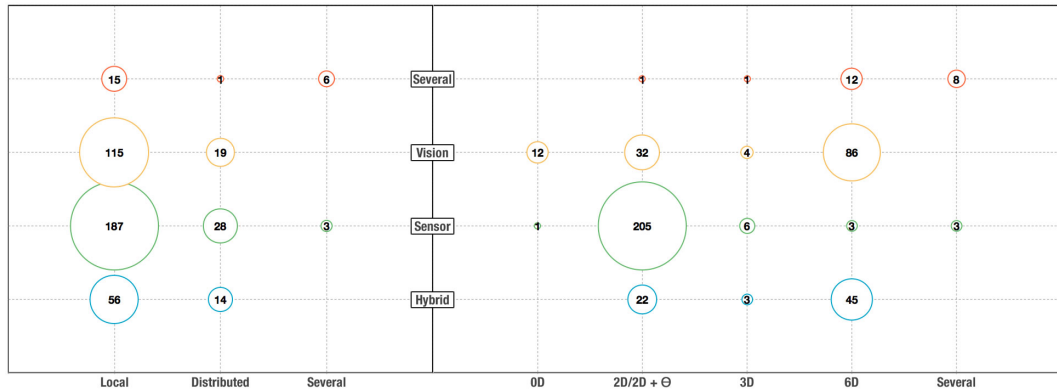


Figure 14 – Two dimensional bubble chart: left side presents the combined tracking type by tracking platform and the right side presents the combined tracking type by degree of freedom with location service systems combined.

The relationship of degree of freedom by tracking platform is shown in Figure 15. The chart shows that for every degree of freedom category more than 80% of the publications are local. Moreover, all works that compute a 0D detection are local, such as (ULLAH et al., 2012).

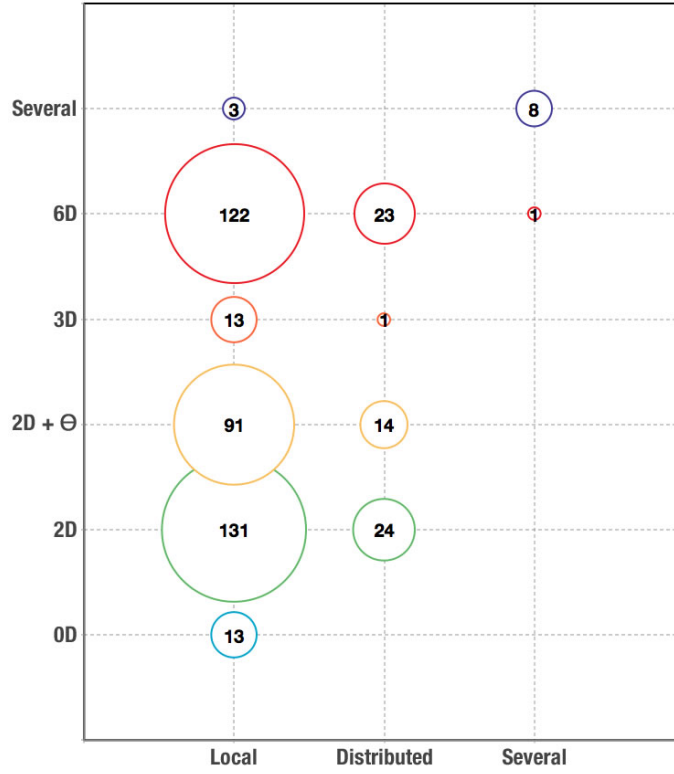


Figure 15 – One dimensional bubble chart of degree of freedom by tracking platform.

2.4 DISCUSSION

Figure 5 shows that, even though the amount of publications in 2015 was smaller than in 2014, overall the number of papers about tracking on mobile devices is increasing over the years. There were almost 6.5 more works in 2015 than in 2009 and it is due to the improvement (HALPERN; ZHU; REDDI, 2016) and popularization (Pew Research Center, 2016) of such devices in recent years.

It is possible to see in Figure 10 (top) that there was an increase of more than two and a half times in the number of publications for all tracking types between 2009 and 2015. There is also a growth of vision-based works, as shown in Figure 10 (bottom). These data indicate that this type of tracking becomes possible with the improvement of the computational power of devices, especially for natural feature tracking.

Figure 10 also shows that sensor fusion tracking had the biggest growth in the analyzed period. It is also possible to note in Figure 6 that the majority of works use this type of tracking. Moreover, 49.1% of all studies do not rely on the camera to perform tracking. This is probably because it is the most suitable approach to compute a pose for location-based solutions, which use 2D and 2D + θ information, and this type of solution is one of the most common type of application for mobile devices. This relationship is emphasized in Figure 13 and Figure 14. Nevertheless, it should be observed that only three works use only sensors to compute a full 6D pose because of their technical limitation, such as noise and error accumulation.

The analysis revealed that 41.3% of all sensor papers use data from only one sensor to compute the device's pose. The other 58.7% perform tracking using a combination of different sensors. This fusion of sensors is important because it allows using the data from one sensor to overcome the weakness of another one. Moreover, all studies that use a single sensor are 2D or $2D + \theta$, as can be seen in Figure 13. The nine papers that use sensors to compute a 3D or 6D pose require a combination of them in order to perform tracking. As seen in Table 2, the two most common sensors, accelerometer and magnetometer, are popular because they can be used in combination with other sensors in both indoor and outdoor situations since they do not require any external infrastructure, such as access points. The two sensors that follow them are related to providing the device's position. Wi-fi is widely used to compute indoor position. Although noisy, GPS is a great way to determine outdoor localization.

Figure 14 shows a clear trend that relates sensor techniques to location-based systems and vision-based approaches to solutions that require a 6D pose. Moreover, it is possible to see in Figure 6 that static model tracking is the favorite among natural feature-based approaches. However, there is a significant amount of systems that use a dynamic model technique. One reason is that some works use learning algorithms to calculate a pose, such as (LAMBRECHT; WALZEL; KRUGER, 2013). These techniques demand a massive processing power in the offline training phase that can be performed previously in a computer but does not use much processing for tracking, which makes them more suitable for mobile devices. One advantage of these techniques is that the trained model is usually refined using the tracking results.

There was a huge increase in the number of publications of $2D + \theta$ works, which had a growth of 22 times between 2009 and 2015, as shown in Figure 11. Papers with 2D and 6D systems grew almost 7 and 4 times in the same period, respectively, which is also a considerable value. All the other categories present a consistently small number of papers throughout the years. Moreover, Figure 7 shows that 2D techniques are the most common ones. The main reason for this is that there is a high demand for location-based applications and the amount of publications reflect this. 6D systems are also very popular because it is the most traditional information required for tracking, especially for augmented reality systems. The fact that the majority of papers of every tracking type except single sensor and sensor fusion calculate a 6D pose reflects the importance of computing the rotation and translation of the device relative to the real world, as observed in Figure 13. This approach is interesting because it combines the benefits of both vision and sensors to perform a more accurate and robust tracking. It is also possible to see in Figure 7 that 0D and 3D approaches are the least common. This is due to the fact that there is a small number of applications that require a 0D or 3D pose.

Even with the devices' limitations, Figure 8 shows that the majority of works execute all the steps to compute the pose at the device. One reason is the lack of a good communication

infrastructure to transfer the data to a remote server. However, it can be noted in Figure 12 that more than half of the distributed works were published in the last three years. This can be an indication that there is a recent improvement in the network infrastructure and researchers are exploring the use of a remote computer, which provides more resources than the mobile device, such as processing power, memory and storage space. Another reason is the possibility of using sensors and cameras that are not available in the device (BAI; GAO; BILLINGHURST, 2013).

Figure 13 indicates that there is no relationship between the tracking type and the execution platform since the proportion of local and distributed works varies little per tracking type. However, the same proportion is not seen when relating tracking platform with degree of freedom. It is possible to see in Figure 15 that a few more than 80% of 6D papers are local while none 0D and only one 3D study is distributed.

As seen in Figure 9, the majority of the papers propose a new technique and there are also several works that use an existing method to create a solution for an open problem. These two research types represent 93% of all studies classified. This is an indication that the demand for systems that use tracking is high. More than that, it is a clear suggestion that the field of tracking for mobile devices still has a lot of open problems to tackle.

2.4.1 Implications for Future Studies

This mapping study not only offers useful information for researchers who are interested in the existing works regarding tracking for mobile devices but also identifies gaps in this research topic.

Most of the works calculate the pose using devices' sensors or computer vision algorithms. However, there is a tendency to combine both approaches to provide a more robust tracking. It can be noticed in Figure 10 that no hybrid work was published in 2009 and in 2015 the 20 papers that use this type of tracking represent 23.8% of the studies in this category. It is the highest percentage in the evaluated period. One reason is the improvement and miniaturization of more complex sensors and cameras. Following this tendency, Apple and Google launched their augmented reality platform at the end of 2017, namely ARKit (Apple Inc., 2017) and ARCore (Google Inc., 2017), respectively. The tracking algorithm on both SDK relies on hybrid techniques. They use the camera image to extract and track natural features and to perform loop closure, while the inertial sensors are responsible to recursively update the device position.

Another sensor that achieved this level of maturity regarding miniaturization and integration into mobile devices is the depth camera. Google's Tango is another augmented reality platform that combines different types of sensor to perform tracking (Google Inc., 2014). However, different from ARKit and ARCore that only use sensors available in most mobile devices, Tango also incorporates a depth and a fisheye camera. These two cameras play an important role in the development of accurate and robust tracking techniques.

As a result, other kinds of devices, such as DAQRI Smart Helmet (DAQRI, 2016) and Microsoft Hololens (Microsoft, 2016), are embedding depth sensors to improve tracking.

This mapping found a few studies focusing on the use of machine learning approaches to compute the pose. But this is a prominent research area because such algorithms learn what are the best features to be used for tracking (TAN; ILIC, 2014). Moreover, as mentioned before, learning techniques transfer most of the computational effort to an offline training phase while the tracking itself demands few processing resources, which makes them suitable for mobile devices. Machine learning is a mature area and its use for tracking is increasing rapidly. Although, there are still several open problems in the area.

Recent improvements in communication networks enable the increasing number of works that use distributed approaches, as shown in this study. In the future, this infrastructure will probably be more reliable and faster (The Next Generation Mobile Networks Ltd., 2015), which creates new opportunities to perform tracking on remote servers, using the mobile device only to capture the input information and display the output results. Moreover, this connectivity is a basic requirement for creating sensitive environments using smart objects. Each one of these connected sensors can be aware of its location and may be used to share this information with a mobile device and perform fully distributed tracking. For instance, smart objects spread over an indoor place can be used to provide or improve the indoor localization of a person using a smartphone connected with them.

It is also important to be aware of the improvements of the hardware capabilities that will be available on mobile devices in near future. New tracking techniques can be proposed or existing ones adapted taking into consideration the use of multiple cores of the device's processor and graphics processing unit (GPU). Beyond that, it is possible that several mobile devices will have chips dedicated exclusively to execute embedded computer vision algorithms, such as Qualcomm's Hexagon digital signal processor (DSP) (Qualcomm Technologies, Inc., 2015). These dedicated chips will allow tracking to be performed faster while consuming less energy.

A research topic that did not appear in the mapping was the inclusion of any kind of semantic in tracking systems, even though it is an important research problem in computer vision. One reason could be that extracting any type of knowledge from the scene demands more resources than a mobile device can handle in a practical time at that point. However, recent tracking techniques for augmented reality, such as ARCore and ARKit are able to identify planes to anchor the augmented content, which can be the floor, a table or a wall. With this simple knowledge of the environment, it is possible to improve rendering, such as adding more natural shadows.

3 PRELIMINARY EXPERIMENTS

The systematic mapping provided valuable theoretical knowledge regarding tracking for mobile devices. These lessons learned were used to elaborate a set of preliminary experimental scenarios aiming to obtain a practical know-how about the specificities of developing tracking techniques for mobile devices. This chapter describes these experiments and discusses their results and findings. The first one evaluates the Google Tango platform to establish a reference of the state-of-the-art trackers. Additionally, another experiment tests the use of parallelism, distributed approach and native implementation in order to find an efficient architecture to execute computer vision algorithms on mobile devices. After that, it is shown the experiments to test different tracking techniques that the systematic mapping indicated to be suitable for mobile devices. One is a face tracking technique using machine learning and local binary features. The other one is a SLAM technique that was developed in desktop and ported to a Tango tablet device.

3.1 EVALUATION OF TANGO PLATFORM

Tango (Google Inc., 2014) is a platform that combines computer vision techniques with state-of-the-art sensors, allowing to perform 6DoF tracking on mobile devices. Even with the release of ARCore and ARKit, Tango is still arguably the best tracker available for them. Mainly because it uses cameras that are not available on most of the handheld devices, such as depth and fisheye ones. While the fisheye camera widens the field of view, which increases the number of characteristics, the depth camera provides a good approximation of the scene structure and scale. These additional data help achieving a more precise and stable tracking.

Therefore, it is relevant to evaluate its accuracy. As far as the author knows, so far no study evaluates the Tango platform. This experiment is important because it provides expertise on how to assess tracking systems and also to find scenarios and situations in which they work well and fail. This section aims to evaluate the precision of both motion tracking and depth sensing technologies of the Tango platform. The evaluation method proposed and the results were discussed in (ROBERTO et al., 2016a), and details are provided in the next subsections.

3.1.1 Tango Platform

As mentioned, the devices that support the Tango platform have some features that provide new ways to navigate in different environments. They use technologies such as motion tracking and depth perception. This subsection explains how these technologies work on the Tango platform as well as describes the characteristics of the sensors available.

Motion tracking means that a Tango device can track its own movement and orientation through 3D space. Tango implements motion tracking using visual-inertial odometry to estimate where a device is relative to where it started. Standard visual odometry uses camera images to determine a change in position by looking at the relative position of different features in those images. Visual-inertial supplements visual odometry with inertial motion sensors capable of tracking a device's rotation and acceleration. This allows a Tango device to estimate both its orientation and movement within a 3D space with even greater accuracy. Unlike GPS, motion tracking using visual-inertial odometry works indoors. In addition to the gyroscope and accelerometers, Tango uses a wide-angle motion tracking camera to add visual information, which helps to estimate rotation and linear acceleration more accurately. To perform motion tracking, the Tango APIs provide the position and orientation of the user's device in full six degrees of freedom. The data is returned with two main parts: a vector in meters for translation and a quaternion for rotation.

Depth Perception gives an application the ability to understand the distance to objects in the real world. Tango tablet implements depth perception with time-of-flight (ToF) technology, which requires the use of an infrared projector and an infrared sensor. Tango tablet depth sensor is designed to work best indoors at moderate distances (0.5 to 4 meters). This configuration gives good depth at a distance while balancing power requirements for infrared illumination and depth processing. It may not be ideal for close-range object scanning or gesture detection.

The Tango APIs provide a function to get depth data in the form of a point cloud. This format gives (x, y, z) coordinates for as many points in the scene as are possible to calculate. Each dimension is a floating point value recording the position of each point in meters in the coordinate frame of the depth-sensing camera.

3.1.2 Evaluation Methodology

Evaluating tracking systems is a challenging task (TAMURA; KATO, 2009). Many efforts have been made in the past years to provide metrics and standards to analyze the aspects related to this problem (PETIT et al., 2011; ROY et al., 2015; HODAN; MATAS; OBDRŽÁLEK, 2016). The main reason is the difficulty to find the ground truth to compare with the obtained results. There are several benchmark datasets available aiming computer vision tracking systems evaluation, each one having many purposes and providing several types of input data. For instance, (LIEBERKNECHT et al., 2009) presents an image dataset to evaluate planar model-based techniques, (SHIBATA et al., 2010) describes a benchmark to measure the quality of 3D model-based algorithms and (STURM et al., 2012) provides RGB and depth information aiming SLAM systems. All of them also provide the expected results, which are used to assert the precision of the algorithm. However, it is hard to use these datasets on mobile devices. One reason is the difficulty to extract information from

different sensors since they are noisy and the precision varies among devices. Therefore, it was proposed a different methodology to evaluate the motion tracking and depth sensing functionalities in the Tango Platform. The selected device was the Yellowstone tablet, which provides all the Tango functionalities.

3.1.2.1 Motion Tracking

The evaluation method is based on moving the Yellowstone tablet between two known positions in the real world. Then, comparing the distance between these positions estimated using the device with the ground truth value. This way, it is possible to evaluate the error the system accumulates during motion tracking from a starting point to an ending position. It was used a graph paper with the precision of one millimeter to ensure the experiment accuracy. The paper was glued to a table so it does not move during the tests. A needle was attached to the base of the Yellowstone tablet in order to have the exact position of the device over the paper. Figure 16 shows the setup. It was designed two different experiments based on this setup, one to evaluate a small augmented reality workspace and another one for large environments.

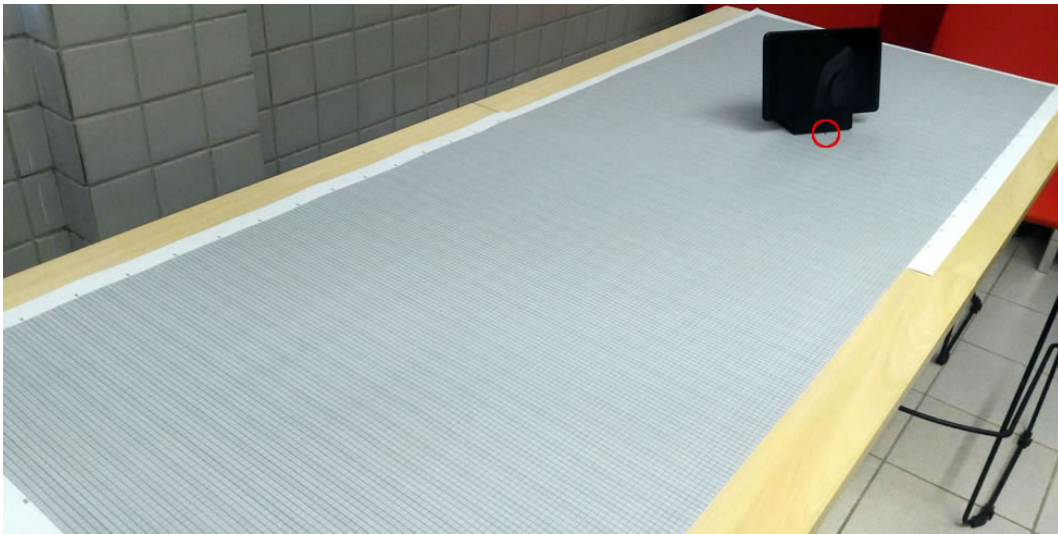


Figure 16 – Evaluation setup consists of a graph paper with precision of one millimeter and measuring 1.5 x 0.55 meters. Red circle highlights the needle used to get the exact position on the paper.

For the first one, the idea is to evaluate how the Yellowstone tablet works on a small workspace, which for this study is a table with an area up to one square meter. Therefore, the device is positioned with the needle on the origin of the graph paper and their axes are aligned. The tablet will be moved freely and placed in any other position on the graph paper. The error is the difference of the Euclidean distances between the origin and the final position computed on the Tango device and measured with the graph paper.

Regarding the large environment, the goal is to measure how the Tango device behaves when performing motion tracking on places such as a regular office and outdoors. The

office is a closed room that has an approximate area of 50 square meters and artificial illumination. As for the outdoor experiment, it was placed in two different courtyards measuring around 100 square meters each. It was also located in the corridor of a building that is open to the outside. All the outdoor measurements were collected using natural illumination during daylight.

It is not possible to have a graph paper that is large enough to cover the entire office or the courtyards. Thus, for this experiment, the device is also positioned with the needle on the origin of the graph paper, their axes are aligned and it is moved freely in these environments. The difference from the previous experiment is that the tablet is returned to the same position where it started. The error is also the Euclidean distance between the position computed on the Tango device after finishing the movement and the initial position.

3.1.2.2 Depth Sensing

Regarding the depth sensor, it was evaluated the accuracy of the 3D points positions obtained from it. The process consists of calculating the Euclidean distance between 3D points reconstructed from the color camera and corresponding ones from the depth camera. The Tango API provides the registration between color and depth cameras. The intrinsic parameters used are the ones from the manufacturer calibration. The 3D points of the color camera are the inner corners of a detected chessboard pattern whose pose is estimated using the Direct Linear Transformation (DLT) method and refined by minimization of reprojection error (HARTLEY; ZISSERMAN, 2003).

Since the depth image generated by the Tango device has a lower resolution when compared to the color image, it has to be upsampled when obtaining the corresponding depth measure of a chessboard corner. Both nearest-neighbor and bilateral interpolation (TOMASI; MANDUCHI, 1998) were evaluated for performing this task.

3.1.3 Results

In order to evaluate the motion tracking capability of the Tango platform, it was used a sample application available on the project GitHub¹ that uses both Tango Area Learning and Motion Capture features and is called “C++ Augmented Reality Example”. This application uses the fisheye camera and the device gyroscope and accelerometer to compute its pose relative to its initial position.

3.1.3.1 Motion Tracking

For the experiment on the small workspace, the Yellowstone tablet was moved freely to any other position over the graph paper. To have statistical power, the sample size for

¹ <<https://github.com/googlesamples/tango-examples-c>>

this experiment was calculated aiming 95% confidence within 1 centimeter precision (JAIN, 1991). Therefore, these measurements were repeated 67 times to ensure that.

Figure 17 shows the error dispersion, in which the smallest was 0.009 meters and the largest was 0.181 meters. On average, the error was 0.067 ± 0.040 meters.

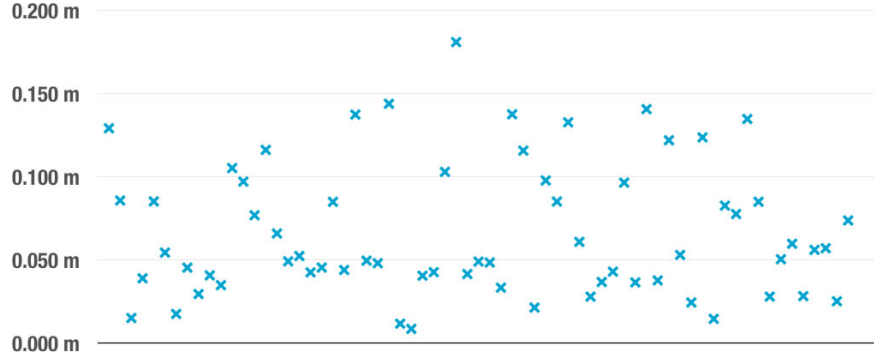


Figure 17 – Error dispersion for the small workspace experiment.

Figure 18 shows the device's position distribution over the graph paper during the experiment.

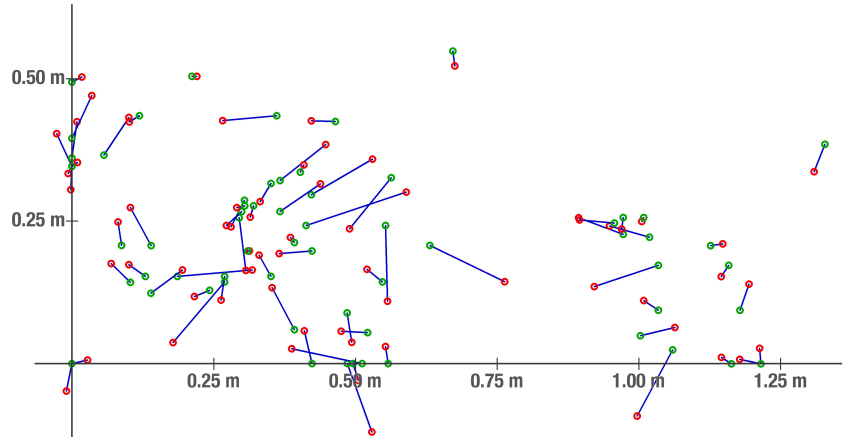


Figure 18 – Distribution of the device positions on the graph paper (green) and their correspondent positions calculated by the Yellowstone tablet (red).

Regarding the evaluation of the motion tracking on large environments, the error is the Euclidean distance between the initial and the final position calculated by the device after moving it freely in this environment and returning to the same location. After a few measurements, it was noted a large difference among the errors from the indoor and outdoor environments. Therefore, it was decided to perform two different evaluations, one for each situation. The number of samples to ensure statistical power emphasize this decision. For the indoor experiments, it was necessary to have 31 samples to have 95% confidence within 2 centimeters precision, which is the double of the small workspace because the covered area was much larger. On the other hand, it was not possible to have such confidence in the outdoor experience. The reason is that the error variation is so high

that it would be necessary to have more than 5000 samples to have 95% of confidence within 2 centimeters precision.

Figure 19 shows the error dispersion in the large indoor scenario. The smallest one was 0.049 meters while the largest was 0.261 meters. On average, the error of the 45 samples measured was 0.142 ± 0.057 meters. Also, the average distance walked with the Yellowstone tablet was 23.608 ± 7.892 meters.

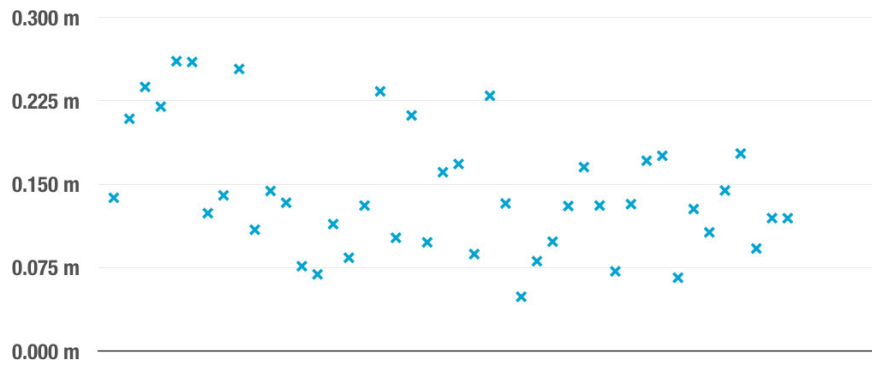


Figure 19 – Error dispersion for the large indoor environment experiment.

Figure 20 illustrates one of the paths walked with the Yellowstone tablet and the difference between the initial and final positions.

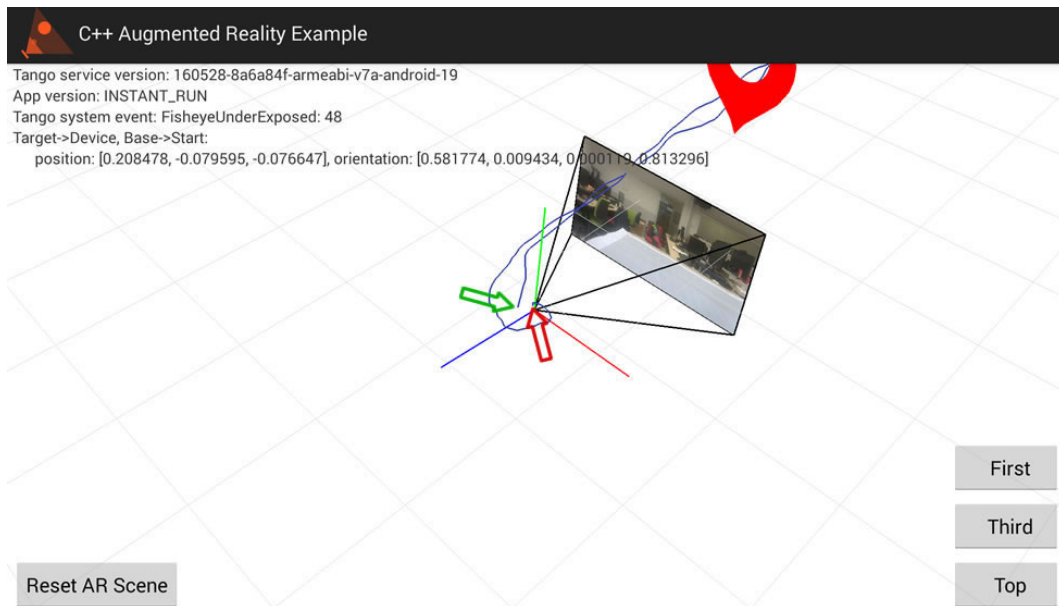


Figure 20 – Screenshot of one of the paths computed using the Yellowstone tablet. The green arrow points to the initial place and the red one to the final position calculated after a free walk. The error is the average Euclidean distance between them.

Regarding large outdoor environments, it was performed 21 repetitions. However, this amount was not enough to have statistical power. Although, the average error of 0.905 ± 0.753 meters indicates that precision of the Yellowstone tablet is much smaller when it is dealing with natural illumination and wide spaces.

3.1.3.2 Depth Sensing

In the first experiments, the chessboard pattern was printed on a paper using black ink. However, the dark squares did not reflect infrared light in a way that would allow robustly estimating the depth of the chessboard corners. Due to this, it was used a mix of cyan, magenta and yellow ink in order to have dark squares that are correctly scanned by the depth camera. This aspect is illustrated in Figure 21. The chessboard was printed on A4 paper with a square side of 28 millimeters.

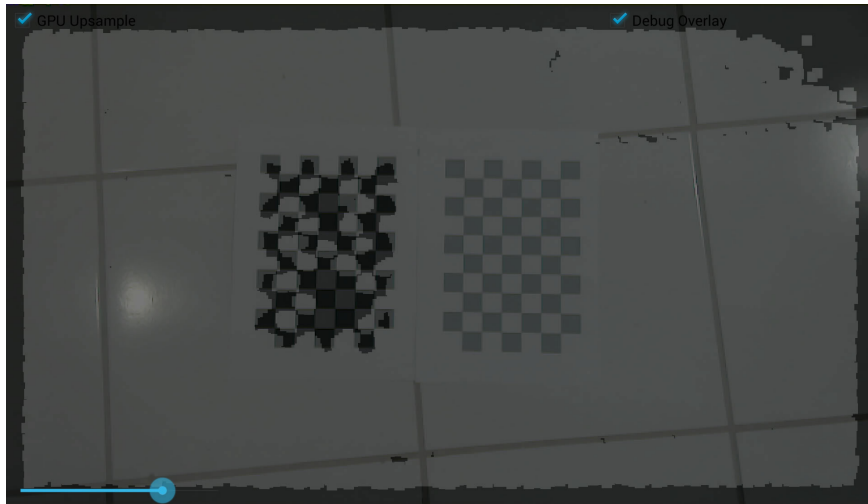


Figure 21 – Screenshot of the depth estimation of two chessboards printed on a paper. On the right, the one printed with a mix of cyan, magenta and yellow inks. On the left, the same pattern printed with a black ink. Note that the sensor is not able to estimate depth on the black squares of the left paper.

Figure 22 shows the mean depth estimation error considering different distances between the device and the chessboard pattern and different depth interpolation strategies. In order to obtain error values accurate within 0.1 millimeters at 95% confidence, 150 samples were collected for each configuration. The average execution times of the nearest-neighbor and bilateral depth interpolation procedures for each point were 0.696 ± 0.127 and 1.881 ± 0.312 milliseconds, respectively.

3.1.4 Discussion

The results showed that the motion tracking of the Yellowstone tablet is 2.3 times more precise on a small workspace than on large indoor environments. However, the presented errors can have an impact on the user experience. An average error of 6 centimeters is often noticed in a small workspace.

On the other hand, even having a bigger error when dealing with large environments, the precision of the motion tracking on indoor spaces is suitable to provide a good user experience for several kinds of augmented reality applications. However, for scenarios in which it is necessary to have accuracy, this error can harm the user experience. Figure 23

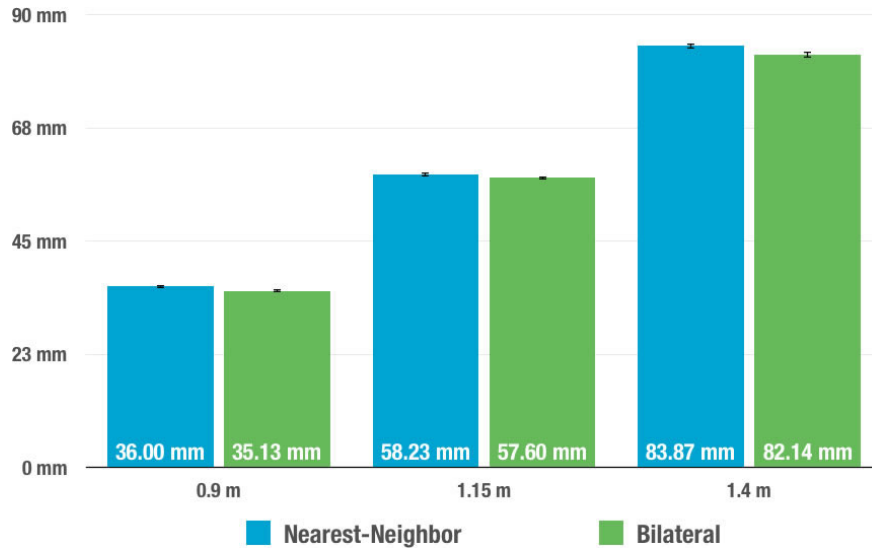


Figure 22 – Mean depth estimation error with respect to distance between Tango device and chessboard pattern using different depth interpolation methods.

(left) shows an example where Yellowstone tablet measure tool is calculating the width of a 0.7 meter door. After moving the device for a few steps away from the door and back, the ruler is placed in a different position, as seen on the right side.

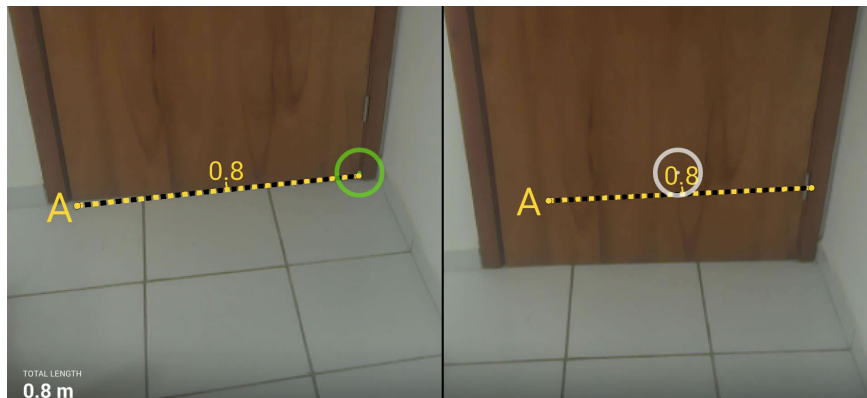


Figure 23 – Door width estimation using the Yellowstone tablet. Left side shows the initial measurement and the right side shows the ruler position after moving the device. Door actual width is 0.7 meters.

During the experiment, it was noticed that the algorithm that Tango uses for motion tracking seems to mistrust the sensors it uses regarding their precision. There is an indication that it has a much stronger confidence in the information provided by the fisheye camera. For instance, sometimes the device was left standing still over the table and when some object moves in front of the camera, the motion tracking algorithm calculates that the tablet was moving in the opposite direction.

The tests on large outdoor environments could not provide results with statistical power because there was a significant variation on the error measured on every sample. However, this disparity suggests that the Yellowstone tablet has some issues to deal with

outdoor illumination and wide spaces. It is emphasized by the fact that no result of the outdoor measurements presented a smaller error than any indoor sample. Moreover, in some cases, the error was greater than 1.0 meter and in the worst case it reached more than 3.0 meters.

Regarding depth sensing, the mean errors presented by Yellowstone tablet were similar to the ones obtained with a desktop depth sensor. The results also suggest that the depth estimation error increases linearly with respect to the distance between the device and the object. Bilateral depth interpolation provided an average precision improvement of 1.82% (1.08 millimeters) with respect to the nearest-neighbor approach. However, it was more than 2.5 times (1.18 milliseconds) slower on average for each point.

3.2 EFFICIENT TRACKING ON MOBILE DEVICES

Computer vision tracking is a task that demands many computational resources. It may be necessary a lot of processing and memory to extract and describe natural characteristics from a scene or to correctly match them with the features from a model. As the tracking system runs for an extended period of time, the number of features may grow, requiring more memory space, both ROM and RAM.

Since mobile devices have limited resources, it is necessary to investigate the best approach to minimize these constraints. This section aims to explore different ways to implement a computer vision tracking system in order to find the best trade-off between processing and memory. The results found in this experiment were published in (LIMA et al., 2015), and details are given in the next subsection.

3.2.1 Android Architectures for Computer Vision Tracking

After analyzing some possibilities to implement an Android system to perform computer vision tracking, two of them appeared to have the potential to combine high performance with low memory consumption. The first one is a Multi-Thread Partial Native implementation, which benefits from using the different cores of the device processor as well as the performance gain of having native code for some tasks, which is known to be faster than Java implementations. The second one is a Client/Server approach that uses the processing power of a remote server to perform the most demanding processing tasks. During their development, another method aroused, which leaves only the necessary functions on the Java side while transfer most of the computation to the native part of the code and is called Full Native implementation.

A simple static model tracking technique was designed in order to provide a quick prototype of these approaches. The first step is to build the model. In order to do so, ORB features (RUBLEE et al., 2011) are extracted from a loaded image that will work as a planar template. After that, a z-coordinate is given to the features, so each one now has a

corresponding 3D point. Then, it is computed their ORB descriptors. The planar static model that will be tracked is represented by the data structure that stores the descriptors and the 3D points of the features extracted from the image.

After creating the static model, it is possible to start tracking this template. First, ORB features are extracted from the frame captured by the device's camera and their descriptors are computed. Then, they are matched with the model's descriptors using a nearest neighbor search approach that finds the Hamming distance between descriptors to determine the best correspondence candidates between these groups of features. However, it is common to have several wrong matches and it is necessary to filter them to improve accuracy. One effective way to do that is by comparing the distance between the closest and the second closest neighbors. This strategy works because correct matches have the nearest neighbor much closer than the second closest one, which is incorrect (LOWE, 2004). After filtering the good matches, the camera rotation and translation are estimated using EPnP (LEPETIT; MORENO-NOGUER; FUA, 2009). RANSAC (FISCHLER; BOLLES, 1981) is also used to reduce the influence of outliers, which are spurious matches that could also remain after the filtering process. This process is repeated until the application finishes.

3.2.1.1 Multi-Thread Partial Native Implementation

This implementation parallelizes part of the tracking steps described above, which are: feature extraction and description, matching with model keypoints, and filtering for good matches.

After building the static model, the tracker enters the main loop. This implementation uses a port of OpenCV (BRADSKI, 2000) to this mobile platform named OpenCV4Android. The library provides the data structures and computer vision functions to perform the tracking task. OpenCV4Android is an interface that provides access to almost every OpenCV function in Java. Additionally, the original C++ version is also available to be used on native developments.

This implementation uses both versions. First, the current frame is captured using a camera interface provided by OpenCV and a Mat structure stores it. The Native part of the system receives this image, which is divided into equal parts according to the number of cores available in the device. Each core will process one part of the frame in parallel. Intel TBB library (Intel Corporation, 2016) has a version that is compatible with Android and was used on this version of the tracker. Then, in parallel, each core extracts features from one part of the image, matches them with the model's keypoints and filters to select the good matches. After every core finishes their tasks, the good matches are combined in only one structure and it is used to compute the device rotation and translation. The augmented reality content is drawn over the original image, which is returned to the Java part of the code so the rendering structure provided by OpenCV can display the modified frame on the screen. Figure 24 illustrates this implementation.

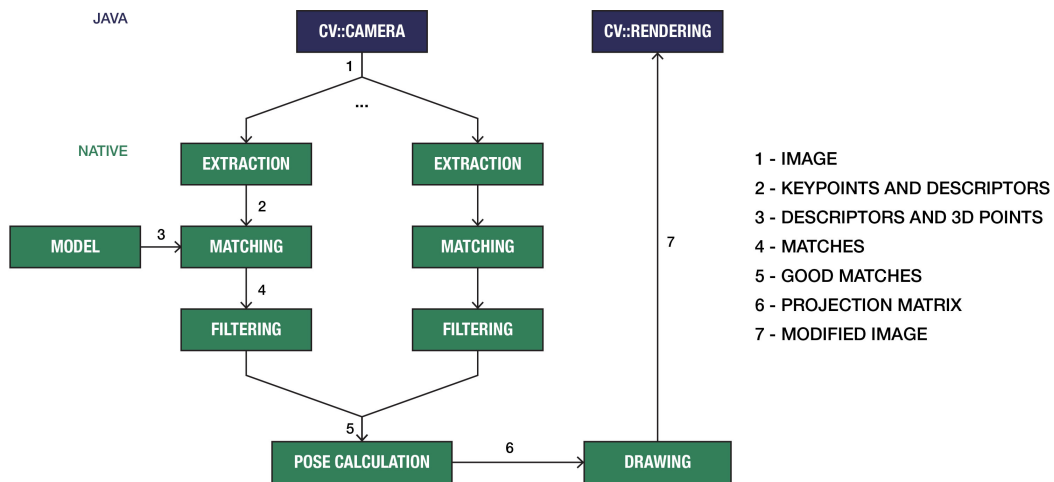


Figure 24 – Multi-Thread Partial Native flow diagram.

3.2.1.2 Client/Server Implementation

This architecture aims to transfer part of the tracking processing to a remote server, which has more computational power to perform this task quicker than on the mobile device. Several tasks can be assigned to be executed at the server. On one hand, it is possible to implement all the tracking pipeline remotely, but it will demand an excellent network infrastructure to transfer the frames in real-time on both directions. On the other hand, the server can execute only selected tasks, which will decrease the network requirements. However, the server processing and memory resources will not be fully exploited.

For this implementation, every tracking step that depends on the image will be processed on the device in order to decrease the dependency of a good network infrastructure. Therefore, as seen in Figure 25, the current frame is captured using a camera interface provided by OpenCV and ORB is used to extract features and compute their descriptors. These information are encapsulated and sent to the server through a wireless local network of 54 Mbps, which is 35% faster than the fastest 4G network available (OpenSignal, Inc, 2016). The server receives this data, which is hundreds of times smaller than the full image, and matches it with the model sent to the server during the initialization. Then, it filters these matches to select the good ones and computes the object rotation and translation. These two vectors are sent back to the device that draws the augmented content over the captured frame and renders it on the screen.

3.2.1.3 Full Native Implementation

It has become clear during the development of these approaches that capturing and rendering a frame using OpenCV functions is not efficient. The alternative found was to use camera structures from Android to capture a frame buffer and render it with OpenGL (Khronos Group, 1997).

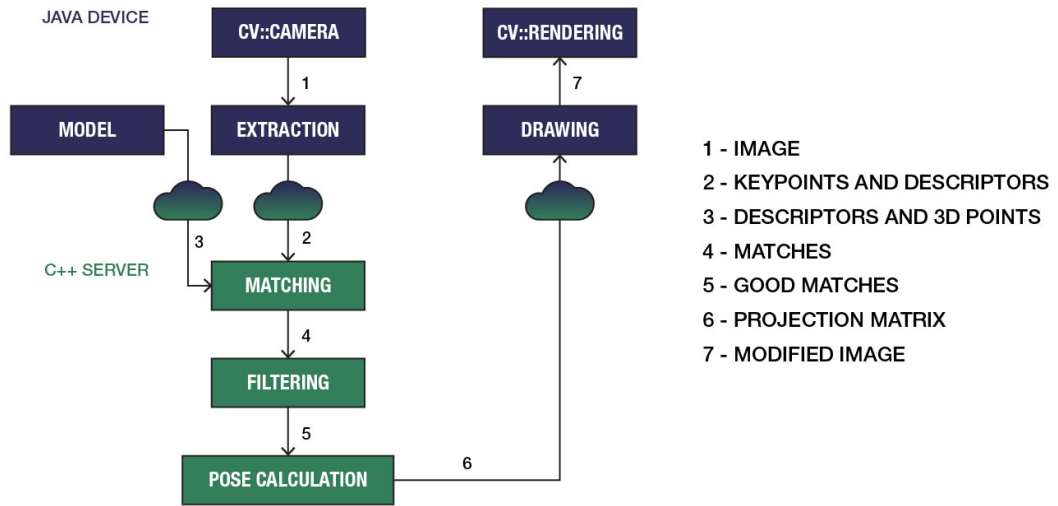


Figure 25 – Client/Server flow diagram.

The problem with this approach is that Android camera and rendering structures work only with YUV color space while the OpenCV library uses the RGB color space (KUEHNI, 2003). In other words, each frame buffer received from the Android camera preview is in the NV21 format, which is the standard picture format on Android camera preview. This way, it was necessary to convert the color space of the image buffer on each frame, and a better approach to do so was to implement a parallel loop rather than using a sequential one. Additionally, it was necessary to implement methods that use OpenCV functions to render on YUV images like it was an RGB frame.

This way, every part of the tracker that uses OpenCV is implemented with native code. Therefore, the frame is captured, stored in a byte array and then sent using Java Native Interface (JNI). The YUV frame buffer is converted to RGB to be tracked sequentially. The augmented content is drawn and the image is converted to a byte array so it could be sent to the Java layer in order to be rendered, as shown in Figure 26.

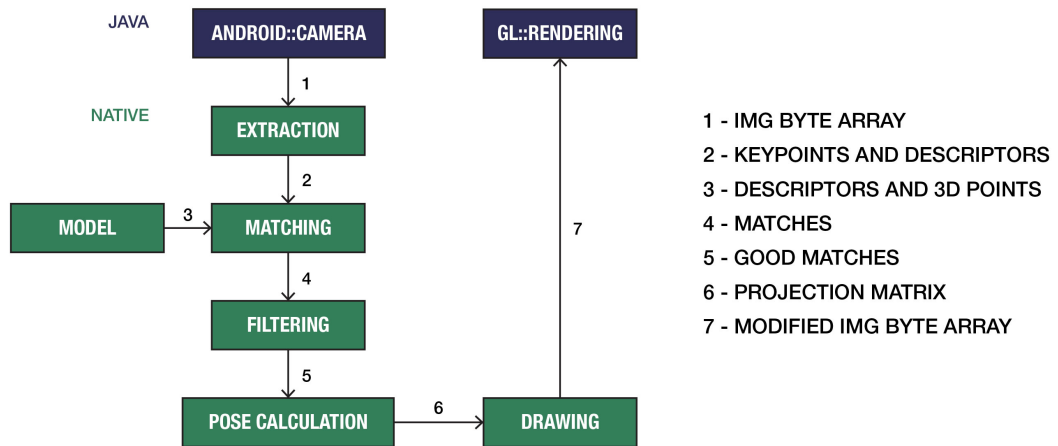


Figure 26 – Full Native flow diagram.

Another benefit of this architecture is that it eliminates the necessity of having OpenCV Manager. This software has to be downloaded from Google Play Store and it provides all the dependencies necessary to execute applications that use the Java version of OpenCV4Android. When programming only with the C++ part, Android compiles all the dependencies and embeds them in the final application package.

3.2.2 Architectures Evaluation

The solutions were compared regarding performance and memory used, both RAM and ROM. Additionally, it was also developed a solution in which all the tracking is performed in Java. This solution was set as a base for comparison with the other implementations. For all of them, the evaluation method consisted in tracking the static model in the real world to project virtual lines over the template board, as seen in Figure 27.

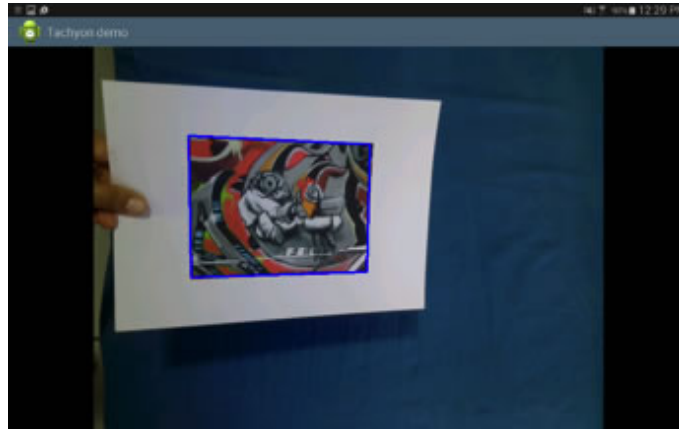


Figure 27 – Screen capture of the system tracking the template. The blue square is the virtual content that is placed on top of the model.

The device used for evaluation was a Samsung Galaxy Note 10.1, which has a 2.3GHz quad-core processor, 3GB and 32GB of RAM and ROM memory, respectively. Figure 28 shows the results of running the tracker with each implementation.

On average, Java implementation presented the worst results, being more than two times slower than the Multi-Thread version. It is clear that the bottleneck for every implementation is the feature extraction step. In the Full Native version, this stage consumes 83.3% of the total processing time. As expected, parallelizing this task provides a good speedup, more than three times when compared to a Java implementation.

In the other tracking phases, there is a significant speedup on using native code when compared with Java. However, parallelizing these tasks decreases their performance when comparing to sequential native implementation. The reason is that there is an overhead to divide and manage the threads. On the server, these tasks combined are executed in almost 2 milliseconds, but the time required to transfer the data to the server and back to the device eliminated all this gain. It is important to mention that the latency of local

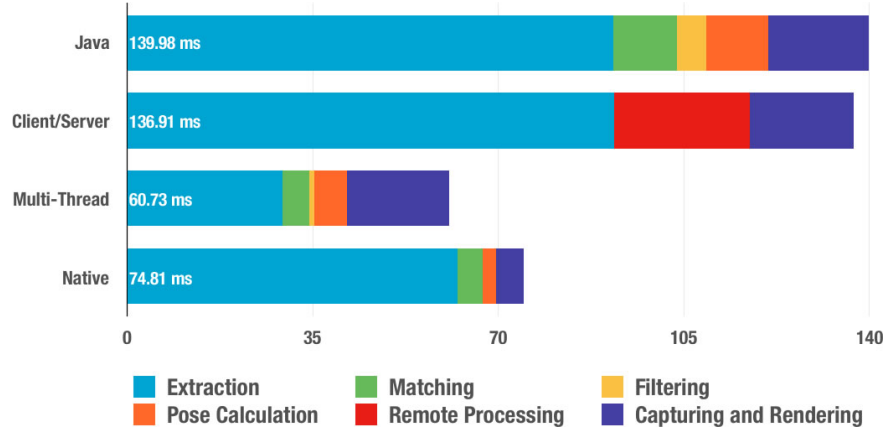


Figure 28 – Execution time in milliseconds on every tracking stage for each implemented architecture. For Client/Server approach, feature matching, matching filtering and pose calculation also includes the time to transfer the data to the server and back to the device.

networks is usually around ten milliseconds, which has a significant impact when dealing with real-time systems.

As expected, the Full Native implementation is able to capture and render the frame using Android functions much faster than the other approaches, which use OpenCV functions to do the same task.

Regarding memory consumption, it is possible to see on Figure 29 that the Client/Server implementation is the most efficient. On the other hand, Multi-Thread Partial Native is the one that requires more RAM to execute. The reason is that the first one deals with less information because part of the steps is processed on the server. However, the second one manipulates all the data on the device and it is also necessary to create new structures to manage the different threads, which impacts on memory consumption.

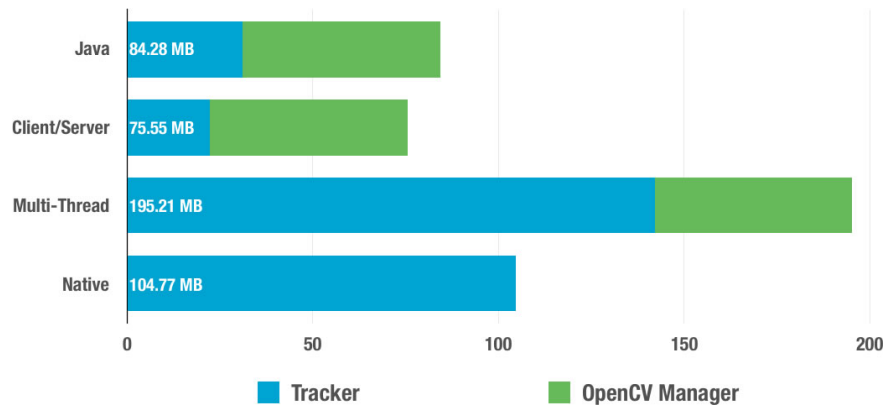


Figure 29 – RAM memory consumption in MB during the execution of each architecture implementation.

The Full Native implementation also has massive memory consumption when compared with the Client/Server approach. It is due the fact that this version includes all the OpenCV dependencies because it is fully implemented using native code. On the other

hand, the Client/Server version requires OpenCV Manager, which consumes twice more memory than the application itself.

As for the space required to store each application, both Client/Server and Multi-Thread Partial Native versions are much smaller than the Full Native implementation. One more time, the reason for that is that they do not embed the OpenCV dependencies because they access them from OpenCV Manager. Figure 30 shows that Full Native application is 7.55 times larger than Client/Server, but has less than half of its total size when OpenCV Manager is also taken into account.

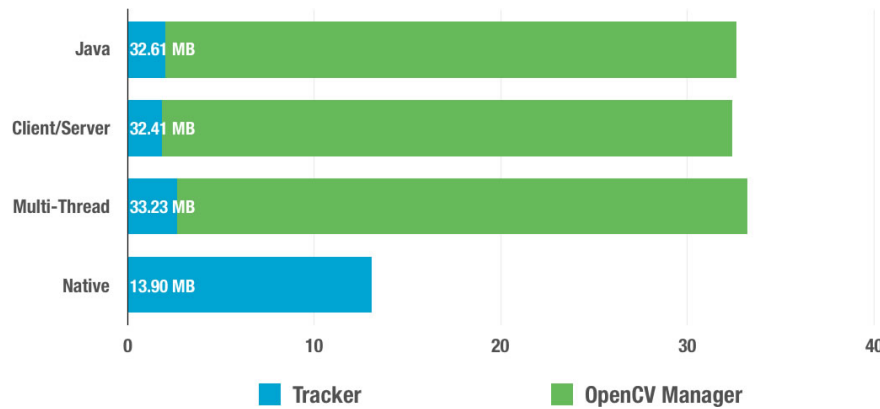


Figure 30 – ROM memory in MB required to storage each architecture implementation.

3.2.3 Discussion

The tests showed that the selection of the best architecture depends on different factors. For instance, when using a mobile device that has a lot of memory and multiple cores, the best approach would be the Multi-Thread Partial Native. However, in case a good network infrastructure is available for a device with few resources, the Client/Server implementation would be a good choice.

It is safe to say that, on average, the Full Native architecture is the most efficient between the evaluated alternatives. It is approximately 23% slower than Multi-Thread Partial Native, but uses about half of the memory. On the other hand, requires around 39% more RAM than Client/Server running more than 45% faster. Other advantage is that it does not require the installation of any additional application. Additionally, depending on the algorithm that will be developed, it is possible to identify parts that could be paralleled to make the system faster.

Perhaps the ideal architecture would be the mixing of all of them. A Full Native approach that supports multi-thread and is able to transfer tasks to a remote server whenever this alternative is suitable. Therefore, the system would constantly and automatically evaluate the resources available on the device as well as the quality of the network to decide which one has the best trade-off for a certain task. There are works that use this principle to select the fastest tracking technique according to the available resources that will achieve

the required level of accuracy (WAGNER; SCHMALSTIEG; BISCHOF, 2009; TEIXEIRA, 2013). Nevertheless, no works were found that evaluate the resources of a mobile device.

3.3 MACHINE LEARNING TRACKING ON MOBILE DEVICES

Because of the characteristics of machine learning approaches, which concentrate most of the computational effort in the offline training phase while the online tracker demands fewer resources, this type of technique can be very useful for tracking. For instance, it is possible to update the camera pose using a very small amount of information, but highly discriminative, which was inferred previously (JURIE; DHOME, 2002).

One application in which machine learning is making an impact is to track a body or its parts, such as face and hand. Face tracking, in particular, has several purposes on mobile devices. For example, it can be used to identify the face of the owner of the phone using the frontal camera to unlock the device or not (HADID et al., 2007). Another example is to provide a makeup tutorial using augmented reality (ALMEIDA et al., 2015). In this sense, it is important to identify not only the position of the face on the screen but also the location of relevant points in the eyes, nose, mouth and chin, called landmarks.

This section describes the research process to develop a real-time face tracking technique for mobile devices that, given an input image, is capable of identifying face landmarks.

3.3.1 Research Methodology

There are several face tracking techniques. Among them is the Constrained Local Model (CLM) (MORENCY, 2012) that performs a local search to find the position of the landmarks. Another one is the work described in (RAMANAN, 2012), which presents an idea similar to the CLM technique, but has better results. Although these techniques provide good precision, they are not optimized for mobile devices. Therefore, it was necessary to perform a literature review in order to determine if there was a more suitable approach.

The snowball sampling method was chosen (GIVEN, 2008), in which a good reference study is used as a seed to find other relevant works. In this sense, both the papers that are referred in the base study and the ones that cite them are gathered for further evaluation. For this study, the reference paper used was (RAMANAN, 2012).

The interest was in finding newer studies that improve the base technique. Therefore, the snowball sampling was applied in only one direction: collecting papers that cite the base approach. In April 2015, a total of 253 studies were selected. Because several papers were about either applications that used face detection algorithms or body tracking, only 22 works about full face tracking techniques remained. Next, those filtered papers were evaluated according to different metrics regarding aspects that would have influence on the performance and precision, as listed below:

- Performance aspects:

- FPS;
- RAM memory;
- ROM memory;
- Energy consumption.
- Precision aspects:
 - Datasets used;
 - Tracking precision;
 - Number of landmarks;
 - Faces in complex environments (In-the-wild).

Table 4 summarizes main aspects extracted from the most relevant papers evaluated. Finally, these features were analyzed and the Local Binary Features (LBF) method described in (REN et al., 2014) was selected. It was the only study that presented an approach that also has a mobile device implementation with real-time results.

Table 4 – Evaluation of main aspects of the most promising works. Cells with asterisk mean that there is not a clear value for the feature.

Study	FPS	Precision	Number of Landmarks	In the wild
Base Work (RAMANAN, 2012)	25 (desktop)	90%	68	Yes
LBF (REN et al., 2014)	300 (mobile) 3000 (desktop)	*	29 - 164	Yes
Real-Time Face Detection in CUDA (CHENG et al., 2014)	45 (desktop)	*	68	Yes
One ms Face Alignment (KAZEMI; SULLIVAN, 2014)	*	0.04 px	194	Yes

3.3.2 Local Binary Features Technique

The selected work presents a regression approach for face alignment that uses a locality principle to independently learn a set of highly discriminative local binary features for each facial landmark. The obtained local binary features are used to jointly learn a linear regression for the final output.

The selected approach predicts facial shape S in a cascaded manner. Beginning with an initial shape S_0 , S is progressively refined by estimating a shape increment ΔS stage-by-stage:

$$S_t = S_{t-1} + \Delta S_t. \quad (3.1)$$

In a generic form, a shape increment ΔS_t at stage t is regressed as

$$\Delta S_t = W_t \phi_t(I, S_{t-1}) \quad (3.2)$$

where I is the input image, S_{t-1} is the shape from the previous stage, ϕ_t is a feature mapping function, and W_t is a linear regression matrix.

3.3.2.1 Learning The Feature Mapping Function ϕ_t

The feature mapping function is composed of a set of local feature mapping functions that are learned individually and then combined to compose $\phi_t = [\phi_{t_1}, \phi_{t_2}, \dots, \phi_{t_L}]$. A standard regression random forest is used to learn each local mapping function for each ϕ_{t_1} . The split nodes in the trees are trained using the pixel-difference feature and the one that gives rise to maximum variance reduction is selected, as shown in Figure 31.

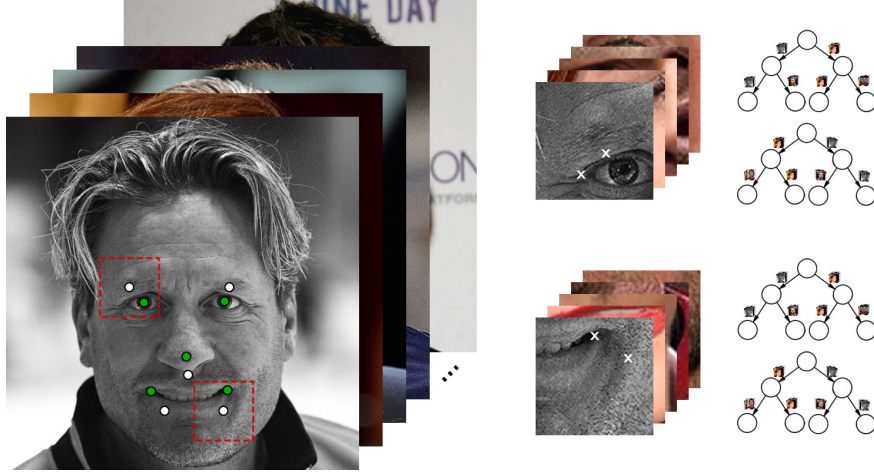


Figure 31 – For every image on the training dataset, each local feature is learned individually from the landmarks (white circles). The intensity difference between two random pixels (white crosses) is used as decision function to split the training images. Each node has a distinct pair of features.

Only pixel features in a local region around the landmark are sampled. In the training phase, the optimal region size is estimated in each stage. The optimal radius region should depend on the distribution of the ground truth landmark of each training image around the respective landmark of the initial shape. Therefore, it decreases as the landmarks converge to the ground truth, as seen in Figure 32. As a consequence, it is necessary to extract fewer features on later stages than on earlier ones for maintaining the ideal distribution.

After creating the random forest, all training images traverse the trees until they reach one leaf node for each tree. The output of the random forest is the combination of the outputs stored in these leaf nodes. Supposing the total number of leaf nodes is D , the ϕ_t will be a $D \times L$ matrix in which a 1 value means the image reaches the corresponding leaf node and 0 otherwise, as illustrated in Figure 33. Therefore, ϕ_t is a highly sparse matrix called global feature mapping function and all ϕ_{t_L} are the local binary features.

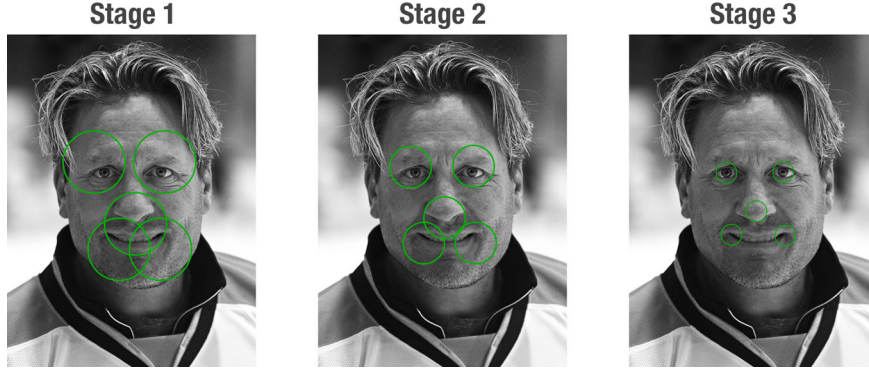


Figure 32 – Local region around the landmark on every stage.

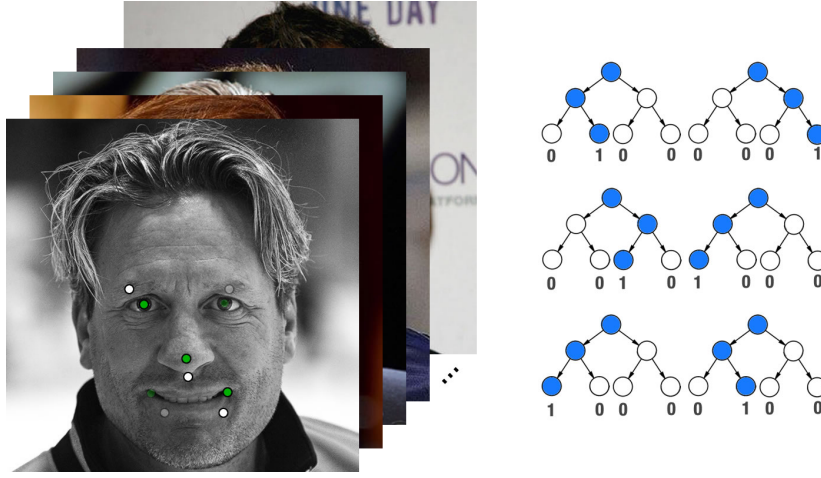


Figure 33 – Traverse of three different landmarks of one of the training images on the generated forest. The sequence of zeros and ones of every landmark is the local binary feature.

3.3.2.2 Learning The Global Linear Regression Matrix W_t

After learning the global feature mapping function, the global linear projection W_t is learned by minimizing the distance of every landmark on the initial shape to the correspondent landmark on the ground truth.

This technique that uses a combination of local learners to discover global functions is called “transfer learning”, which the authors claim that significantly improves performance for two reasons: the locally learned output by random forest is noisy because the number of training samples in a leaf node may be insufficient and the global regression can effectively enforce a global shape constraint and reduce local errors caused by occlusion and ambiguous local appearance.

3.3.2.3 Tracking a New Face

The training phase provides the global feature mapping function and the global linear projection, which are loaded by the tracker. First, a regular face detector is used to determine the position of the face in the image. This position is used to choose a location

to S_0 that is at least near the correct position of S . After that, the image traverses the trees until it reaches one leaf node for each tree, which will have 1 while all other leaf nodes will be 0. The global linear regressor is applied to the image's local binary feature to determine the ΔS_0 increment and compute S_1 , which is used as input to repeat the process until S is found. Figure 34 illustrates this process.



Figure 34 – Cascade shape regressor, in which the landmark position is incremented every stage.

3.3.3 Implementation on Android

There is an open source C++ implementation of the technique available, which was not developed by the original paper's authors². Figure 35 shows the results of this implementation.

One advantage of having a C++ implementation is that it is easier to port to the efficient Android architecture described in the last section than other programming languages since it is not necessary to rewrite all the code. Additionally, OpenCV library is the only dependence of this implementation. Another benefit is the fact that it is possible to port only the tracking part of the system to Android and maintain the training phase in desktop without having incompatibilities regarding the training file.

This desktop platform implementation of LBF has some parts that do not compile or execute on the Android platform, such as the OpenCV functions for rendering images on a window. Because of that, it was created identifiers for both platforms, Windows and Android, in order to make the same code of the LBF system work properly on each of them. Hence, a Visual Studio project was created inside the JNI folder, in order to have an application for desktop environment that runs the same code of the Android environment. As a benefit, it is possible to execute, test and debug the C++ code using the Visual Studio IDE, once the Android developer environment does not provide tools for debugging the C++ code at this moment.

Additionally, it was necessary to compile some of the libraries used in the project itself. For example, in the case of the OpenCV library, only the modules utilized by the

² <<https://github.com/yulequan/face-alignment-in-3000fps>>

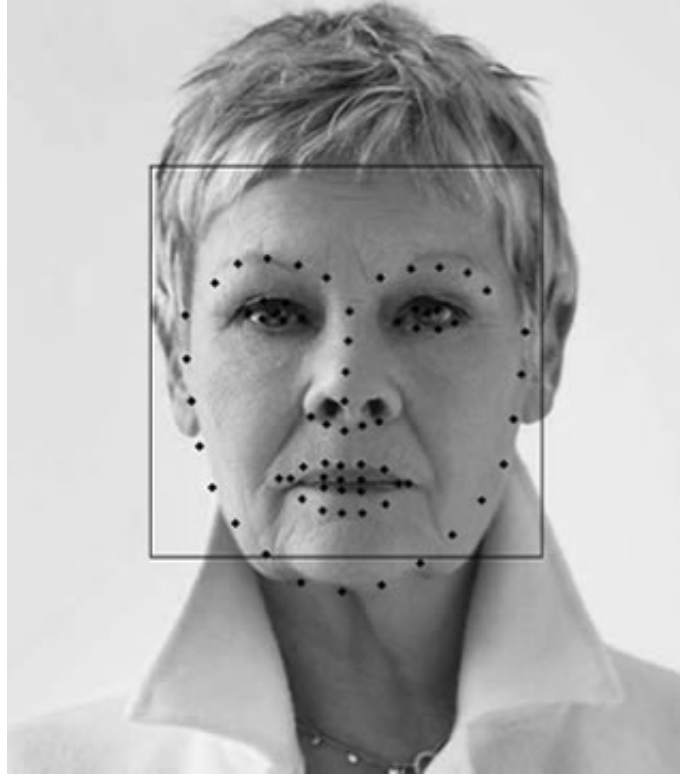


Figure 35 – Result from C++ implementation of the LBF technique.

application were compiled, discarding unused modules, and this brought gains in the size of the system.

It was possible to test several parameters in desktop to determine which ones were the most suitable for Android. Four parameters were modified:

- Number of landmarks: $L = 31$ or $L = 68$;
- Number of trees per landmarks: $N \in \mathbb{Z} \mid 4 \leq N \leq 10$;
- Depth of each tree: $D \in \mathbb{Z} \mid 4 \leq D \leq 7$;
- Number of stages: $T \in \mathbb{Z} \mid 4 \leq T \leq 7$;
- Features per stage: $(F_1, F_T) \in \{(80, 200), (160, 400), (240, 600), (320, 800), (400, 1000)\}$.

The number of features on intermediate stages varies evenly from the minimal value to the maximum depending on the number of stages. Regarding the number of landmarks, 68 is the most common configuration used by most face tracking datasets. This is the only parameter that has influence in both training file size and execution time. Therefore, it was selected the smallest subset of landmarks that is necessary to identify the parts of a face. Figure 36 shows the standard landmark configuration with 68 points and the 31 landmarks selected. The number of stages, trees and their depth impact only the size

of the training file, which grows as these parameters increase. Therefore, the code was modified to train and test the 1,120 combinations of these parameters automatically and evaluate them according to the average error of the estimated position of every landmark with respect to the ground truth, the size of the training file and execution time.

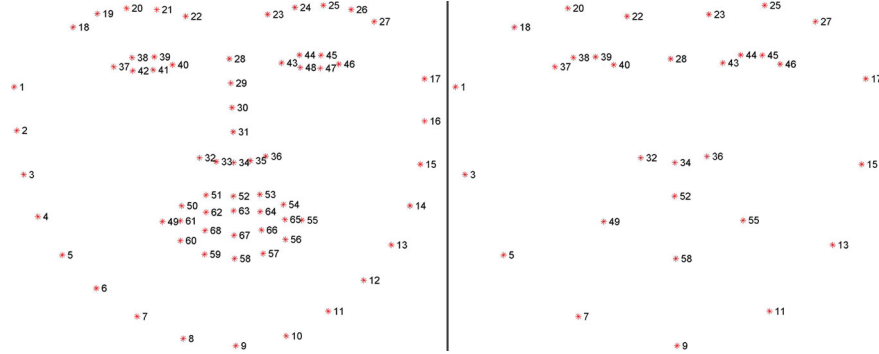


Figure 36 – Standard landmark configuration with 68 points (left) and the 31 landmarks selected (right).

The set of parameters with the best trade-off between precision, execution time and training file size was $P = \{L = 31, N = 5, D = 4, D = 4 \text{ and } F = (400, 1000)\}$. On desktop, this set of parameters presented a precision of 7.095 pixels, a training file of 1.60 MB and a mean execution time of 2.16 ms (462.96 FPS).

The training file generated with this set of parameters was used as input for the Android version of the LBF tracker. The evaluation was performed using an LG G3 device. The face was tracked using the frontal camera with 1920x1080 resolution. The binary application package file generated had 9.86 MB and used 45.4 MB of the RAM memory while executing. The execution time was 8.06 ± 3.57 ms (124.07 FPS). Regarding precision, since the Android version uses the same code as the desktop version, it is natural to say that they share the same error. Indeed, visually both implementations present similar results.

3.3.3.1 Improvements

The port of LBF to the Android platform brought some challenges regarding efficient implementation on mobile devices. One example is the data structure that should be used. The desktop implementation used double precision for floating point numbers. It has an impact on training file size. Changing these numbers from double to float decreased the size of the training file and did not have any implications in the tracking precision.

Another challenge is to know the resources provided by the Android platform and the libraries used. For instance, the first Android version was taking about 125 ms on average to track a single frame (8 FPS). The bottleneck was the OpenCV face detector, which takes more than 100 milliseconds to find the face position on each frame. The alternative was to use the Android native face detector, which provides the list of every face in a

frame by the time it is available. This feature is accelerated in hardware and has no impact on the frame capture rate, which remains at 30 FPS with or without the face detector.

These are examples of situations and problems that can occur as the tracking system becomes more complex. The solution remains as lessons learned for future implementations.

Regarding improvements on the LBF itself, the original technique was designed for desktop and supports only images in which the face is upright. On the other hand, a mobile device can be held in any orientation while tracking the users' face. Thus, the first mobile version fails to track faces if the device is rotated. In order to solve this limitation, the first approach was to use the devices' sensor for obtaining its orientation. Then, the image is rotated to make the face upright. Finally, tracking is performed using the corrected image. However, the rotation of a full HD image is very slow. Rather than rotating the entire image, a more efficient approach is to rotate only the 31 landmarks of the initial guess and put them in the same orientation of the device, as shown in Figure 37.

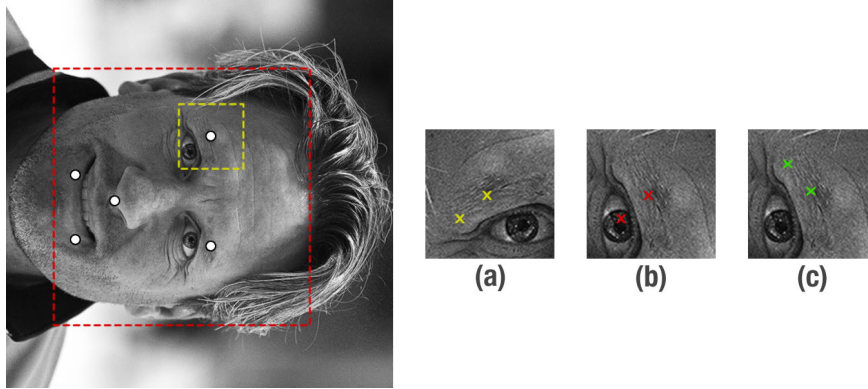


Figure 37 – Initial guess chosen as if the image is in the upright position (a). If it is used on the rotated image, it will lead to an incorrect result (b). Therefore, if the initial guess is rotated using the device's orientation, it will be on the correct position (c).

However, the patch around each landmark that is used to transverse the tree is also rotated. Thus, the two pixels used in the decision node are also rotated in order to access the correct position in the image. This modification allows LBF to track in every orientation of the device and achieve the same result as in the upright position, as can be seen in Figure 38.

3.3.4 Discussion

The use of machine learning to perform tracking seems to be very promising. The main concern is regarding the space necessary to store the training file because it can grow as the objects or scene to be tracked get more complex. For instance, in order to track only 31 landmarks in a face, LBF needs to load a 1.60 MB training file and to track a planar static model similar to the one used on the last experiment, the training file would have more than 25 MB (SIMões, 2016). This size will be even bigger when dealing with more



Figure 38 – Tracking results with different device orientations³.

complex environments, such as a small table with 3D objects and this is critical on mobile devices.

This fact is reflected on the systematic mapping presented in the previous chapter in which the machine learning approaches are used to track only parts of the body, such as face and hands. However, it is possible to make traditional methods more robust when combining them with learning techniques. One example is (VINEET et al., 2015), which uses a random forest classifier in order to evaluate the features extracted from an image to perform a semantic evaluation.

3.4 SIMPLE SLAM SYSTEM ON MOBILE DEVICES

Simultaneous Localization And Mapping, popular known as SLAM, is often related as an algorithm, but in fact it is the problem of building a map while localizing the device within that map at the same time (LEONARD; DURRANT-WHYTE, 1991). It is often referred as a chicken and egg problem because a good map is required to localize the device while a precise location is needed to build the map. Moreover, these two problems cannot be solved independently of each other and consume a lot of computational resources.

However, SLAM algorithms have been recently deployed on mobile devices since they are continuously improving regarding processing power and memory. This fact makes them powerful enough to perform such complex tasks. This scenario favors the creation of numerous types of applications since these kind of devices create several opportunities that are only possible when the user can be mobile.

This section describes the initial stages regarding the port of an initial SLAM technique to a mobile device platform using Google's Project Tango tablet (Google Inc., 2014). The main motivation to use this device is to benefit from the technology embedded on it, especially the depth sensor. The results found on this implementation were also published in (ARAUJO et al., 2016).

³ A video with this result is also available at <<https://goo.gl/kGRj4K>>

3.4.1 SLAM Systems Works on Mobile Devices

SLAM systems are traditionally used for robot navigation (SIM; ROY, 2005). Recently, their use has become more and more popular in augmented reality applications (KLEIN; MURRAY, 2007). Since then, several works were developed using mobile devices so they can benefit from the mobility provided by such tools. One example is (KLEIN; MURRAY, 2009), which describes an initial prototype of a keyframe-based SLAM system running on an iPhone 3G. They have adapted PTAM (KLEIN; MURRAY, 2007) to generate small maps in order to mitigate the impact of the device's limitations regarding processing power, memory and storage capabilities.

Recent studies have proposed different types of solutions to overcome those limitations. For instance, (PIRCHHEIM; REITMAYR, 2011) uses the plane-induced homography between keyframes to implement an efficient mapping algorithm that decreases processing time. Li and Mourikis (LI; MOURIKIS, 2012) focus on visual-inertial odometry to create a real-time navigation system. Martin et al. (MARTIN et al., 2014) optimize SLAM on a mobile device by decoupling localization and mapping steps. Each one runs in a different thread and their results are combined at the end of the process.

On the other hand, mobile devices have some features that can help SLAM systems and are not available on a desktop platform, such as the device's sensors. Usually, they combine data measured by sensors with information extracted from the camera (KAO; HUY, 2013; SCHÖPS et al., 2017). There are also works that use geolocation from the GPS combined with a remote server to reduce the amount of data processed on the device (VENTURA et al., 2014). They create a 6DoF map locally on the mobile device. The global localization method, which runs on a server, processes the globally-registered map and returns a refined global registration correction to the mobile client.

Another way to efficiently run SLAM systems on mobile devices is using a panorama map (PIRCHHEIM; SCHMALSTIEG; REITMAYR, 2013). They are registered on a 3D map in order to be able to maintain tracking during rotational camera motions, which is also a limitation for SLAM systems in general. Therefore, they can handle cameras with pure rotational motion while creating larger and denser maps. There are also works that integrate camera information with device's sensors to create the panorama map and continuously update it (VENTURA; HOLLERER, 2012). Because of that, the system can be used in large outdoor spaces.

3.4.2 Simple Tracking and Mapping

The proposed algorithm, called Simple Tracking and Mapping (STAM), is a monocular SLAM technique. In general, such kind of SLAM uses a standard camera that does not provide odometry, which gives the device position. Thus, it has to be performed visually. For STAM, visual odometry is done by tracking features using optical flow and feature

descriptors. This information is used to triangulate new points and increment the map.

In the initialization phase, it is extracted the 2D coordinates of the corners of a chessboard pattern, and they are associated with their corresponding 3D points, which are known based on the size of each square on the pattern. Such keypoints are named *square features*. Additionally, it is extracted SURF (BAY; TUYTELAARS; GOOL, 2006) features from the entire initial frame, except the area inside the chessboard. Since this pattern is repetitive, this action minimizes the extraction of non-discriminative keypoints. These extracted features have no corresponding 3D points and are named *triangle features*. Then, it is computed the SURF descriptors for both square and triangle keypoints in the first frame. After that, the projection matrix of the first camera is calculated using only the square points since they are the most reliable points available. Finally, the first frame is stored as a keyframe and the square features utilized for the pose computation are kept in a track set, which consists of the collection of features in the map currently being used for tracking.

In the tracking phase, the points in the track set are tracked in the current frame using Kanade-Lucas optical flow (LUCAS; KANADE, 1981). Successfully tracked square features from the track set are then used to compute the projection matrix for the current frame based on the obtained 2D-3D correspondences. The track set is updated leaving in this set only the features successfully tracked so far. Then, it is performed a baseline check between the current frame and the last keyframe by verifying if the distance between camera optical centers is higher than a threshold.

Whenever there is enough baseline, the mapping procedure is executed, in which SURF descriptors from keypoints extracted in the current frame are matched against the triangle features of the previous keyframe. The matching features are then triangulated, which means they now have a 3D point associated with them and are labeled as *square features*. They are also added to the map and the current track set. A feature matching between a square feature and a feature in the current frame is considered a re-detection. Square features that are re-detected are also added to the track set.

This process, illustrated in Figure 39, is repeated while there are new frames to be processed. A sparse bundle adjustment procedure that optimizes the camera trajectory is also performed after a pre-established number of new keyframes are added to the keyframe list.

3.4.3 STAM Evaluation

As mentioned before, evaluating tracking techniques is a challenging task and many efforts have been made to provide metrics to analyze such systems. Since 2008, the International Symposium on Mixed and Augmented Reality (ISMAR) promotes the ISMAR Tracking Competition, a contest aiming to challenge state-of-the-art trackers through real-world problems. All the scenarios prepared for the competition try to replicate

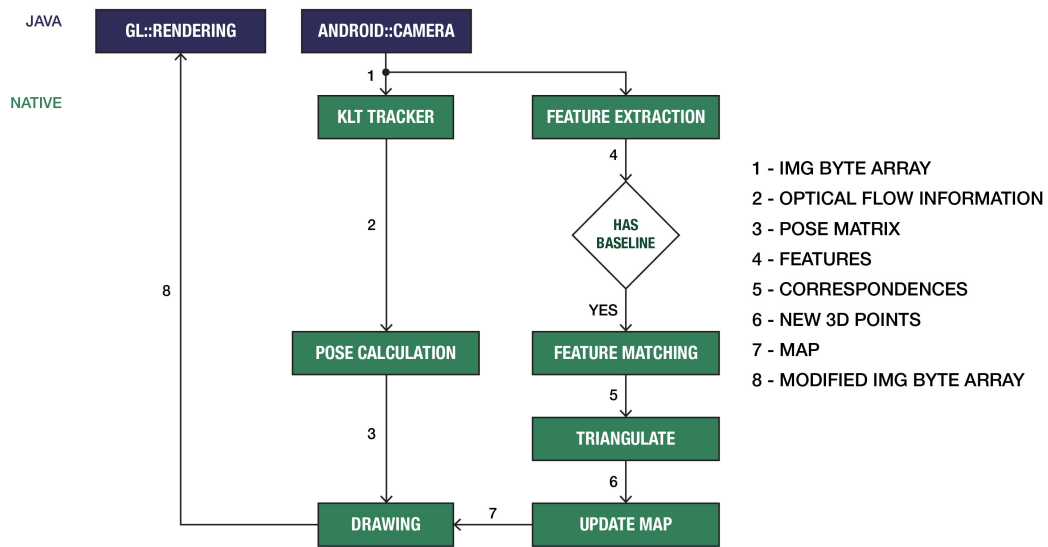


Figure 39 – STAM flow diagram.

real problems for tracking systems, such as lighting conditions, task specificities, user constraints, levels of texture information available, objects relative size, camera resolution and others.

In 2015, ISMAR introduced a different competition style in which a SLAM system should be incrementally built in order to be able to compete on different levels. There were two categories: off-site and on-site competition (International Symposium on Mixed and Augmented Reality, 2015). The former evaluates the system using images and participants would submit their results online. In the later, the system should be utilized in a real scenario that simulates a small office and participants should identify certain elements based on given 3D coordinates in an unknown area.

3.4.3.1 Off-Site Competition

This category was divided into three levels, which evaluated the processes of camera pose tracking step by step:

- **Level 1 - Point Matching:** finding the locations of known feature points in an image;
- **Level 2 - Point Tracking:** tracking known feature points in successive frames;
- **Level 3 - Mapping:** tracking and mapping new feature points in successive frames.

The first goal of *Level 1* was to find the 2D positions of the reference points, which were given by 2D image patches extracted from the same image it was used to find the 2D point, as seen in Figure 40. Every patch had a 3D point associated with it. Thus, the second objective was to calculate the projective transformation matrix from the correspondences between the 2D positions of the reference points in the image and their 3D coordinates.

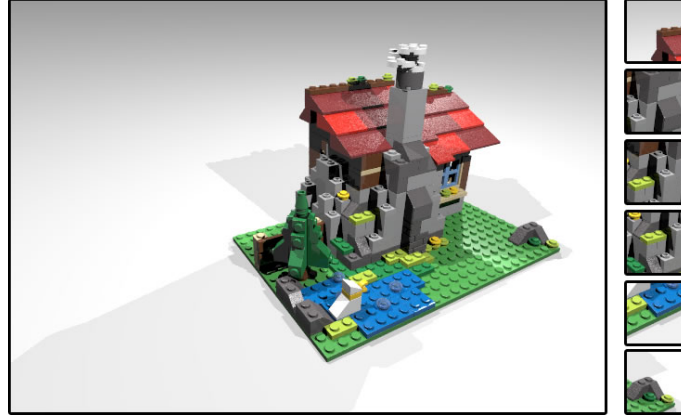


Figure 40 – Frame from which the patches were extracted. Six samples of these patches are on the right. The 2D position of each patch should be found in the same image.

On *Level 2*, the goal was to find the 2D positions of the reference points that were also given by image patches in the initial image, then track them through the sequence of frames. Finally, calculate the projective transformation matrix from correspondences between the 2D positions of the reference points in the image sequence and their 3D coordinates. As seen in Figure 41, all the images of the sequence contain the same scene part that can be found in the initial frame.

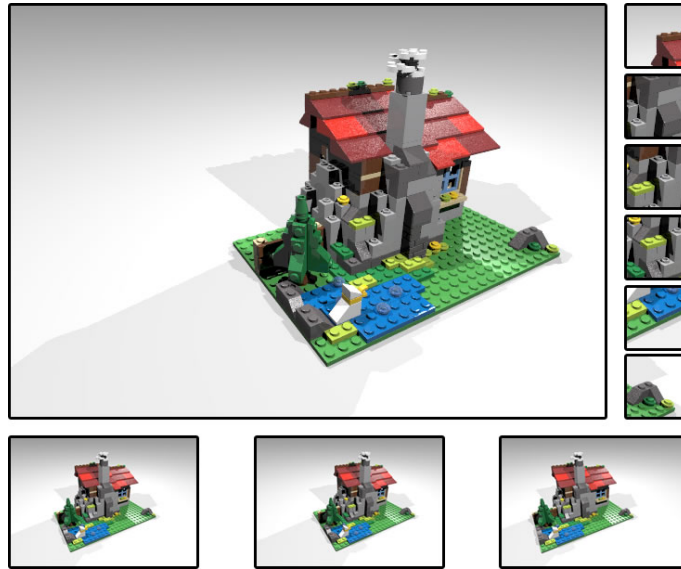


Figure 41 – Six samples of patches (top right) were extracted from the first frame (top left). These patches should be tracked along the image sequence (bottom row).

The goal on *Level 3* was to find the projective transformation matrix from the correspondences between the 2D positions only in the last image of the sequence and their 3D coordinates. This task was similar to the previous level, except that none of the patches found in the first frame appear in the last image, as shown in Figure 42. Therefore, for this level, only tracking the given points was not enough and it was necessary to map new

features.

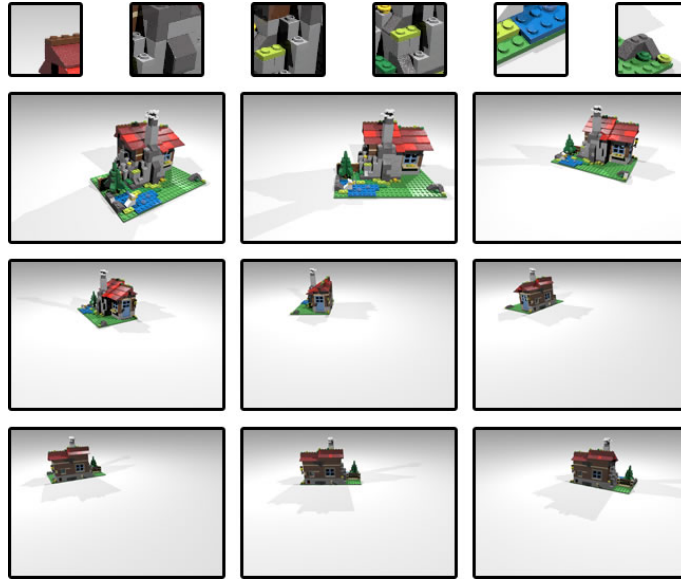


Figure 42 – Six samples of patches (top row) were extracted from the first frame and nine sample images from the sequence to be tracked. Note that the patches are disappearing along the image sequence.

3.4.3.2 On-Site Competition

In this category, the contestant should use the system to guide him in an environment that simulates a regular office in order to mark the contest points in one of the posters on the walls. The task is similar to *Level 3* of off-site category, and it was performed in two sessions in which the contest points changed between them. The difference is that no initial points were given and the tracking area was much larger, an 8 by 8 meters room, as seen in Figure 43.

There was a chessboard marker printed on A0 paper on the starting area in which the size of each square was 10 cm, as shown in Figure 44. It was used to calibrate the coordinate system. Each rectangle on Figure 43 represents a table that had different types of objects that should be used by the tracking system. Each table had objects that presented different challenges for tracking systems, such as small and reflective objects or notebooks showing a video.

As noted in Figure 44, there were some posters with black and white squares on the wall where the contest points were located. Therefore, the system should display the virtual point correct position so the contestant could mark it using a pen. Additionally, the user should perform a specific path. From the start area, the contestant should enter the first hallway, assign the first and second points in session one and the first one on session two. Then, the contestant should enter the second hallway to mark the third or second point depending if it was the first or the second session. Later, the contestant should go back

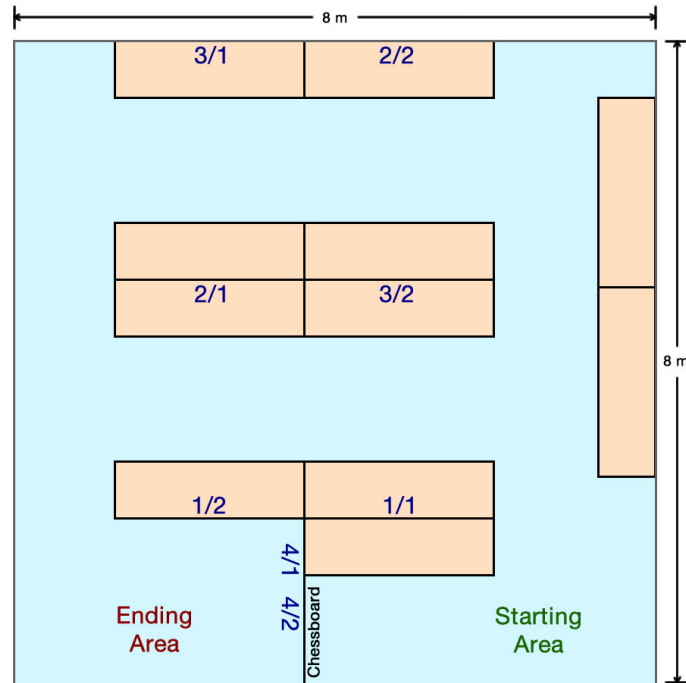


Figure 43 – Schematic of the on-site tracking area. Each rectangle represents a table with different objects. The numbers X/Y indicate the poster used to assign the four competition points in which X represents the point order and Y the competition session.



Figure 44 – Images from the tracking area. The chessboard used to calibrate the tracking system is seen on the top left image. The other images show the tables with the trackable objects and the posters to mark the 3D points.

to the first hallway to collect the third point for session 2. Finally, the contestant should enter the ending area to assign the final points.

3.4.3.3 Results

The first version of STAM was developed in C++ using data structures and basic functions from the OpenCV library. The competition rules stated that the tracker system

must run using any device that only one person can carry and still be able to mark the points on the posters. Thus, the device used was a Microsoft Surface Pro 3 with an Intel Core i5 processor having 2 1.9GHz cores, 4GB of RAM memory and an Intel 4400 graphics card. The rear camera provided the video for tracking with 640x480 pixels resolution.

The STAM system was used to compete on Level 3 of the off-site category along with five other teams. The system was also evaluated in the on-site competition, in which there were two other teams. Regarding the off-site category, the organizers released the results anonymously and identified only the winner team. Since STAM was the winning system in that category, it is possible to attest its performance on Level 3.

Table 5 shows the result of Level 3, which is the average distance of all 3D points that were mapped on the last frame when projected using the matrix computed by the system to the projection of the same 3D points using the ground truth matrix. It is possible to see that only three teams were able to track the complete scene on every scenario. On average, STAM was more than four times more accurate than the second place, which indicates that the algorithm used to optimize camera trajectory was able to minimize the reprojection error during the scenes.

Table 5 – Mean reprojection error in millimeter of all points on the last frame of off-site category Level 3⁴.

Team ID	Scenario 1	Scenario 2	Scenario 3	Average
301	552.833	2245.485	-	-
303	253.760	1250.542	402.034	635.445
304	217.836	-	-	-
305	1673.195	-	-	-
STAM	104.209	44.990	49.320	66.173
309	195.672	448.788	226.012	290.157

The results for the on-site category are shown in Table 6. Note in Figure 44 that there were several repeated objects on the table. Some of them were small, others were reflective and a few would change during the time. Additionally, the spaces between tables were almost textureless. All these facts make very hard for a system to map correctly such environment. Therefore, it is possible to see that only one team and in one session was able to map and track the whole area. As for STAM, it was able to identify only the closest point from the starting place. That is due to the fact that the system accumulates too much error during tracking in such complex environment.

3.4.4 Implementation on Tango Device

As mentioned in Section 3.1, Tango is a platform that combines computer vision techniques with state-of-the-art sensors, allowing to perform tracking on mobile devices.

Table 6 – Mean reprojection error in millimeters of on-site category.

Team ID	Session 1				Session 2			
	1/1	2/1	3/1	4/1	1/2	2/2	3/2	4/2
A	55	328	239	-	197	291	-	-
STAM	42	-	-	-	231	-	-	-
C	54	291	53	-	227	271	245	262

For this implementation, it was also selected the Yellowstone tablet.

3.4.4.1 Android Programming

Since STAM was developed in C++, it can be ported to the efficient Android architecture in a process similar to the one used to port LBF to mobile devices. That includes the identifiers for both platforms, Windows and Android, so the same code of the STAM system works properly on each of them.

Additionally, it was necessary to compile some of the libraries used in the project itself. For example, the OpenCV library was recompiled to use only the modules necessary to the application and also to add extra modules required by the STAM system that are not available on the standard OpenCV version, such as xfeatures2d. Another example was the library for performing bundle adjustment. In desktop, it was cvsba (LOURAKIS; ARGYROS, 2009), which does not have an Android version. Thus, it was replaced by the Ceres Solver library (AGARWAL; MIERLE et al.,), which is an open source C++ library that supports Android and was also recompiled for the mobile platform.

3.4.4.2 Results

It was performed a preliminary evaluation of the Tango version of STAM using the scenario shown in Figure 45. It aims to compare performance and precision between the desktop and mobile implementations. Although the Android version shares the same code of the desktop development, the optimization with bundle adjustment uses a different library. Therefore, it is necessary to evaluate if this modification has any impact on the precision.

As for execution time, it was measured the time to perform the three main steps of the STAM technique, which are calibration, frame by frame tracking using optical flow and optimization with bundle adjustment. Figure 46 compares the results between the desktop⁵ and the Tango implementation.

The performance difference between platforms on the calibration step is the smallest one. It is because this procedure only executes OpenCV functions and stores only few data.

⁴ A video with this result is also available at <<https://goo.gl/3mdkcL>>

⁵ Intel Core i7 with 3.60GHz and 8GB of RAM



Figure 45 – Evaluation environment with the calibration chessboard in the middle and objects with different types of texture around it.

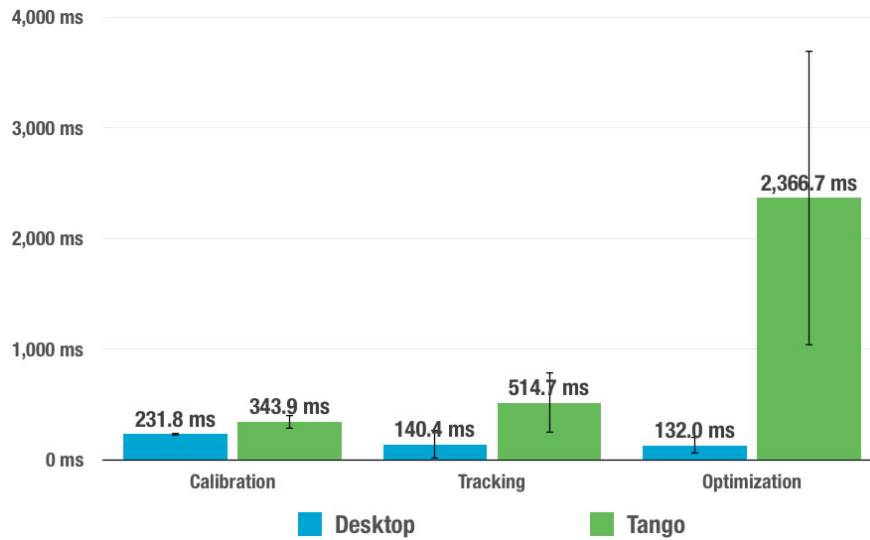


Figure 46 – Execution time in milliseconds to run the main steps of STAM.

The performance difference increases for the tracking step due to two main reasons. First, as the number of points on the map increases, the Android version felt more the lack of processing and memory resources than the desktop version. Then, as tracking processing occurs, the Tango device starts to heat up fast, making the execution even slower.

The biggest difference between platforms is on the optimization process. It is because STAM uses different bundle adjustment libraries on both environments. Moreover, the *cvsba* library only optimizes the camera trajectory while the *Ceres Solver* library optimizes both camera path and point cloud. Therefore, the optimization step on Android deals with much more variables than the desktop implementation.

Regarding tracking precision of STAM system on the mobile device, the reprojection error on both implementations was compared. It was used a chessboard printed on a graph paper to calibrate the coordinate system. Then, the system starts to track a 3D point

located on the outer corner of one of the squares, shown in Figure 47. After one minute of tracking in which the device was moving freely, the user marked the position of the point at the paper. The same procedure was performed with the desktop version. The average distance of every assigned position to its correct location is the tracking precision. The error of STAM running on Tango was 10.50 millimeters (± 3.91) while the average error on desktop was 19.75 millimeters (± 6.98). This was expected because the Android version has a more robust optimization.

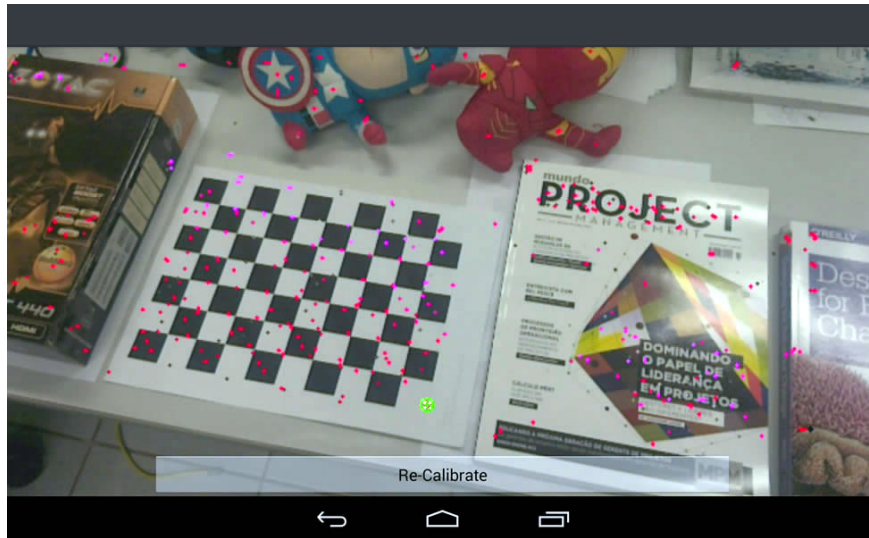


Figure 47 – Tracking procedure running on Tango device. Small points are the keypoints extracted from the mapping while the large green dot shows the tracking point, which is a known point in the real world based on the chessboard template⁶.

⁶ A video with this result is also available at <<https://goo.gl/ME9JHJ>>

4 INCREMENTAL SEMANTIC TRACKING

From the preliminary experiments, it was acknowledged the importance of finding high-level semantic information from a point cloud. This is a challenging task and it can be used in various applications. Also known as called semantic modeling, it is useful to compactly represent the scene structure and efficiently understand the scene context. And such knowledge results in several improvements when applied to tracking techniques. It was mentioned in this thesis that two of the most used augmented reality libraries incorporate simple semantics into their tracking techniques. Google’s ARCore and Apple’s ARKit estimate planes to place the augmented content, which can be the floor, a table or vertical planes such as walls. With the detection of planes, these trackers allow a more realistic rendering and stable positioning of the virtual objects.

The benefits of detecting other kinds of shapes are even higher. For instance, it can be used to provide haptic feedback on augmented reality applications (HETTIARACHCHI; WIGDOR, 2016). Besides that, objects are usually overrepresented when defined using a point cloud because it is not necessary to have so many points to describe it. Therefore, an implicit representation can replace redundant points, which is particularly helpful when targeting devices with memory restrictions, such as robots or UAVs. Furthermore, a tracking system can use these primitives to denoise the reconstructed map or constrain its optimization (RAMADASAN; CHATEAU; CHEVALDONNÉ, 2015), which can reduce tracking errors. For instance, these semantic constraints can make a difference in the optimization phase of the experiment presented in Subsection 3.4.4, in which the execution time increases as the number of points grows.

This chapter explains a method for incrementally modeling and tracking primitives on a sparse point cloud. It uses the generating process of point clouds on SLAM effectively and relies on geometric and statistical analyses to filter unreliable shapes. This approach was presented in (ROBERTO et al., 2018) and is the extension of the technique published in (ROBERTO et al., 2017).

4.1 SEMANTIC MODELING

Automatic reconstruction of 3D object shapes is useful for several applications, such as blueprint generation for architecture. Since it is a difficult task to accomplish, it has been a relevant research topic for years. Several methods have been proposed to achieve it by using laser scanners (ZHU et al., 2011) and cameras (MA et al., 2005). Due to low-cost RGB-D sensors, namely Microsoft Kinect and mobile devices with Google Tango, 3D data acquisition became more common and runs in real time even on such devices. These sensors acquire the depth of an object on a pixel-by-pixel basis and, then, describe the

obtained shapes as a point cloud. Although they are very useful for 3D measurement and visualization, there are some aspects to be improved. For instance, the scene is usually represented using a point cloud or a mesh computed from it. The latter simply consists of connected points with little information about the semantic structure.

Several methods have been proposed in the literature to determine semantic in point clouds, most of them dense. One conventional approach is to use reverse engineering techniques to estimate geometric primitives, such as region growing (HOLZ et al., 2012). It can efficiently deal with large amounts of data because it makes simple comparisons using the normals to determine if a set of points belongs to the same group. However, this approach is not robust to noisy point clouds because it can lead to a wrong classification. The robustness has been improved using both Hough Transform (DROST; ILIC, 2015) and RANSAC (LOPEZ-ESCOGIDO; FRAGA, 2014). Another approach is based on machine learning techniques, which combine local features and AdaBoost to detect complex objects (PANG; NEUMANN, 2013). Support Vector Machine (SVM), Fast Point Feature Histogram (FPFH) descriptors and RANSAC are also used to extract semantics from a point cloud (HUANG; YOU, 2013). Recent works used Convolutional Neural Network (CNN) to perform a semantic classification using the keyframes of a dense monocular SLAM and then apply this result to the point cloud (TATENO et al., 2017).

Some methods are able to detect only specific primitives, such as planes (NGUYEN; REITMAYR; SCHMALSTIEG, 2015; OEHLER et al., 2011) or cylinders (LIU et al., 2013; QIU; ZHOU; NEUMANN, 2014). Although limited, detecting only one class of shape still has several applications. For example, (KIM et al., 2012) estimates the floor plan of houses from the planes detected in a dense 3D point cloud generated using the Kinect sensor. On the other hand, other methods deal with different classes of shapes at once. For instance, GlobFit (LI et al., 2011) can estimate planes, cylinders, cones and spheres. Some even detect pre-modeled complex shapes along with planes and cylinders in industrial scenarios (PANG et al., 2015).

One advantage of having dense data extracted from laser or infrared sensors is that the acquired point cloud contains several information that the algorithm can use for the detection. Moreover, the data is relatively noiseless, which means that a large number of points can stably fit primitive shapes. Therefore, it is more challenging to extract primitive shapes in a noisy and sparse point cloud of a partially-observed object computed from image-based approaches with mobile devices. Some studies have tackled this issue by limiting specific situations, such as detecting only one shape class, such as ARKit and ARCore. Another example is (SINHA et al., 2008), which estimated planes based on their reconstruction to create textured models. In this context, RANSAC based methods are very promising to work with sparse point clouds. It is because they estimate primitives by initially picking a minimal group of points for each shape and detecting the one that approximates the maximum number of points (SCHNABEL; WAHL; KLEIN, 2007). Besides,

they can also work with data containing a large number of outliers (ROTH; LEVINE, 1993). However, the performance of such approach for sparse point clouds was never investigated.

Another aspect of existing semantic methods is that they usually work in batch, which means that an input data is analyzed all together only once. It is consistent with the generation method. Usually, the dense sensors generate the entire point cloud at once. However, the performance of visual SLAM system regarding both the accuracy of the reconstruction and the computational cost for real-time applications improved drastically in recent years (MUR-ARTAL; MONTIEL; TARD?S, 2015; UCHIYAMA et al., 2015). One characteristic of several of these SLAM methods is that the 3D map is generated incrementally. This generating process can provide valuable information for a semantic approach in addition to the point cloud itself.

4.2 RANSAC-BASED METHOD ON SPARSE POINT CLOUD

The first step in this research to perform semantic modeling and tracking on sparse point cloud was to evaluate how dense methods perform using sparse maps. This could provide valuable information to develop a new approach specific for sparse point clouds or adapt existing ones. As mentioned in the previous section, RANSAC methods seem to be promising for that task. From all approaches found in the literature, Efficient RANSAC (SCHNABEL; WAHL; KLEIN, 2007) presents the best results. Both regarding execution time and the precision of the primitive estimated. Moreover, it detects multiple classes of shape, which are planes, cylinders, spheres, cones and torus.

4.2.1 Method Overview

Efficient RANSAC requires a point-cloud \mathcal{P} with normals for each point as input. The output is a set of primitive shapes Ψ in which a point $p_N \in \mathcal{P}$ will be assigned to only one shape $\psi_n \in \Psi$ or it will remain unassigned.

In summary, (SCHNABEL; WAHL; KLEIN, 2007) searches the primitive with maximal score m for each iteration of the technique. The score is a function based on the number of points that can be part of a shape candidate, which consider the distance of a compatible point as well as the deviation of its normal from the one of the primitive. Therefore, for each iteration, the algorithm randomly selects one point of \mathcal{P} and then collects a minimal subset of points that are closer to the first one. It is because (SCHNABEL; WAHL; KLEIN, 2007) explores the fact that shapes are local phenomena, which means that the probability that two points belong to the same primitive is higher the smaller the distance between the points.

Candidates of all considered shape types are generated for every minimal set and all candidates are collected in the set \mathcal{C} . All of the candidates need only three points to describe each type of shape. The score m for each candidate is computed using a statistical

approximation that increases the algorithm performance. The candidate is only selected if the probability $P(|m|, |\mathcal{C}|)$ that there is no better candidate is high. In this case, $|m|$ and $|\mathcal{C}|$ are the number of points in the shape and the number of candidates, respectively. When a candidate is accepted, the corresponding set of points is removed from \mathcal{P} . The algorithm repeats until the probability $P(\tau, |\mathcal{C}|)$ of not finding new shapes is high given a τ value that is defined by the user.

4.2.2 Evaluation

There is an implementation of Efficient RANSAC available¹. This version was used to test how it deals with noisy and sparse point clouds. It was created a test scene in which two consecutive keyframes can be seen in Figure 48 (a) and (d). The point cloud, shown in Figure 48 (b) and (e), were generated using a SLAM system (UCHIYAMA et al., 2015) that was modified to increase the number of generated 3D points. These clouds have 1,427 and 2,180 points respectively. The reconstruction of the clouds was calibrated using a chessboard pattern so that the clouds were represented in metric scale.

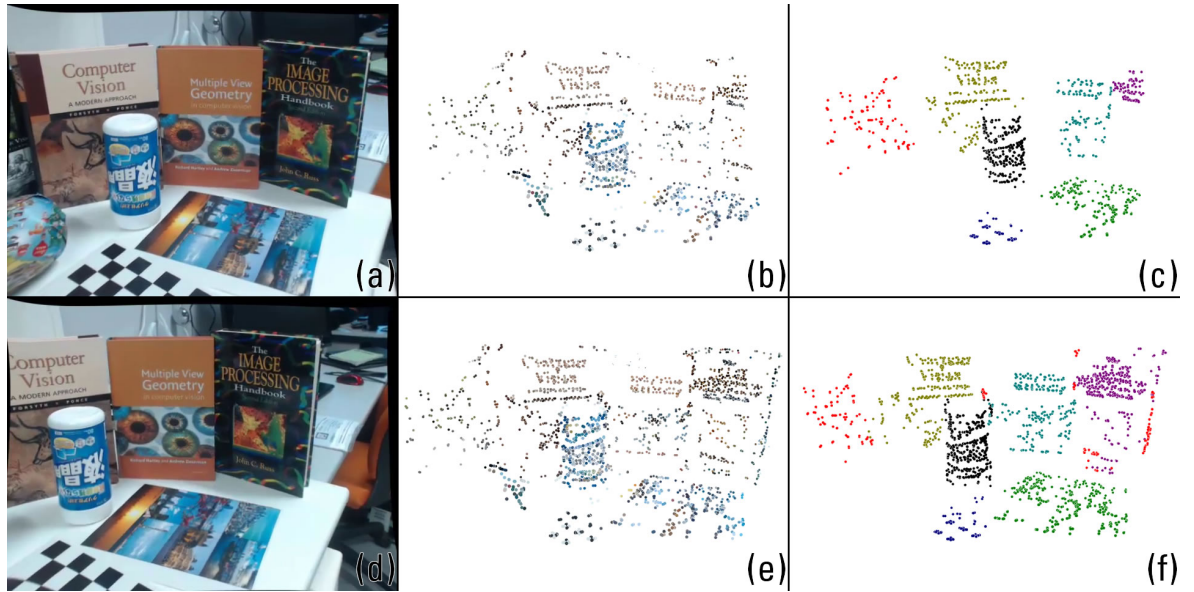


Figure 48 – Test scene (a) and (d) and their reconstructed point cloud (b) and (e). Eight shapes were detected on the first keyframe (c) and twelve primitives were found on the second one (f). Even with the point cloud being very similar, the red points represent the six primitives that were differently detected between keyframes.

Efficient RANSAC has five heuristic parameters. To deal with a sparse point cloud, more than 180 thousand different combinations of these parameters were automatically tested. The goal was to find the best set that maximized the number of assigned points, which were the ones from the input point cloud that were modified to be fitted to a shape.

¹ Available at <<http://cg.cs.uni-bonn.de/en/publications/paper-details/schnabel-2007-efficient/>>

These parameters should also minimize the distance between these points and the ones on the input point cloud. Figure 48 (c) shows the result for one of the best set of parameters, which assigned 95.937% of the points with an average distance of 1.949 millimeters for each assigned point. Efficient RANSAC was able to detect most of the shapes correctly. However, one of the detected primitives was different from the expected result since a plane was identified as a cylinder.

Another issue noted was the inconsistency between keyframes. Figure 48 (f) shows the shapes detected using the data from the subsequent keyframe. In the result, 93.896% of the points were assigned with an average distance of 2.098 millimeters. Four detected primitives were wrong, including two planes identified as cylinders. Compared to the previous keyframe, six shapes were different. This unstable result occurred often due to the noisy data combined with the small number of points.

4.3 GEOMETRIC AND STATISTICAL INCREMENTAL SEMANTIC TRACKING

From the evaluation, it was clear that Efficient RANSAC achieves good results when detecting primitives in a sparse point cloud. However, it still requires improvements of consistency and precision for various applications. Therefore, it is possible to use the primitive estimation from Efficient RANSAC and the generation process of point cloud from visual SLAM systems to perform an incremental semantic modeling. This approach can improve both the precision and stability of the primitive detection. Moreover, it also allows the tracking of these primitives through the scene.

In summary, this method runs Efficient RANSAC using the sparse point clouds that are incrementally generated from a visual SLAM system. Then, it uses the history information of the estimated primitives and their parameters to match the shapes over time. Also, it estimates the reliability of the detected primitive using the geometry of the shape. When this estimation is not reliable enough, it performs a statistical evaluation using the detection history to eliminate random detected shapes. Figure 49 illustrates the flow of the method, which is detailed in the following subsections.

4.3.1 Efficient RANSAC

When a visual SLAM system reaches a keyframe, it updates the map adding new points to it. Then, the Geometric and Statistical Incremental Semantic Tracking method, or simply GS-IST, runs Efficient RANSAC to detect the shapes for every new map. However, it was necessary to make some modifications in both the code and the configuration in the available code of Efficient RANSAC. The configuration changes aim to reduce the number of types of primitives detected. Instead of five, three classes of shapes are tracked: planes, sphere and cylinder. They were selected because they can be used to model most of the

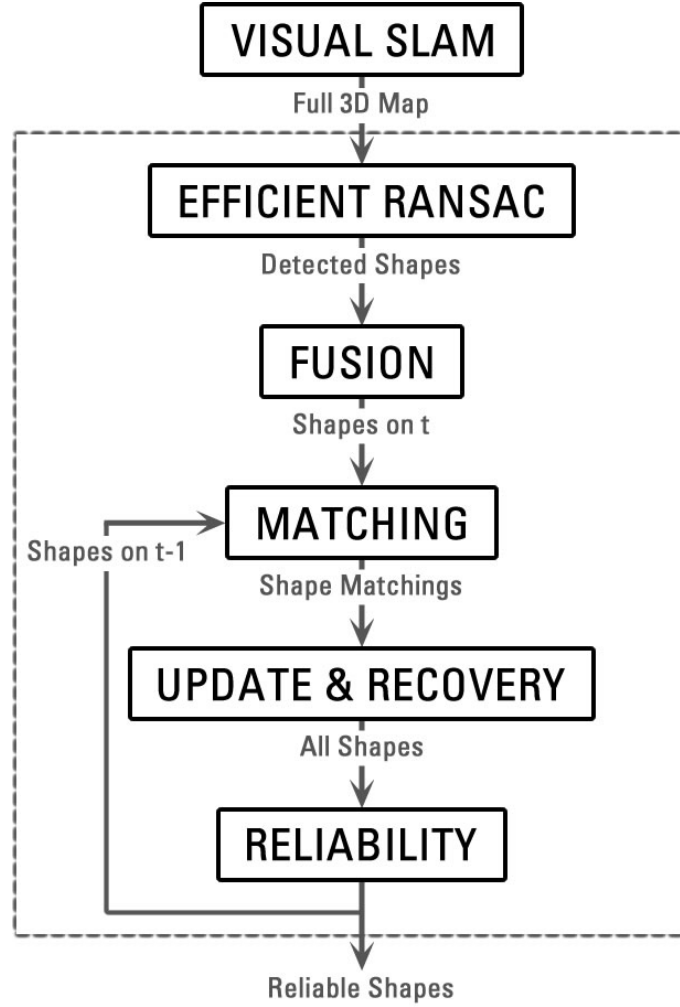


Figure 49 – Flow of the Geometric and Statistical Incremental Semantic Tracking (GS-IST) approach.

objects. For instance, when modeling industrial scenarios, almost 80% of the scene can be represented only by planes and cylinders (HUANG; YOU, 2013).

There were two code modifications. The first one intends to replace the random selection of the initial seed to use the same pseudo-random point for every execution, which helps in the evaluation process. The other change was necessary to save the value of intermediary computation in the Efficient RANSAC algorithm that will be used later by GS-IST. This data is the average Euclidean distance between the points in the input point cloud and their respective projection on the estimated shape. Since Efficient RANSAC computes this projection as part of its algorithm, returning this information would be more effective than recalculating it later. Originally, this library returns the list of primitives with their respective parameters and all input points used to estimate them projected on the shapes. After the modification, it also provides for each primitive the sum of the distances between all input points and their projections. This way, the necessary information will be available for a quick access when needed.

All the primitives are passed to the subsequent modules, which will identify and

eliminate the unreliable shapes and track the remaining ones.

4.3.2 Shape Fusion

In keypoint-based visual SLAM systems, most of the characteristics are clustered at highly textured areas. For example, in Figure 48 (a), two different patterns over the table were observed and formed two distinct clusters of points as illustrated in Figure 48 (b). In this case, Efficient RANSAC detected them as two separate planes even though they belong to the same plane on the table.

In order to improve the results from Efficient RANSAC, different shapes that belong to the same primitive are first fused. It not only provides a more representative primitive but also increases the overall information regarding the shape. Since a primitive with a small number of points can be unreliable, it is better to discard it. However, if more than one shape with similar parameters is detected, even if they are small, it is safe to assume that they correspond to the same element in the scene. The reason is that it is unlikely that two primitives are wrongly detected with the same parameters. Therefore, the fusion of primitives based on parameters helps in the history evaluation. This process uses the similarity of the shape parameters to decide whether to fuse two shapes or not:

- **Plane:** the planes are parallel given an angle threshold α_t and the distance between them is smaller than the distance threshold d_t ;
- **Sphere:** the distance between their centers is lower than d_t and the difference between their radii is less than the radius threshold r_t ;
- **Cylinder:** the angle between the axis direction of both cylinder is smaller than α_t , the distance between them is less than d_t and the difference between their radii is smaller than r_t .

While d_t and r_t are controlled for each case such that they are 2% of the largest size of the point cloud bounding box, α_t is always set to 5° .

Additionally, it is considered the proximity between the primitives to restrict or widen the similarity thresholds. The principle is that two distant shapes have a smaller possibility to be the same than closer ones. Experimentally, the thresholds are widened by 25% when evaluating the fusion of primitives that have an intersection. Otherwise, it is restricted by 25%. This means that distant shapes have to be more alike to be merged. On the other hand, closer primitives are more likely to be the same and the threshold can be less restricted.

4.3.2.1 Parameter Computation

Similar shapes according to the aforementioned criteria are fused. GS-IST sets the parameters of the resulting fused shape as a weighted average between the parameters

of both primitives. The weight is based on the geometric analysis of the points in the detected shape. The main idea is to use the average Euclidean distance of each input point that was used to estimate the primitive to the resulting shape. Even for noisy data, this distance will be smaller on correct estimations than on wrong ones. For instance, considering a globe being tracked that was correctly modeled as a sphere or incorrectly detected as a plane. The average distance of the 3D points in the globe to their projection in the sphere will be smaller when compared to the distance of the same input points to their projection in the wrong plane. Figure 50 illustrates this idea.

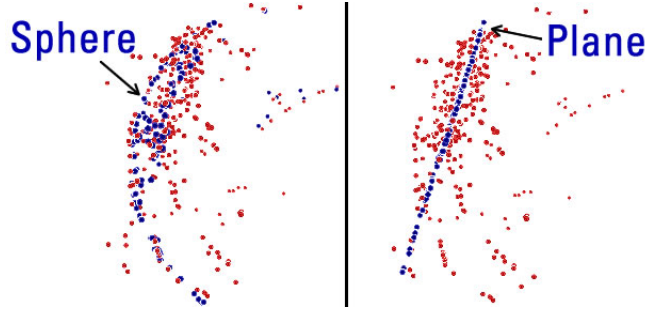


Figure 50 – Difference between the input points to their correspondent projected points on a shape estimated correctly (left) and incorrectly (right).

Thus, the parameters P_f of the fused shape will be:

$$P_f = \sum_{i=1}^n w_i P_i, \quad (4.1)$$

where n is the number of similar shapes to be fused and P_i are their parameters. The weight w_i is:

$$w_i = \frac{np}{\sum_{j=1}^{np} d_j}, \quad (4.2)$$

where d_j is the Euclidean distance between the input points to its projection in the estimated shape and np is the number of points in the estimated primitive. The weights are normalized and $\sum_{i=1}^n w_i = 1$.

Using Equation 4.1, every fused shape will influence the resulting primitive. However, it will be closer to the one with the smaller error. This average can be applied to every parameter except the plane and the cylinder position. For the plane position, it is only valid if they are all projections on the other planes. Thus, the position of one point is projected in all others and then the weighted average is computed, as illustrated in Figure 51. This will also work in case of parallel planes. As for the cylinder position, this parameter will be the axis intersection. In case of concurrent or parallel axes, the cylinders are fused in pairs if there is more than one. First, it is selected the points on each axis that is closer to the other and the resulting position will be their weight average.

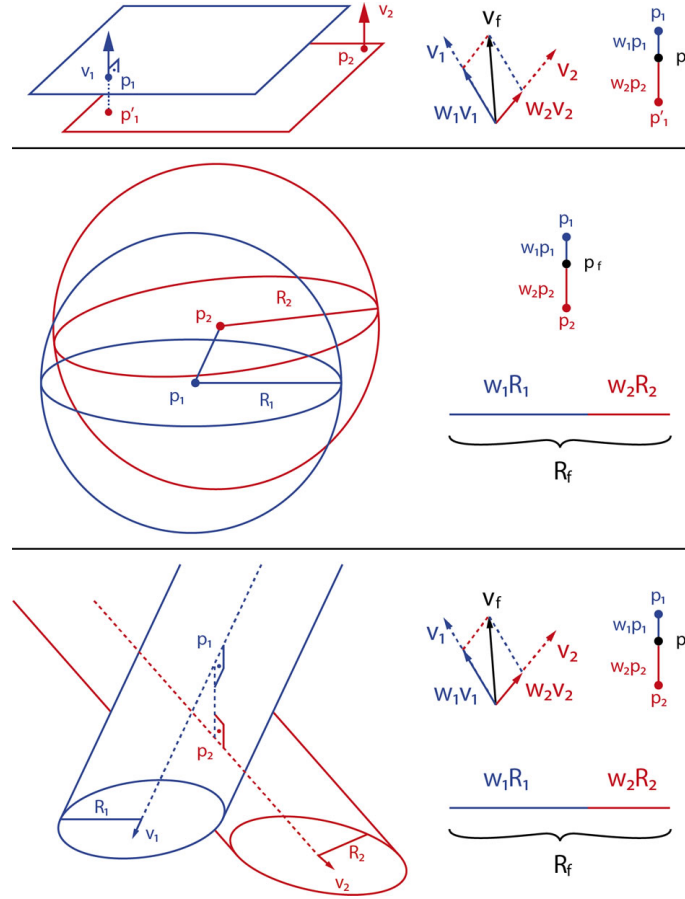


Figure 51 – Fusion of parameters for different classes of shapes.

4.3.2.2 Inclusion Criteria

In case a shape is not fused with any other, GS-IST performs an initial reliability evaluation to decide whether to keep this detected primitive or not. It is used four geometric characteristics of the primitive to make this assessment, as described below:

- **Number of Points:** shapes with more points are usually more reliable because the estimation is based on a large amount of data. The number of points in good primitives is *larger than 2%* of the whole point cloud;
- **Dispersion:** it measures how spread are the points in the primitive. Since the keypoints are clustered around highly textured areas, small regions tend to concentrate most of the shape points. The dispersion of reliable shapes is *smaller than 20%* of this value for the whole point cloud;
- **Distance:** it is the same Euclidean distance mentioned previously. Reliable shapes have an average distance *smaller than 5%* of the largest size of the entire point cloud bounding box;
- **Radius:** the sphere and cylinder radius can also provide a hint regarding the shape's reliability. A noisy plane can be estimated as one of these two primitives with a

considerable radius. Therefore, good spheres and cylinders have radius *smaller than* the largest size of the entire point cloud bounding box. It should be noted that this criterion is not applied to planes.

The system only keeps shapes that pass in all of these criteria. These values were determined experimentally and they are based on the dimension of the input point cloud because it puts all thresholds in proportional to the scene scale.

4.3.3 Shape Matching

To match the primitives between consecutive keyframes, GS-IST uses the intersection of the 3D bounding box and the distance between the center of mass. Due to several factors, such as inconsistency or the shape volume, a primitive on a given keyframe can match with several others on the previous one, including shapes of a different type. Therefore, it is necessary to detect the shape on the previous keyframe that is the most likely to be the correspondent on the current one. It is computed a s score for each primitive that a given shape on current detection intersects on the previous estimation. This score is proportional to the intersection volume and inversely proportional to the distance between the center of mass:

$$s = \frac{\sum_i^{ns} |\psi_i| p_i}{\sum_i^{ns} |\psi_i| d_i} \frac{|\psi_s|}{\sum_i^{ns} |\psi_i|}, \quad (4.3)$$

where ns are the indexes of the shapes with intersection on previous keyframes, $|\psi_i|$ is the number of points of that shape, p_i and d_i are the intersection ratio and distance between centers of mass, respectively, and $|\psi_s|$ is the number of points in the primitive on the current keyframe. It is selected the one with the maximum score as correspondence.

4.3.4 Shape Update and Recovery

The shape detected on current frame inherits the history data of the one it matched on previous detections. These data contain the primitive class that was detected on every keyframe, as well as the average distance to the original point cloud at that detection. With that information, GS-IST can verify if the current estimation is following the historical data.

The system checks the class that this primitive was detected as over time. If the current shape has the same type of the one that appears in more than half of them, including the current detection, its parameter is updated. This new parameter will be the weighted average of each detection over time. With this update, every previous detection will influence the final shape, which results in primitives that are more stable through time rather than using just the parameters from the last one.

The process to update the new parameter P_u is similar to the one explained in Subsection 4.3.2:

$$P_u = \sum_{i=1}^{n-1} w_i P_{u-1} + w_n P_f, \quad (4.4)$$

where n is the number of shapes in the past, including current detection and w_i are their respective weights, which are normalized. P_{u-1} and P_f are, respectively, the parameters in previous detections and the current parameter after fusion.

On the other hand, if the current shape has a type that is different from the one that appears most of the time, it is changed to that class of primitive. As for the parameters, it will be the same as the previous P_u from that type.

In this step, GS-IST also evaluates shapes that were not detected on the current keyframe but appeared previously. It recovers these shapes with the same parameters from the last appearance. The history data will be updated using the average distance of the recovered shape, but it will not have a primitive class associated at this particular moment. This shape will eventually disappear when not detected anymore because there will be no class of primitive that appear in more than 50% of the time.

4.3.5 Reliability Computation

At this point, GS-IST has a set of detected shapes and it is necessary to determine which of them are reliable. This reliability computation is based on a geometric and statistical evaluation.

4.3.5.1 Geometric Analysis

Each shape has a history information since its first appearance, which is a list of each primitive it matched in the past keyframes. However, a given shape may have different classes over time due to imprecision or inconsistency. Therefore, to perform the geometric analysis, it is computed the weight w_c for each class of primitive that appears in the history data:

$$w_c = \frac{1}{\sum_{i=1}^h d_i}, \quad (4.5)$$

where h is the number of times that each class of primitive appears in the history data and d_i is the average distance of the points in that shape to the original point cloud.

The weights are normalized and the one with the maximum value is the dominant class. GS-IST judges shapes whose dominant primitives have a weight higher than 0.75 as reliable. On the other hand, it considers unreliable those in which all weights are smaller than 0.5. When the weight of the dominant shape is between these two values, its classification will be determined by the statistical analysis. If this evaluation shows that the detection class

through history is random, the primitive will be unreliable. Otherwise, it will be set as reliable.

4.3.5.2 Statistical Analysis

GS-IST performs a runs test for randomness to determine if the estimation history is random (SHESKIN, 2011). Basically, this non-parametric test uses the expected value and standard deviation to estimate the minimum number of runs that a sample can have to be considered random. A run means a sequence of consecutive estimates of one particular class of primitive. However, it was decided to look at the history of classification as binary data because the convergence is faster. Therefore, it was denoted a + for the first primitive detected. Then, the sign is repeated if the shape class is the same as the previous one. Otherwise, it is inverted. For a 5% level of significance, the sample is random if the number of runs is greater than:

$$N(R) = \mu - 1.65\sigma, \quad (4.6)$$

where the expected value μ and the standard deviation σ for the total number of samples n are:

$$\mu = \frac{2n - 1}{3}, \quad (4.7)$$

$$\sigma = \sqrt{\frac{16n - 29}{90}}. \quad (4.8)$$

Table 7 shows the example of a history information with a sequence of four spheres, followed by one cylinder, one plane and then by two other spheres. There are $R = 4$ runs and $n = 8$ samples. In this case, the maximum number of runs for a nonrandom sample is $N(R) = 3.269$, which indicates a random detection.

Table 7 – History of the estimated shape from a primitive. For each sample that represents a keyframe K_i , it was classified as plane (P), sphere (S) or cylinder (C).

	K_1	K_2	K_3	K_4	K_5	K_6	K_7	K_8
Primitive	S	S	S	S	C	P	S	S
Label	+	+	+	+	−	+	−	−

4.4 EVALUATION

GS-IST was implemented in C++ using OpenCV² and Efficient RANSAC as libraries. For this evaluation, it was compared GS-IST with Efficient RANSAC regarding precision,

² Available at <<http://opencv.org/>>

recall and $F_{0.5}$ -Score. The choice of this metric instead of F_1 -Score was because it highlights the precision, which is the focus of the method. The only public dataset found is designed to detect primitives based on single-view images (XIAO; RUSSELL; TORRALBA, 2012), which was not suitable for this incremental approach. Since it was not found any dataset that has the generating process of point clouds, it was created one with five different scenarios to evaluate semantic modeling and tracking, which is available for download³. This dataset has the RGB images, a text file containing the camera parameters and the rotation and translation for each frame, the list of keyframes and the point cloud generated on each keyframe. It has five scenes targeting distinct types of primitives and different numbers of keyframes, as seen in Table 8 and illustrated with some screenshots in Figure 54. The number of points in the last keyframe indicates how sparse are the point clouds. It is worth to mention that Efficient RANSAC does not track the primitives, it only detects the shapes at each keyframe because it is a batch-based approach.

Table 8 – Details of the dataset used for the evaluation, in which P, S and C stands for Plane, Sphere and Cylinder, respectively.

Test Case	Number of Frames	Number of Keyframes	Number of Points in Last Keyframe	Primitives in Scene
<i>Case 1</i>	1,660	31	16,698	P, S, C
<i>Case 2</i>	1,346	24	3,874	C
<i>Case 3</i>	849	20	4,257	S, C
<i>Case 4</i>	405	7	2,781	P
<i>Case 5</i>	499	17	4,445	P

It is possible to see in Figure 52 that GS-IST obtained 100% precision in all cases while Efficient RANSAC never achieves more than 82%. It was noticed that there are more wrong estimations in the initial keyframes, as seen in the chart on Figure 53. It uses the precision of *Case 1* over time to illustrate this behavior, which is expected since the point cloud has few points in the initial reconstructions. The geometric and statistical analyses are able to identify these early incorrect detections. For instance, in the first three keyframes of *Case 2*, the bottle in the right side is assigned as a sphere, then as a cylinder and later as a sphere again because of the small number of noisy points. For the Efficient RANSAC, the bottle is incorrectly estimated as a sphere twice in three consecutive detections. Using this tracker, the bottle is assigned to the same primitives in the first three keyframes but, each one has a weight based on the geometric analysis. After normalization, the weights of detection history are 0.186 (sphere in the first keyframe), 0.537 (cylinder in the second keyframe) and 0.277 (sphere in the third keyframe). Thus, for GS-IST, the bottle will be assigned as a cylinder because its weight is higher than the 0.463 of the sphere.

³ Available at <<https://github.com/rarrafel/vSLAM-dataset>>

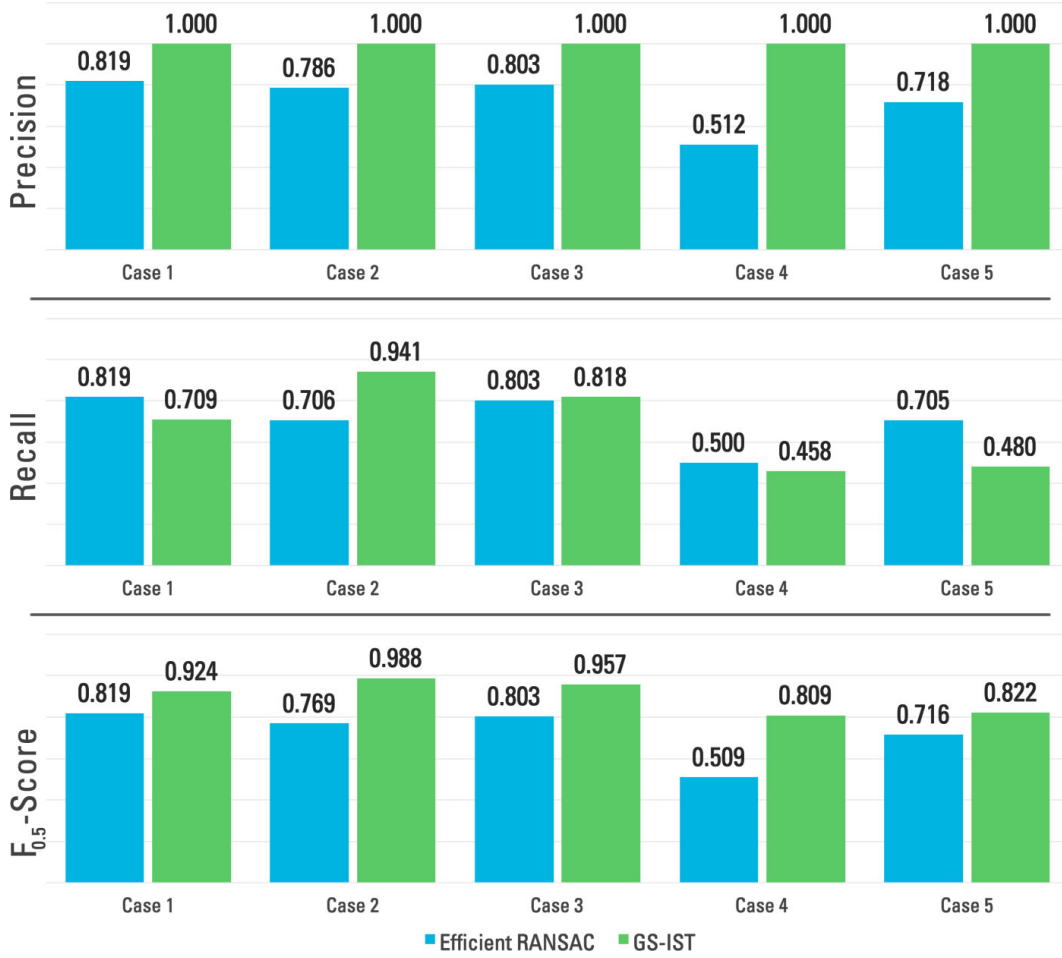


Figure 52 – Comparison of precision, recall and $F_{0.5}$ -Score between Efficient RANSAC (SCHNABEL; WAHL; KLEIN, 2007) and GS-IST.

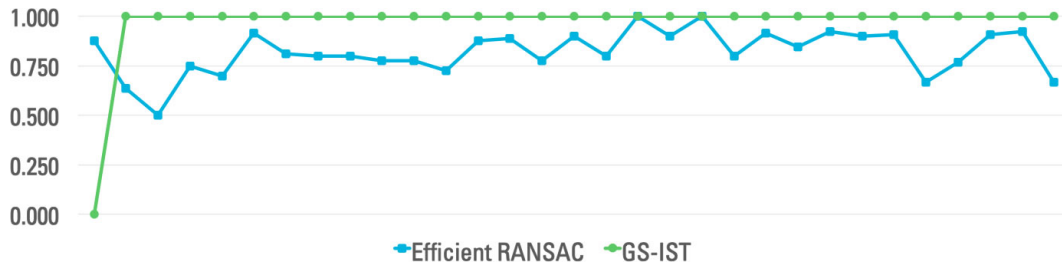


Figure 53 – Precision over time of Efficient RANSAC (SCHNABEL; WAHL; KLEIN, 2007) and GS-IST on Case 1.

Regarding the recall, Figure 52 displays that GS-IST is worse than Efficient RANSAC in three of the five cases. This happens because Efficient RANSAC outputs twice more shapes on average than GS-IST, even though some of them are incorrect. Thus, this method compromises recall in order to be entirely sure that the most reliable shapes are selected. Using the same bottle as an example, in the third keyframe the cylinder weight is 0.537, which is below the reliability threshold of 0.75. However, since it is above the 0.5 unreliability mark, it is performed a statistical analysis to verify the randomness of this detection. According to Equation 4.6, this estimation history is random and the shape is

assigned as unreliable.

Concerning $F_{0.5}$ -Score, Figure 52 shows that GS-IST presents a better result in all cases. The most significant improvement is in *Case 4*, which is very challenging because the reconstruction is very noisy. This evaluation indicates that the restriction imposed improved the precision, but with the cost of having few shapes detected. Therefore, it is possible to adjust the parameters to have more primitives and decrease the precision. These changes will depend on the target application. Table 9 provides some examples of possible modifications to make and the outcome for *Case 1*.

Table 9 – The influence of modifications in GS-IST on the final precision and recall in *Case 1*.

Condition Changed	Precision	Recall
<i>Remove geometrical analysis</i>	-1.409%	+1.846%
<i>Remove statistical analysis</i>	-1.409%	+3.139%
<i>Double elimination thresholds</i>	-0.704%	+2.602%

Figure 54 compares the results of both approaches in each test case. It shows situations in which the same object was detected as distinct shapes at different moments by Efficient RANSAC (rows 1 and 4). This phenomenon does not happen with GS-IST (rows 2 and 5). The detected shapes are represented by the projection of the input points used to compute the primitive. Rows 3 and 6 displays one view of the input point cloud in red and some of the estimated shapes with GS-IST in blue. From the last row, it is possible to see how challenging is *Case 4*. Although the books are aligned in real life, the points from the left one are not aligned with the other two.

4.4.1 Metric Evaluation

Case 1 has a chessboard pattern, which means that the reconstruction can be calibrated to the metric scale. This allowed an accuracy evaluation of object pose and parameters for this case in particular since there is no such pattern in the other cases. Considering the absence of ground truth for pose estimation, it was measured the average distance of each input point to its projection on the estimated primitive to assess how close they were. Figure 55 compares this distance over time between Efficient RANSAC and GS-IST. The leap in the distance is expected due to the error accumulation of the SLAM method. The average distance was 2.925 ± 0.370 mm for GS-IST while for Efficient RANSAC it was 3.247 ± 0.611 mm. It is worth to mention that this accuracy depends on the quality of the map. In this case, the error accumulation of the SLAM method was not precisely measured but it was around 15 mm, which is compatible with the error of the other systems evaluated in this Ph.D. thesis.

Concerning the parameter accuracy, it was used the radius of the globe and the wipe container, which are 50 and 40 mm respectively. The average radius of the sphere detected

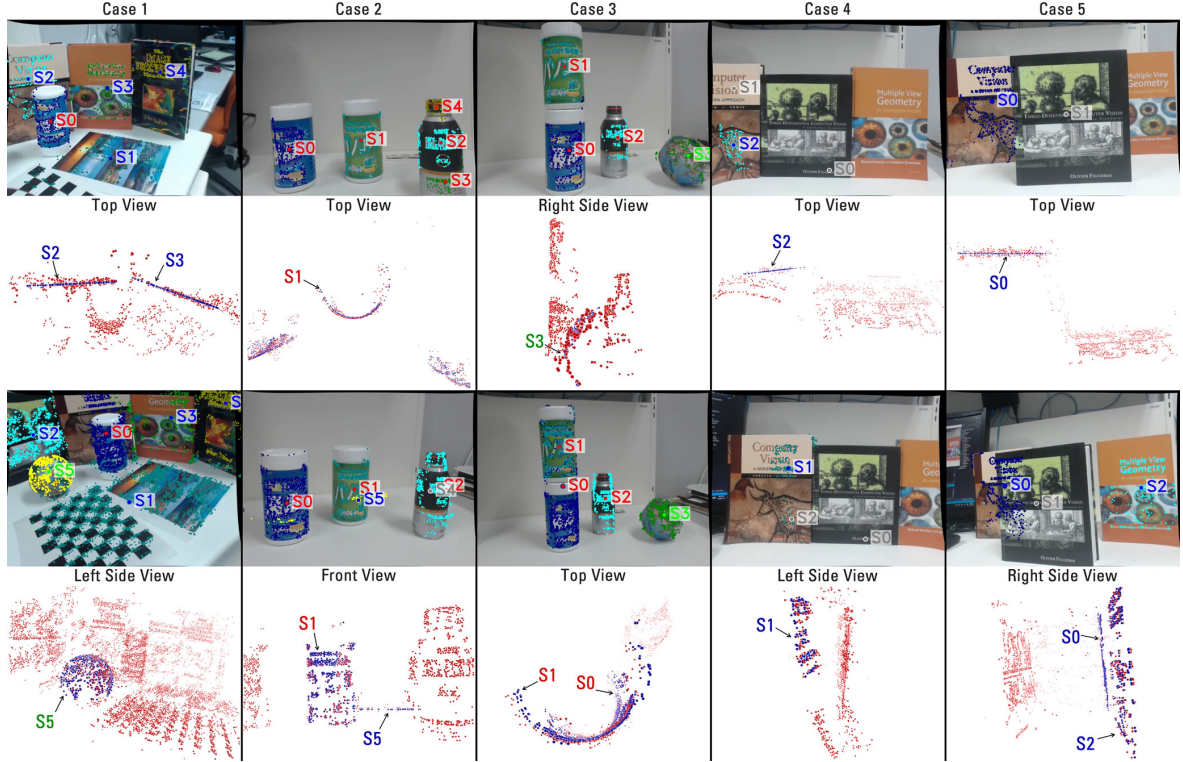


Figure 54 – The first and third rows show the result of GS-IST on each test case. Blue labels represent planes, green ones are for spheres and red for cylinders. The second and fourth rows show one particular view of the input point cloud (in red) and some of the estimated primitives (in blue)⁴.

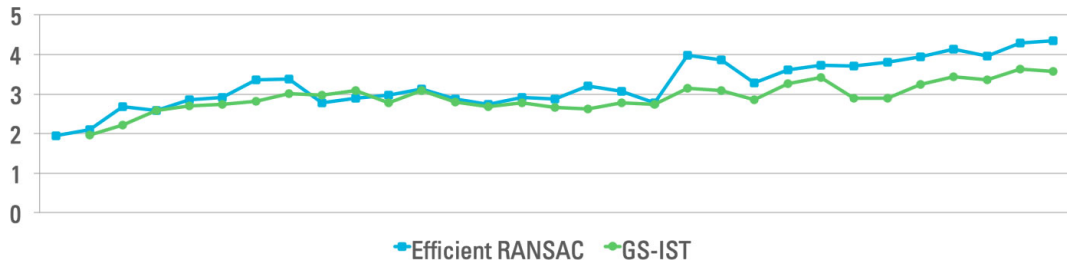


Figure 55 – Average distance in millimeters of each input point to its projection on the estimated primitive over time for *Case 1*.

with the globe points was 43.147 ± 0.318 mm, resulting in an error of 6.853 mm. As for the cylinder estimated based on the wipe container, the 33.723 ± 1.001 mm radius is 6.277 mm smaller than the real object. Figure 56 shows that in the first detection the cylinder radius is 31.051 mm and it gradually increases closer to the actual measurement over time, ending with 35.952 mm. These are the largest and smallest error in comparison with the ground truth for both the cylinder and the sphere. It can be credited to parameter update over time and the increase in the number of points that, even with the error accumulation, adds more data for the shape extractor. The sphere radius, on the other hand, decreases around 1 mm from the first to the last keyframe, going to the opposite direction of the actual radius. This can also be credited to error accumulation.

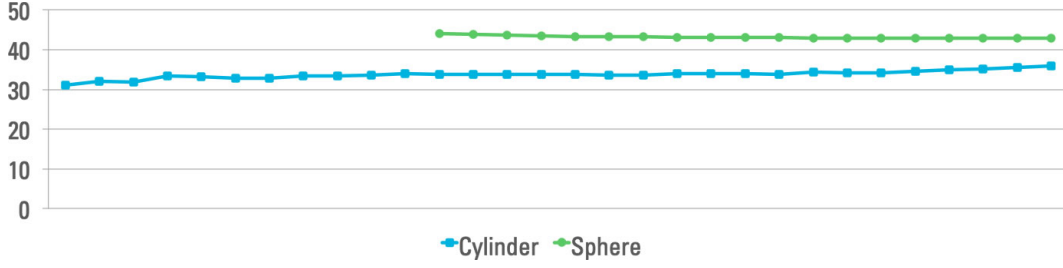


Figure 56 – Error in millimeters over time of the cylinder and the sphere radii detected in *Case 1*.

4.4.2 Runtime Evaluation

Concerning the computational cost, Efficient RANSAC takes on average 25.474 ± 10.380 ms to estimate the primitives in a computer with a Core i7-6820 (2.70 GHz) and 16GB of RAM. The other steps combined run in 14.995 ± 3.019 ms on average. The bottleneck is the *shape fusion* step, which takes 9.982 ± 2.957 ms of that time. It is worth mentioning that the execution time is related to the number of input points. Therefore, the measurements, which are the averages of all five test cases, were normalized to a group of thousand points.

4.4.3 Segmentation Evaluation

It was also evaluated how GS-IST segments the point cloud, which is a natural outcome of semantic modeling. Several objects have the form of the basic primitives tracked with this method. Looking at an average from all five test cases, 70.85% of the points can be assigned to a plane, sphere or cylinder. Even though the dataset deals with scenes designed with this type of primitives, the number is similar to the study that claims that 78% of all elements in an industrial scenario can be modeled using these three shapes (HUANG; YOU, 2013). Moreover, the chart on Figure 57 shows that only 6.30% of the remaining points were not labeled as any primitive. The other 22.85% points come from primitives that were discarded because they were unreliable.

4.4.4 Point Cloud Representation Evaluation

Finally, it was compared the scene representation using the point cloud and the modeled primitives. The scene is usually overrepresented when it is described using the points because there are many redundant points. It was measured the memory necessary to represent the reconstruction of each test case using the point cloud and it was compared with the description of the same map using the data structure of the primitives modeled with GS-IST. Figure 58 shows this difference in KB between then, ordered by the number of points from the reconstructed map. The most significant difference is in the last keyframe

⁴ A video with this result is also available at <<https://goo.gl/5RGrYm>>

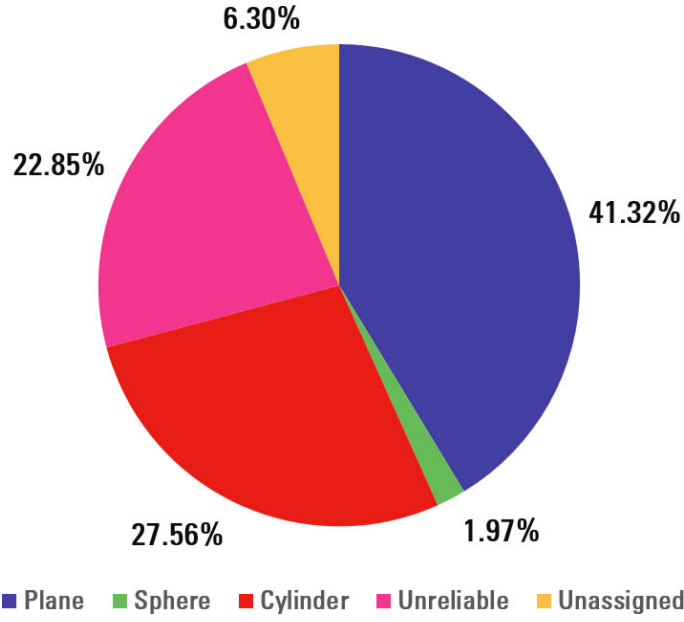
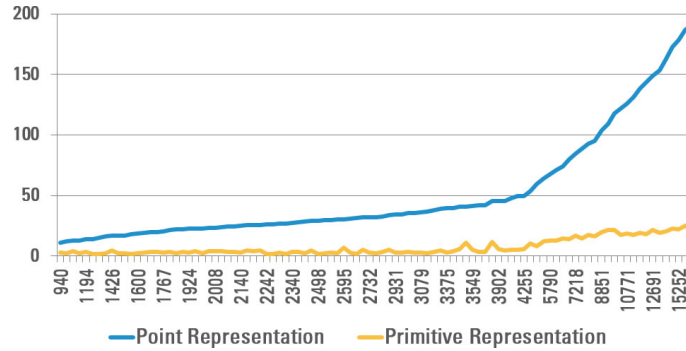


Figure 57 – Percentage of points that were labeled to each primitive.

of *Case 1*, which has 16,302 points. Using the point cloud requires 8.69 times more memory than describing the same map using the six detected primitives plus the 3,915 unreliable and unassigned points.

Figure 58 – Memory (in KB) required to describe a scene using the point cloud and the data structure of the detected primitives for *Case 1*.

Moreover, describing a scene using points commonly results in over and underrepresentation at the same time. For instance, considering only the cylinder detected in the last keyframe of *Case 1*, it can be noticed that there are more points than necessary to describe the textured front side but none to represent the back side. Thus, it is possible to use much less information to define this shape while filling the missing parts. Using this cylinder as an example, it is necessary 24 times more memory to describe it using the points than using the data structure of the detected primitive.

4.4.5 Proof of Concept

It is intended to see how GS-IST responds to an application that benefits from having semantic knowledge of the environment. It was developed Shape Hunt, a system to help children to identify some of the primitives they are learning in school. The idea of this proof of concept is that it draws a shape and he/she has to find real objects with the same geometric form. Since GS-IST has the scene map and can track the objects in the scene, it can identify all selected shapes and the child always has to find a new one. Figure 59 shows this proof of concept.

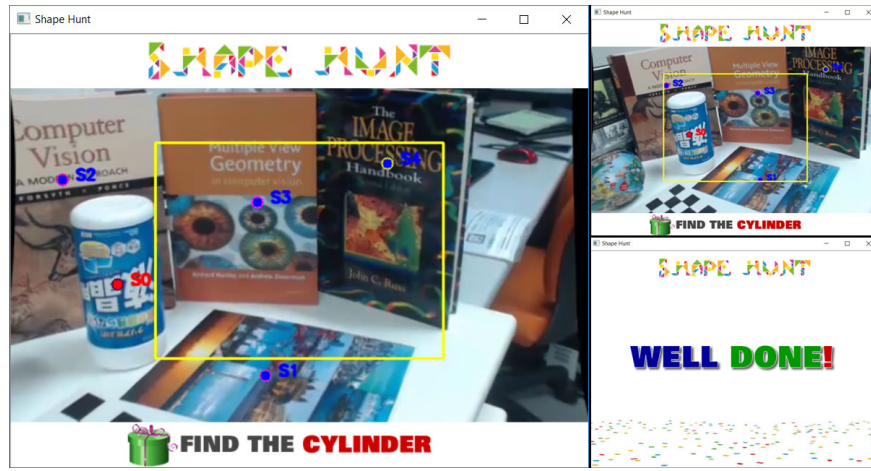


Figure 59 – The application indicates the shape the user has to find (left image). When the correspondent shape is centered (top-right), it gives a positive feedback (bottom-right) and moves to the next primitive⁵.

For this proof of concept, it was used the images from the dataset as input for this application. The scenes are limited to a small workspace, which was sufficient for this test. The application behaves as expected. Each primitive was detected only once and none was misclassified.

4.4.6 Evaluation on Dense Point Cloud

Although designed for sparse point clouds, the foundations of this tracker were adapted to work in a dense one. For this situation, it was also used the generating process of the point cloud to create a semantic map. The incremental process is also explored to improve the precision and stability of the shape detection over time. Moreover, the history of the fusing process is used to reduce the influence of error accumulation in the RGB-D SLAM system. Additionally, since the scale is available, it is possible to have the metric information from the modeled primitives. The evaluation showed that error on measuring the radius of spheres and cylinders varies between 0.1 and 4.6 millimeters and that it is difficult to model spheres that have the radius smaller than 30.0 mm. Similar to GS-IST, the execution

⁵ A video with this result is also available at <<https://goo.gl/d76Sfw>>

time is proportional to the number of points. However, the modeling process was executed in less than 100 milliseconds, on average, in a standard desktop. The complete details about the method adaptation for dense point clouds and its evaluation were published in (OLIVIER et al., 2018).

5 MOBILE EVALUATION OF INCREMENTAL SEMANTIC TRACKING

Mobile devices are reducing the gap to desktop computers in terms of processing power and memory space. From 2009 to 2015 the desktop CPU clock frequency increased by 33% while the mobile CPU clock grew 252% in the same period (HALPERN; ZHU; REDDI, 2016). This increment in the processing capabilities of mobile devices allowed the development of more complex algorithms, including tracking techniques, which is one of the foundations of augmented reality. This can be linked to the improvement and popularization of augmented reality solutions for such devices.

Some of these recent advancements in tracking techniques involves removing the necessity of any external marker, improving the execution time to achieve real-time performance and having a more stable tracking that is not harmed by jitter or drift. Most of them are thanks to the improvement of the device's sensors that allowed the development of visual-inertial trackers. However, less progress was made in regard to extracting any kind of semantics from the 3D map of the scene. This chapter focus on the evaluation of the incremental semantic tracker from Chapter 4 on mobile devices, showing that it is feasible to extract and track basic primitives on such platforms.

5.1 SEMANTIC MODELING ON MOBILE DEVICES

As seen in Chapter 2, tracking techniques on mobile devices had a substantial improvement until 2015. Not only on the number of publications but, most importantly, on tracking capabilities. Since then, the development of visual-inertial techniques become popular and received a lot of attention from the community, which comes in different flavors. There are distributed approaches that automatically selects between a fast visual tracker based on optical flow and a visual-inertial odometry depending on certain quality criteria (PIAO; KIM, 2017). Others are improving the traditional visual-inertial methods, such as adding a lightweight tightly-coupled fusion approach which integrates nonlinear optimization and loop detection (LI et al., 2017). More recently, (SOLIN et al., 2018) proposed a probabilistic inertial-visual odometry technique that is robust to occlusion and large-scale navigation without the need to perform loop closures. They achieve that by propagating the uncertainty of the measurements to every aspect of the tracking procedure, which includes the camera motion, the geometry and the minimization problem. However, none of these techniques are able to retrieve any type of semantics from the scene.

The recent development of machine learning techniques allowed the possibility to extract semantics from the image based on network architectures designed to deal with constraints of such devices (SANDLER et al., 2018). For instance, (TOBIAS et al., 2016) uses deep learning to perform domain-specific object recognition. They achieve high classification and near

real-time execution time running on an iPad. In another example, a Convolutional Neural Network (CNN) was incorporated to a mobile augmented reality application to perform object detection (RAO et al., 2017). They use the inertial sensors and the GPS to track an outdoor environment, detect geographical landmarks using this CNN and render their information with 3D coherence due to the inertial odometry. DeepCham (LI et al., 2016) use CNN to recognize objects as well. They rely on distributed redundancy, enhanced bounding box and generic deep model to facilitate the creation of training instances. There are examples of geometric techniques for object recognition too. One example is (CHEN et al., 2015), which created a distributed approach that uses a cache of frames scheme to improve the method performance. They are able to identify from faces to traffic signs.

Although most of the studies find semantics from images, it is also possible to retrieve these information from point clouds. This is more complex to achieve with mobile devices due to the overload of data to process, especially when dealing with dense point sets. However, recently some systems were developed that are able to extract information from such data. These techniques usually generate the point cloud using Google Tango devices, such as (RUNCEANU et al., 2017) that apply an iterative RANSAC approach to segment planes and model walls of indoor structures. Using a Tango as well, (SANKAR; SEITZ, 2017) also detect planes and model indoor environments. Nevertheless, they also detect planes to segment more complex objects and based on their set of points and visual appearance the authors are able to match the object with a 3D model available in the library.

Extracting semantic information using sparse point clouds should be simpler when concerning the amount of data to process. However, this is the same reason why this task is more challenging. Nevertheless, both Google's ARCore and Apple's ARKit incorporate semantic modeling as part of their scene understanding feature. They both extract planes based on the scene 3D reconstruction for creating surfaces in order to have a more stable positioning of the virtual objects.

5.2 MOBILE IMPLEMENTATION

In order to evaluate how GS-IST would perform running on mobile devices, the technique presented in the last chapter was ported to the Android platform. Since it was developed in C++, this facilitates the port using the efficient Android architecture in a similar process used to port LBF and STAM to mobile devices. It was also necessary to compile the libraries used in the project to Android: OpenCV and Boost¹. Efficient RANSAC was treated as a library as well but the compilation process was a little bit different due to the fact that it was added to the project in order to facilitate modifications in the source code.

¹ Available at <<https://www.boost.org>>

5.2.1 Evaluation

The two most critical aspects of mobile devices are the energy consumption and the temperature. Although the number of cores and CPU clock have increased lately, the processors' architecture compromise on speed to be more efficient on these two facets (REINER, 2012). Most of the studies aiming mobile devices focus their evaluation only on execution time along with any qualitative assessment that is adequate for the proposed method. Using the 10 papers cited in the previous section and the 25 relevant studies listed in Table 3 at Chapter 2 as a sample, 62.9% measured execution time and 45.7% evaluated only this criteria. Energy consumption was evaluated in 17.1%, RAM memory usage in 8.6% and 28.6% performed qualitative assessments alone. Only (LI et al., 2016) evaluates execution time, energy consumption and memory usage. However, they did not detail the methodology used to perform these measurements. These three criteria were used to evaluate GS-IST along with CPU load, which is a good indicator of the potential of parallel execution of a certain application.

This evaluation was executed in devices with different capabilities and from distinct manufacturers, which is important to assess how GS-IST performs in dissimilar conditions. The chosen devices were the Samsung Galaxy S8 and the ASUS ZenFone 3. They were selected using ARCore as a reference. The Galaxy S8 is in the list of the supported devices² while the ZenFone 3 was selected to stress GS-IST since it has a configuration inferior to those in that list. Table 10 show some of the technical specifications of these smartphones. All the tests were executed with the device fully charged, on airplane mode, with all other applications closed and connected to the computer via the USB cable. The only exception was the energy measurement in which the device was disconnected from the computer.

Table 10 – Summary of the technical specifications of the devices used for evaluation.

Features	Galaxy S8	ZenFone 3
<i>Android Version</i>	8.0 (Oreo)	7.0 (Nougat)
<i>Display</i>	5.8" (1440 x 2960)	5.5" (1080 x 1920)
<i>Chipset</i>	Qualcomm MSM8998 Snapdragon 835	Qualcomm MSM8953 Snapdragon 625
<i>CPU</i>	Kyro Octa-core (4x2.35 GHz and 4x1.9 GHz)	Cortex-A53 Octa-core (2.0 GHz)
<i>Memory</i>	64 GB, 4 GB RAM	32 GB, 3 GB RAM
<i>Battery</i>	3000 mAh	3000 mAh

Regarding the dataset, it was used the same five scenes generated in the previous chapter. The images and pose files were stored in the device and loaded on every frame. The same happened for the map files, which were loaded on every keyframe. Since the

² Available at <<https://developers.google.com/ar/discover/supported-devices>>

codes are identical, precision and recall are equal to the ones presented in the last chapter. Figure 60 shows a few keyframes of GS-IST running on both phones.

5.2.1.1 Execution Time

The execution time is proportional to the number of points processed and all time measurements, which are the averages of all five test cases, were normalized to a group of thousand points. Figure 61 shows that the average execution time of GS-IST on ZenFone 3 is 9.9 times slower in comparison with the desktop implementation and it is 8.5 times slower on Galaxy S8. For this test, the same desktop computer with Core i7-6820 (2.70 GHz) processor and 16GB of RAM was used. The *shape fusion* was slower than the average on mobile devices, being the desktop implementation 16.5 (ZenFone 3) and 15.0 (Galaxy S8) times faster.

The CPU load is an important measure because it indicates how much room the GS-IST leaves to perform other processing, such as the SLAM technique. For that evaluation, it was used Qualcomm's Treppn Profiler³. This application, available in the Play Store, samples the desired information in a constant time interval. In this test, both the CPU load and the Normalized CPU load were sampled every 100 milliseconds. The operating system imposes a limit on how much processing an application can use. The CPU load represents how much of that limit is being used by the application while the Normalized CPU load indicates how much processing is being used in relation to the total processing power of the device.

It is possible to observe in Figure 62 that GS-IST presents some execution peaks. These apexes coincide with the keyframes, which are moments in which the technique extracts the primitives. The maximum normalized load for the ZenFone 3 was 90% of CPU, similar to the 89% value for the Galaxy S8. The average normalized load was also similar for both devices, in which the Samsung device was $30.545\% \pm 25.119\%$ of Normalized CPU load while the ASUS mobile phone presented $27.128\% \pm 16.353\%$. This happens because for most of the time the execution is in between keyframes, a moment in which the device is not processing much data. The median of the Normalized CPU load is a numerical indication for that and it was 19% and 21% for the Galaxy S8 and ZenFone 3 respectively. However, there was a difference in the average CPU load, which was $39.283\% \pm 27.131\%$ for the Samsung phone and $46.483\% \pm 25.040\%$ for the ASUS device. Smaller differences between the CPU load and the Normalized CPU load suggests that the device is allowed to use the full potential of the processor. In that case, this value was 8.738% for the Galaxy S8 but it was 19.355% for the ZenFone 3.

³ Available at <<https://developer.qualcomm.com/software/treppn-power-profiler>>

⁴ A video with this result is also available at <<https://goo.gl/A6T51d>>



Figure 60 – Results of GS-IST running on a Samsung Galaxy S8 and an ASUS ZenFone 3. Blue labels represent planes, green ones are for spheres and red for cylinders⁴.

5.2.1.2 Energy Consumption

The most precise method to evaluate energy consumption is by using external instruments that can directly measure the current drained by the device. However, these equipment require opening the device to be attached to the physical battery, which is difficult for most smartphones nowadays since their batteries are not easily accessible. An alternative

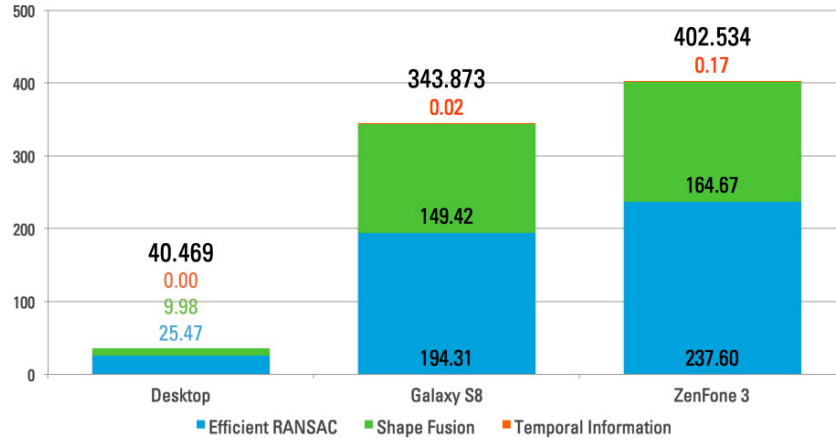


Figure 61 – GS-IST execution time in milliseconds divided by stages on desktop and two different mobile devices.

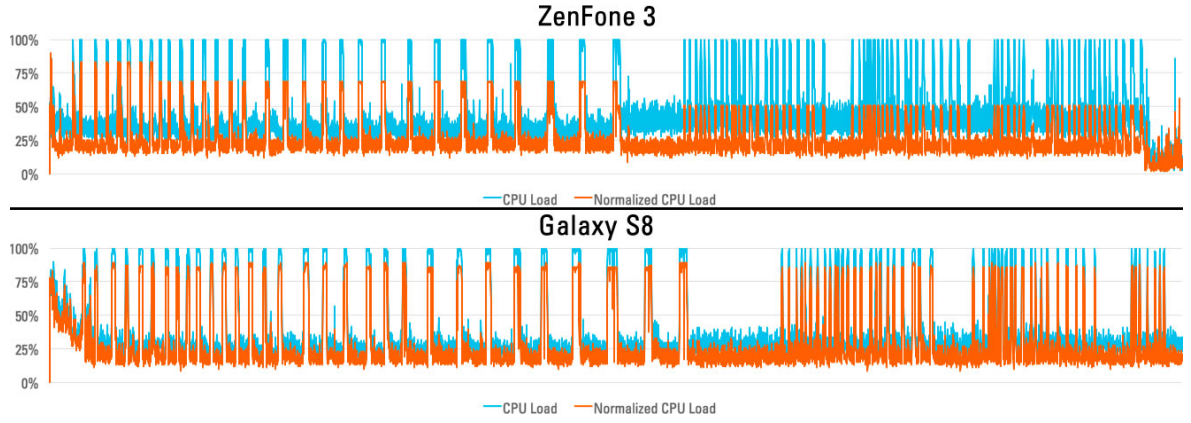


Figure 62 – CPU load and Normalized CPU load over time of all five test cases on a ZenFone 3 (top) and Galaxy S8 (bottom).

is to use profiler tools that use the battery API to assess the voltage and the state of charge at certain intervals. This procedure is much more accessible but it is not as accurate as using these external instruments. The latter approach was selected for this evaluation and Trepro Profiler was also used for this task. Qualcomm’s tool has an accuracy of 99%, which is reported to be one of the highest for such profilers (HOQUE et al., 2015).

As expected, Figure 63 shows that energy consumption on both mobile devices follow the same pattern of the CPU load since more energy is required in those most computational-intense moments. The average consumption on the ZenFone 3 was 1.776 ± 1.162 W every 100 ms. Since the battery capacity of this device is 3000 mAh with 3.85 V, this means that this device could run GS-IST for 5 hours and 12 minutes before drain all the battery when it is fully charged considering an energy efficiency of 80% (VALOEN; SHOESMITH, 2007). The Galaxy S8 battery has the same characteristics (3000 mAh and 3.85 V), which determines that its average 2.271 ± 1.998 W consumption would drain a fully charged battery in 4 hours and 04 minutes.

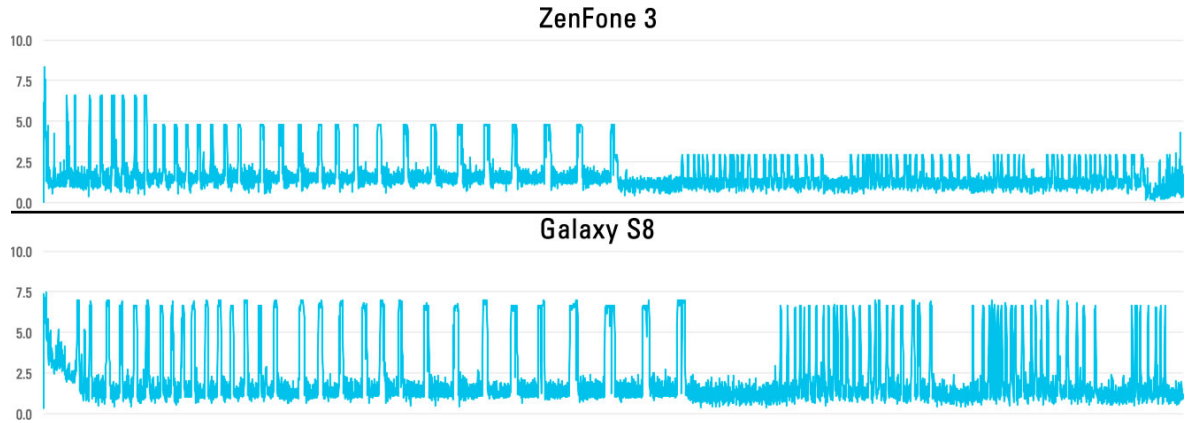


Figure 63 – Energy consumption (in W) over time of all five test cases on ZenFone 3 (top) and Galaxy S8 (bottom).

5.2.1.3 Memory Usage

Concerning the memory usage, two measurements can be done. One is the storage space the sample APK requires when installed and the other is the RAM memory it uses when running. The former value can easily be found in the device settings. In the ZenFone 3, GS-IST used 45.14 MB of the storage space while in the Galaxy S8 it occupied 47.38 MB.

To evaluate the latter it was used the Android Profiler available on Android Studio, which builds a chart with the RAM memory usage as the application is executed. Similar to the previous evaluation, Figure 64 shows that the RAM memory has some peaks when extracting the primitives. For the ZenFone 3, the moment with most memory usage is in the 28th keyframe of *Case 1* with 195.39 MB, which represents 6.4% of the device total memory. When not processing the keyframes, the sample app consumes between 26.14 MB and 38.19 MB.

As seen in Figure 65, the memory usage for the Galaxy S8 is similar, although with higher absolute values. The memory peak was 322.97 MB in the 25th keyframe of *Case 1*, being 7.9% of the phone RAM memory. The memory consumption when not extracting primitives was above 119.07 MB and below 151.16 MB.

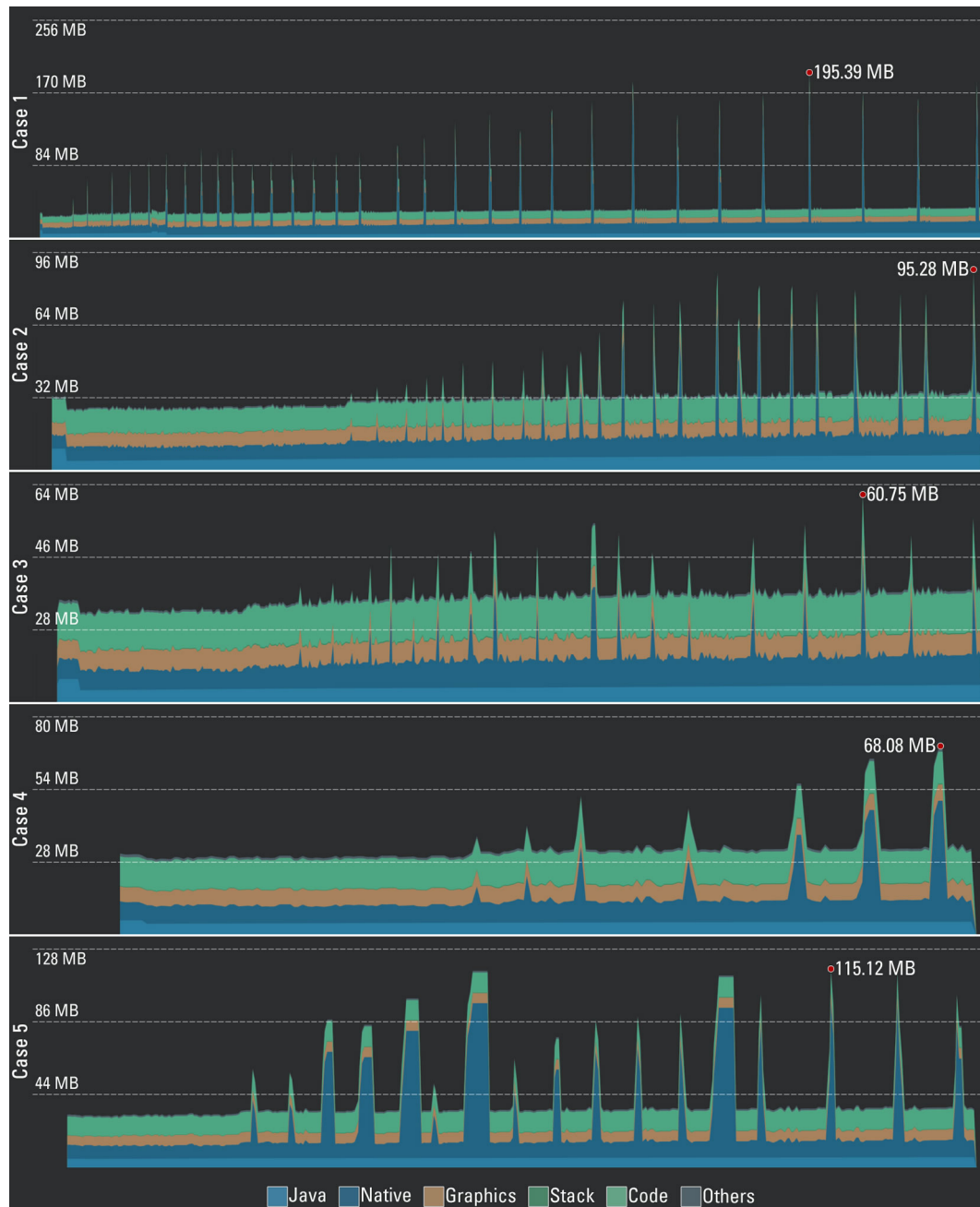


Figure 64 – Memory usage over time of GS-IST running on ZenFone 3. Each test case has a different time scale.



Figure 65 – Memory usage over time of GS-IST running on Galaxy S8. Each test case has a different time scale.

5.3 DISCUSSION

From Figure 61 it is possible to see that GS-IST does not run in real-time on any of the devices it was tested. It is especially true when the number of points increases. However, this approach uses less than a third of the CPU power regardless of the device. This means that GS-IST and a SLAM technique can run in different threads in order that the semantic modeling system and the SLAM method can interact with each other without compromising their performance. Since this mobile version is running on Android devices, it is possible to use ARCore as a SLAM Method and extract the primitives from the point cloud the SDK generates. Nevertheless, it is necessary to perform an evaluation to see if the ARCore map is too sparse. Additionally, there is room for optimization in the implementation.

Moreover, Figure 62 shows that Android 7.0 on ZenFone 3 imposed a more restrict limit on how much processing GS-IST could use. In fact, that limit decreased over time. For instance, GS-IST could use 90% of the CPU at the beginning of the execution and that boundary decreased to approximately 75% in the 9th keyframe of *Case 1* and to 50% from *Case 2* and beyond. On the other hand, Android 8.0 on Galaxy S8 allowed GS-IST to use around 90% of the CPU from beginning to end.

Regarding energy consumption, the 4 hours GS-IST needs to drain the battery in the worst case is sufficient to use this technique without having extra concerns about the battery. Moreover, considering the habits and the average time users spend on mobile devices (Hacker Noon, 2017), it is unlikely that a person would use a specific app for a period so long. In order to put this in perspective, Table 11 compares the time other common applications take to fully discharge the battery, such as watching a 1080p video, navigating using Google Maps and playing FIFA Soccer. Besides the video activity, GS-IST has an energy consumption similar to the other applications.

Table 11 – Time that different applications take to fully drain a battery fully charged on ZenFone 3 and Galaxy S8 devices.

Application	ZenFone 3	Galaxy S8
<i>1080p video</i>	8h07min	6h31min
<i>Google Maps (using 4G network)</i>	6h32min	4h29min
<i>GS-IST</i>	5h12min	4h04min
<i>Fifa Soccer (using WiFi connection)</i>	5h09min	4h47min

The evaluation showed a noticeable difference in RAM memory usage between both devices. The ASUS phone used approximately 100 MB less memory than the Samsung one. It is possible to see in Figure 64 and Figure 65 that graphics uses much more memory in Galaxy S8 than in ZenFone 3. The graphics are responsible for more than 70% of that difference. It was not found any details for that. However, based on observation using

other devices, one possibility is that the ASUS device delegate the rendering activity to the GPU, therefore graphics-related structures uses the GPU memory.

The memory usage of GS-IST was not a concern since, in the worst case, it used less than 8% of the total RAM memory of the device. This is due to the fact that modern smartphones have a fair amount of RAM memory. For them, replacing the representation from point clouds to primitives does not have any impact on the execution time. However, this change in representation can be important in some situations. Complex algorithms can use a lot of memory and perform several computations for each element in the scene, which can overload the powerful devices even for a sparse map. Tracking optimization methods have that characteristic. In fact, some developers reported that the current version of ARCore (v1.2.1) can run out of memory when it has to perform bundle adjustment with a map whose size is that of a large room. Rendering algorithms are another case in which having a large number of elements can cause the usage of most of the device's resources. Therefore, a more efficient representation can be very helpful in circumstances like that.

6 FINAL CONSIDERATIONS

Tracking on mobile devices evolved a lot since this Ph.D. started four and half years ago. At that time the phones were not so powerful, which raised big concerns regarding the processing power and memory capabilities required to run real-time tracking of such devices. The perception was that there was a growing interest on having tracking techniques on mobile devices but a few were able to track more than a planar template. The systematic mapping conducted confirmed that perception. Moreover, it also showed that until 2015, most of the works used the device's sensors to compute $2D$ or $2D + \theta$ pose on location-based applications. Nowadays, the devices improved in great extension and new techniques have emerged, which make possible the development of increasingly popular SDKs that are able to perform real-time tracking with $6D$ poses.

Based on the findings of the systematic mapping, it was created experimental scenarios aiming to evaluate different tracking approaches for mobile devices. One of these experiments was the proposition of a method to assess the Google Tango platform aiming to establish a reference of the state-of-the-art trackers. It was observed that the motion tracking errors were around 6 and 14 centimeters for small and large indoor scenarios, which is suitable to provide a good user experience, including for augmented reality applications. This experiment was followed by another one that evaluated the use of different forms to develop computer vision systems on mobile devices, such as using parallelism and distributed approach. It indicated that each solution has strengths and limitations depending on the situation. However, native development was the most efficient on average. It was performed experiments to evaluate different tracking techniques that had the potential to be suitable for mobile devices. The first was a face tracking technique using machine learning and local binary features, which was adapted to consider the characteristics of mobile devices, such as camera orientation. The second was a SLAM technique that was originally developed in desktop and ported to a Tango tablet device. The mobile version presented some issues regarding performance, especially when it needs to process a large number of data in situations like bundle adjustment.

One of the main lessons learned in this Ph.D. study was the importance of finding high-level semantic information from a scene. Therefore, it was developed a technique that detects and tracks primitives, called GS-IST. This method uses the generating process of sparse point clouds of visual SLAM systems and applies geometrical and statistical analyses to incrementally estimate and track planes, spheres and cylinders. The evaluation indicated that GS-IST improved precision in all test cases, which outperformed existing methods in this criteria. The developed approach focuses on precision and for that, it compromises recall to assure we have the correct shapes. However, we can modify the parameters to increase recall when necessary. Additionally, this technique was ported

to the Android platform and evaluated to assess how it performed running on mobile devices. The evaluation showed that the mobile version is slower when compared with the desktop implementation but it can be executed on a separate thread of the SLAM technique because the CPU load is not so high. Finally, the energy consumption and memory usage were not a concern.

6.1 FUTURE WORK

There is still some work that needs to be done after this Ph.D course. One is to investigate the possibility of integrating GS-IST with Google's ARCore, which would allow the development of use cases based on real-world problems that can benefit from having a semantic knowledge of the scene. One way to do that is using a well-structured methodology that is based on the design thinking theory that combines interdisciplinary teams to conceive innovative solutions (ROBERTO et al., 2016). There is a team at Voxar Labs creating bridges between academic research and innovative solutions with high impact and the goal is to team up with them to run an instance of their process.

After integrating the semantic tracker with ARCore, it is possible to create new test cases that allow the evaluation on more complex environments. These new scenes can be industrial factories that have several machines, mechanical workshops loaded with equipment, or warehouses, which have various shelves and boxes. Although very challenging, these scenarios have several objects that can be modeled with the primitives GS-IST detects. Therefore, this would allow stressing the technique in order to find points for further improvements.

Further activities include performing a more extensive evaluation of the accuracy of object tracking and its parameters. In order to accomplish that, it is necessary to have a dataset with ground truth pose and measurements. This can be achieved with the creation of other scenes that would include a chessboard pattern or other means to recover the scale. It is also possible to obtaining ground truth for object pose using markers tracked with libraries such as ArUco (GARRIDO-JURADO et al., 2014). An additional possibility is to create a synthetic case, downsample and apply noisy to the point cloud in order to simulate sparse reconstructions.

Another idea for future work is to use the semantic knowledge of the scene to improve the tracking results of the visual SLAM system. There are some ways to achieve this goal. One is to constrain the map 3D points during the bundle adjustment to move only over the surface of the shape it belongs, which would optimize the point cloud respecting the semantic structure and leading to a faster convergence. A different strategy is to optimize the map using the primitive parameters instead of the points. The idea is to save on computation by minimizing the error of a few parameters rather than do the same operation for hundreds of points that represent the same elements on the scene. This would have an impact on the scalability of the environment to be tracked.

6.2 CONTRIBUTIONS

This Ph.D. research produced some contributions to the community:

- A systematic mapping that extensively cataloged and classified the area of tracking for mobile devices, providing a reference for new researchers in the field to have a quick overview of the area;
- A paper catalog adapted from an open-source system, which is an easy way to display the results of the mapping;
- An open-source paper crawler¹, which can help other researchers to gather scientific papers;
- An evaluation of different architectures aiming to efficiently develop tracking techniques on mobile devices;
- A face tracker system that was part of an application developed in a project in partnership with a major mobile phone manufacturer;
- A SLAM technique that was developed and used to compete in the 2015 ISMAR Tracking Competition, winning the first place on the “On-Site Category: Level 3”;
- The method used to evaluate Google Tango, which is an easy way to perform preliminary evaluations on mobile motion trackers;
- A technique that uses geometric and statistical evaluation to incrementally perform semantic modeling and tracking of primitives on sparse point clouds;
- A dataset with sparse point clouds of primitives that is helpful to evaluate semantic modeling and tracking²;
- The port of the incremental semantic tracker to mobile devices;
- A guideline to evaluate computer vision techniques on mobile devices.

6.3 PUBLICATIONS

Some scientific papers were published during this Ph.D. Seven were directly related to this study, as mentioned in the text. One was related to the systematic mapping (ROBERTO; LIMA; TEICHRIB, 2016). Three were about the experiments, one being the evaluation of Google Tango (ROBERTO et al., 2016a), other was the partial results of the experiment about the architecture of computer vision systems on Android (LIMA et al., 2015) and the

¹ Available at <http://www.cin.ufpe.br/~rar3/tracking_sm/paper-analysis.tar.gz>

² Available at <<https://github.com/rarrafel/vSLAM-dataset>>

third one was a publication about the port of STAM to Google’s Project Tango (ARAUJO et al., 2016). Finally, three more papers were published concerning GS-IST (ROBERTO et al., 2017; ROBERTO et al., 2018; OLIVIER et al., 2018). Additionally, this Ph.D. study was selected to be presented and discussed with the computer vision community in the *PhD Forum* of the IEEE Winter Conference on Applications of Computer Vision (WACV) in 2018.

There were also four publications targeting tracking (LIMA et al., 2017) or mobile devices (LINS et al., 2014a; LINS et al., 2014b; LIMA et al., 2014) that are not directly related to this work but were important to gain experience in the Ph.D. topics. Additionally, there were seven papers regarding topics not directly related to this study (MOTA et al., 2014; SILVA et al., 2015; SILVA; ROBERTO; TEICHRIEB, 2015; MOTA; ROBERTO; TEICHRIEB, 2015; ROBERTO et al., 2016; SILVA et al., 2016; ROBERTO et al., 2016b).

REFERENCES

- AGARWAL, S.; MIERLE, K.; OTHERS. *Ceres Solver*. <<http://ceres-solver.org>>.
- ALMEIDA, D. R. O. d.; GUEDES, P. A.; SILVA, M. M. O. d.; SILVA, A. L. B. V. e.; LIMA, J. P. S. d. M.; TEICHRIEB, V. Interactive makeup tutorial using face tracking and augmented reality on mobile devices. In: *Virtual and Augmented Reality (SVR), 2015 XVII Symposium on*. [S.l.: s.n.], 2015. p. 220–226.
- ALTMAN, D. *Practical Statistics for Medical Research*. Taylor & Francis, London, 1990. (Chapman & Hall/CRC Texts in Statistical Science). ISBN 9780412276309. Disponível em: <<https://books.google.com.br/books?id=v-walRnRxWQC>>.
- AN, J. H.; HONG, K. S. Finger gesture-based mobile user interface using a rear-facing camera. In: *Consumer Electronics (ICCE), 2011 IEEE International Conference on*. [S.l.: s.n.], 2011. p. 303–304. ISSN 2158-3994.
- ANDO, B.; BAGLIO, S.; LOMBARDO, C.; MARLETTA, V. An advanced tracking solution fully based on native sensing features of smartphone. In: *Sensors Applications Symposium (SAS), 2014 IEEE*. [S.l.: s.n.], 2014. p. 141–144.
- Apple Inc. *ARKit - Apple Developer*. 2017. [Online; last access: 29-Dec-2017]. Disponível em: <<https://developer.apple.com/arkit/>>.
- ARAUJO, T.; ROBERTO, R.; TEIXEIRA, J. M.; SIMÕES, F.; TEICHRIEB, V.; LIMA, J. P.; ARRUDA, E. Life cycle of a slam system: Implementation, evaluation and port to the project tango device. In: *Virtual and Augmented Reality (SVR), 2016 XVIII Symposium on*. [S.l.: s.n.], 2016.
- ARTH, C.; WAGNER, D.; KLOPSCHITZ, M.; IRSCHARA, A.; SCHMALSTIEG, D. Wide area localization on mobile phones. In: *Mixed and Augmented Reality, 2009. ISMAR 2009. 8th IEEE International Symposium on*. [S.l.: s.n.], 2009. p. 73–82.
- BAI, H.; GAO, L.; BILLINGHURST, M. Poster: Markerless fingertip-based 3d interaction for handheld augmented reality in a small workspace. In: *3D User Interfaces (3DUI), 2013 IEEE Symposium on*. [S.l.: s.n.], 2013. p. 129–130.
- BAY, H.; TUYTELAARS, T.; GOOL, L. V. Surf: Speeded up robust features. In: LEONARDIS, A.; BISCHOF, H.; PINZ, A. (Ed.). *Computer Vision – ECCV 2006*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006. p. 404–417. ISBN 978-3-540-33833-8.
- BRADSKI, G. The OpenCV Library. *Dr. Dobbs's Journal of Software Tools*, 2000.
- CHEN, D.; TSAI, S.; VEDANTHAM, R.; GRZESZCZUK, R.; GIROD, B. Streaming mobile augmented reality on mobile phones. In: *Mixed and Augmented Reality, 2009. ISMAR 2009. 8th IEEE International Symposium on*. [S.l.: s.n.], 2009. p. 181–182.
- CHEN, T. Y.-H.; RAVINDRANATH, L.; DENG, S.; BAHL, P.; BALAKRISHNAN, H. Glimpse: Continuous, real-time object recognition on mobile devices. In: *Proceedings of the 13th ACM Conference on Embedded Networked Sensor Systems*. New York, NY, USA: ACM, 2015. (SenSys '15), p. 155–168. ISBN 978-1-4503-3631-4. Disponível em: <<http://doi.acm.org/10.1145/2809695.2809711>>.

- CHENG, S.; ASTHANA, A.; ZAFEIRIOU, S.; SHEN, J.; PANTIC, M. Real-time generic face tracking in the wild with cuda. In: *Proceedings of the 5th ACM Multimedia Systems Conference*. New York, NY, USA: ACM, 2014. (MMSys '14), p. 148–151. ISBN 978-1-4503-2705-3. Disponível em: <<http://doi.acm.org/10.1145/2557642.2579369>>.
- CHON, Y.; TALIPOV, E.; CHA, H. Autonomous management of everyday places for a personalized location provider. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, v. 42, n. 4, p. 518–531, July 2012. ISSN 1094-6977.
- COHEN, J. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, Durham, v. 20, n. 1, p. 37–46, 1960.
- DAQRI. *DAQRI Smart Helmet – DAQRI*. 2016. [Online; last access: 07-May-2016]. Disponível em: <<http://daqri.com/home/product/daqri-smart-helmet/>>.
- DROST, B.; ILIC, S. Local hough transform for 3d primitive detection. In: *2015 International Conference on 3D Vision*. [S.l.: s.n.], 2015. p. 398–406.
- ELLOUMI, W.; GUISSOUS, K.; CHETOUANI, A.; CANALS, R.; LECONGE, R.; EMILE, B.; TREUILLET, S. Indoor navigation assistance with a smartphone camera based on vanishing points. In: *Indoor Positioning and Indoor Navigation (IPIN), 2013 International Conference on*. [S.l.: s.n.], 2013. p. 1–9.
- ENGELKE, T.; KEIL, J.; ROJTBERG, P.; WIENTAPPER, F.; WEBEL, S.; BOCKHOLT, U. Content first - a concept for industrial augmented reality maintenance applications using mobile devices. In: *Mixed and Augmented Reality (ISMAR), 2013 IEEE International Symposium on*. [S.l.: s.n.], 2013. p. 251–252.
- FIGUEIREDO, L. S.; PINHEIRO, M.; NETO, E. V.; CHAVES, T.; TEICHRIEB, V. Human-computer interaction – interact 2015: 15th ifip tc 13 international conference, bamberg, germany, september 14-18, 2015, proceedings, part ii. In: _____. Cham: Springer International Publishing, 2015. cap. Sci-Fi Gestures Catalog, p. 395–411. ISBN 978-3-319-22668-2.
- FISCHLER, M. A.; BOLLES, R. C. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, ACM, New York, NY, USA, v. 24, n. 6, p. 381–395, jun. 1981. ISSN 0001-0782. Disponível em: <<http://doi.acm.org/10.1145/358669.358692>>.
- GARRIDO-JURADO, S.; MUÑOZ-SALINAS, R.; MADRID-CUEVAS, F.; MARÍN-JIMÉNEZ, M. Automatic generation and detection of highly reliable fiducial markers under occlusion. *Pattern Recognition*, v. 47, n. 6, p. 2280 – 2292, 2014. ISSN 0031-3203. Disponível em: <<http://www.sciencedirect.com/science/article/pii/S0031320314000235>>.
- GHERGHINA, A.; OLTEANU, A.; TAPUS, N. A marker-based augmented reality system for mobile devices. In: *Roedunet International Conference (RoEduNet), 2013 11th*. [S.l.: s.n.], 2013. p. 1–6. ISSN 2068-1038.
- GIVEN, L. M. *The Sage encyclopedia of qualitative research methods*. [S.l.]: Sage Publications, 2008.

- Google Inc. *Tango* / *Google Developers*. 2014. [Online; last access: 29-Dec-2017]. Disponível em: <<https://developers.google.com/tango/>>.
- Google Inc. *ARCore - Google Developer* / *ARCore* / *Google Developers*. 2017. [Online; last access: 29-Dec-2017]. Disponível em: <<https://developers.google.com/ar/>>.
- GOZICK, B.; SUBBU, K.; DANTU, R.; MAESHIRO, T. Magnetic maps for indoor navigation. *Instrumentation and Measurement, IEEE Transactions on*, v. 60, n. 12, p. 3883–3891, Dec 2011. ISSN 0018-9456.
- GRUBERT, J.; LANGLLOTZ, T.; GRASSET, R. Augmented reality browser survey. *Institute for Computer Graphics and Vision, University of Technology Graz, Technical Report*, n. 1101, 2011.
- HA, J.; CHO, K.; ROJAS, F.; YANG, H. Real-time scalable recognition and tracking based on the server-client model for mobile augmented reality. In: *VR Innovation (ISVRI), 2011 IEEE International Symposium on*. [S.l.: s.n.], 2011. p. 267–272.
- Hacker Noon. *How Much Time Do People Spend on Their Mobile Phones in 2017?* 2017. [Online; last access: 21-May-2018]. Disponível em: <<https://hackernoon.com/how-much-time-do-people-spend-on-their-mobile-phones-in-2017-e5f90a0b10a6>>.
- HADID, A.; HEIKKILA, J. Y.; SILVEN, O.; PIETIKAINEN, M. Face and eye detection for person authentication in mobile phones. In: *2007 First ACM/IEEE International Conference on Distributed Smart Cameras*. [S.l.: s.n.], 2007. p. 101–108.
- HAGBI, N.; BERGIG, O.; EL-SANA, J.; BILLINGHURST, M. Shape recognition and pose estimation for mobile augmented reality. *Visualization and Computer Graphics, IEEE Transactions on*, v. 17, n. 10, p. 1369–1379, Oct 2011. ISSN 1077-2626.
- HALPERN, M.; ZHU, Y.; REDDI, V. J. Mobile cpu's rise to power: Quantifying the impact of generational mobile cpu design trends on performance, energy, and user satisfaction. In: *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. [S.l.: s.n.], 2016. p. 64–76.
- HARTLEY, R.; ZISSERMAN, A. *Multiple view geometry in computer vision*. [S.l.]: Cambridge university press, 2003.
- HERLING, J.; BROLL, W. Random model variation for universal feature tracking. In: *Proceedings of the 18th ACM Symposium on Virtual Reality Software and Technology*. New York, NY, USA: ACM, 2012. (VRST '12), p. 169–176. ISBN 978-1-4503-1469-5.
- HETTIARACHCHI, A.; WIGDOR, D. Annexing reality: Enabling opportunistic use of everyday objects as tangible proxies in augmented reality. In: *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. New York, NY, USA: ACM, 2016. (CHI '16), p. 1957–1967. ISBN 978-1-4503-3362-7. Disponível em: <<http://doi.acm.org/10.1145/2858036.2858134>>.
- HODANĚ, T.; MATAS, J.; OBDRŽÁLEK, Š. On evaluation of 6d object pose estimation. In: SPRINGER. *European Conference on Computer Vision*. [S.l.], 2016. p. 606–619.

- HOLZ, D.; HOLZER, S.; RUSU, R. B.; BEHNKE, S. Real-time plane segmentation using rgb-d cameras. In: _____. *RoboCup 2011: Robot Soccer World Cup XV*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012. p. 306–317. ISBN 978-3-642-32060-6. Disponível em: <https://doi.org/10.1007/978-3-642-32060-6_26>.
- HOQUE, M. A.; SIEKKINEN, M.; KHAN, K. N.; XIAO, Y.; TARKOMA, S. Modeling, profiling, and debugging the energy consumption of mobile devices. *ACM Comput. Surv.*, ACM, New York, NY, USA, v. 48, n. 3, p. 39:1–39:40, dez. 2015. ISSN 0360-0300. Disponível em: <<http://doi.acm.org/10.1145/2840723>>.
- HU, C.-L.; CHO, C.-A.; LIN, C.-J.; FAN, C.-W. A location tracking and messaging system for mobile group communication in ip networks. In: *Consumer Electronics (ICCE), 2010 Digest of Technical Papers International Conference on*. [S.l.: s.n.], 2010. p. 155–156.
- HU, W.; LIAN, S.; SONG, X.; LI, T. Mobile camera based cross-screen interaction by object matching and tracking. *Consumer Electronics, IEEE Transactions on*, v. 59, n. 3, p. 452–459, August 2013. ISSN 0098-3063.
- HUANG, J.; YOU, S. Detecting objects in scene point cloud: A combinational approach. In: *2013 International Conference on 3D Vision - 3DV 2013*. [S.l.: s.n.], 2013. p. 175–182. ISSN 1550-6185.
- HUANG, Z.; HUI, P.; PEYLO, C.; CHATZOPOULOS, D. Mobile augmented reality survey: A bottom-up approach. *arXiv preprint arXiv:1309.4413*, 2013.
- Intel Corporation. *Threading Building Blocks*. 2016. [Online; last access: 14-May-2016]. Disponível em: <<https://www.threadingbuildingblocks.org>>.
- International Symposium on Mixed and Augmented Reality. *ISMAR 2015 Tracking Competition*. 2015. [Online; last access: 22-Jun-2016]. Disponível em: <<http://ypcex.naist.jp/trakmark/tracking-competition/>>.
- ISSARTEL, P.; GUENIAT, F.; AMMI, M. Slicing techniques for handheld augmented reality. In: *3D User Interfaces (3DUI), 2014 IEEE Symposium on*. [S.l.: s.n.], 2014. p. 39–42.
- JAIN, R. *The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling*. Wiley, 1991. ISBN 9780471503361. Disponível em: <<https://books.google.com.br/books?id=CN1QAAAAMAAJ>>.
- JUNG, H.; LEE, Y.-D.; JEONG, D.-U. A novel method for energy expenditure using multisensor based activity monitoring. In: *Computer Sciences and Convergence Information Technology (ICCIT), 2011 6th International Conference on*. [S.l.: s.n.], 2011. p. 103–106.
- JURIE, F.; DHOME, M. Real time robust template matching. In: *BMVC*. [S.l.: s.n.], 2002. p. 1–10.
- KAO, W.-W.; HUY, B. Q. Indoor navigation with smartphone-based visual slam and bluetooth-connected wheel-robot. In: *Automatic Control Conference (CACS), 2013 CACS International*. [S.l.: s.n.], 2013. p. 395–400.

- KAZEMI, V.; SULLIVAN, J. One millisecond face alignment with an ensemble of regression trees. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition*. [S.l.: s.n.], 2014. p. 1867–1874. ISSN 1063-6919.
- KEELE, S. *Guidelines for performing systematic literature reviews in software engineering*. [S.l.], 2007.
- Khronos Group. *OpenGL - The Industry Standard for High Performance Graphics*. 1997. [Online; last access: 22-Jun-2016]. Disponível em: <<https://www.opengl.org>>.
- KIM, Y. M.; DOLSON, J.; SOKOLSKY, M.; KOLTUN, V.; THRUN, S. Interactive acquisition of residential floor plans. In: *2012 IEEE International Conference on Robotics and Automation*. [S.l.: s.n.], 2012. p. 3055–3062. ISSN 1050-4729.
- KITCHENHAM, B.; BRERETON, P.; BUDGEN, D. The educational value of mapping studies of software engineering literature. In: *Software Engineering, 2010 ACM/IEEE 32nd International Conference on*. [S.l.: s.n.], 2010. v. 1, p. 589–598. ISSN 0270-5257.
- KLEIN, G.; MURRAY, D. Parallel tracking and mapping for small ar workspaces. In: *Proceedings of the 2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*. Washington, DC, USA: IEEE Computer Society, 2007. (ISMAR '07), p. 1–10. ISBN 978-1-4244-1749-0. Disponível em: <<http://dx.doi.org/10.1109/ISMAR.2007.4538852>>.
- KLEIN, G.; MURRAY, D. Parallel tracking and mapping on a camera phone. In: *Mixed and Augmented Reality, 2009. ISMAR 2009. 8th IEEE International Symposium on*. [S.l.: s.n.], 2009. p. 83–86.
- KUEHNI, R. G. *Color space and its divisions: color order from antiquity to the present*. [S.l.]: John Wiley & Sons, 2003.
- KURZ, D.; BENHIMANE, S. Gravity-aware handheld augmented reality. In: *Mixed and Augmented Reality (ISMAR), 2011 10th IEEE International Symposium on*. [S.l.: s.n.], 2011. p. 111–120.
- KURZ, D.; BENHIMANE, S. Handheld augmented reality involving gravity measurements. *Computers & Graphics*, v. 36, n. 7, p. 866 – 883, 2012. ISSN 0097-8493. Augmented Reality Computer Graphics in China.
- LAMBRECHT, J.; WALZEL, H.; KRUGER, J. Robust finger gesture recognition on handheld devices for spatial programming of industrial robots. In: *RO-MAN, 2013 IEEE*. [S.l.: s.n.], 2013. p. 99–106. ISSN 1944-9445.
- LANE, N.; MILUZZO, E.; LU, H.; PEEBLES, D.; CHOUDHURY, T.; CAMPBELL, A. A survey of mobile phone sensing. *Communications Magazine, IEEE*, v. 48, n. 9, p. 140–150, Sept 2010. ISSN 0163-6804.
- LEONARD, J. J.; DURRANT-WHYTE, H. F. Mobile robot localization by tracking geometric beacons. *IEEE Transactions on Robotics and Automation*, v. 7, n. 3, p. 376–382, Jun 1991. ISSN 1042-296X.

- LEPETIT, V.; MORENO-NOGUER, F.; FUA, P. Epnp: An accurate $O(n)$ solution to the pnp problem. *Int. J. Comput. Vision*, Kluwer Academic Publishers, Hingham, MA, USA, v. 81, n. 2, p. 155–166, fev. 2009. ISSN 0920-5691. Disponível em: <<http://dx.doi.org/10.1007/s11263-008-0152-6>>.
- LESHED, G.; VELDEN, T.; RIEGER, O.; KOT, B.; SENGERS, P. In-Car GPS navigation: Engagement with and disengagement from the environment. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. New York, NY, USA: ACM, 2008. (CHI '08), p. 1675–1684. ISBN 978-1-60558-011-1.
- LI, D.; SALONIDIS, T.; DESAI, N. V.; CHUAH, M. C. Deepcham: Collaborative edge-mediated adaptive deep learning for mobile object recognition. In: *2016 IEEE/ACM Symposium on Edge Computing (SEC)*. [S.l.: s.n.], 2016. p. 64–76.
- LI, M.; KIM, B. H.; MOURIKIS, A. Real-time motion tracking on a cellphone using inertial sensing and a rolling-shutter camera. In: *Robotics and Automation (ICRA), 2013 IEEE International Conference on*. [S.l.: s.n.], 2013. p. 4712–4719. ISSN 1050-4729.
- LI, M.; MOURIKIS, A. I. Vision-aided inertial navigation for resource-constrained systems. In: *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*. [S.l.: s.n.], 2012. p. 1057–1063. ISSN 2153-0858.
- LI, P.; QIN, T.; HU, B.; ZHU, F.; SHEN, S. Monocular visual-inertial state estimation for mobile augmented reality. In: *2017 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. [S.l.: s.n.], 2017. p. 11–21.
- LI, Y.; WU, X.; CHRYSATHOU, Y.; SHARF, A.; COHEN-OR, D.; MITRA, N. J. Globfit: Consistently fitting primitives by discovering global relations. *ACM Trans. Graph.*, ACM, New York, NY, USA, v. 30, n. 4, p. 52:1–52:12, jul. 2011. ISSN 0730-0301. Disponível em: <<http://doi.acm.org/10.1145/2010324.1964947>>.
- LIEBERKNECHT, S.; BENHIMANE, S.; MEIER, P.; NAVAB, N. A dataset and evaluation methodology for template-based tracking algorithms. In: *Mixed and Augmented Reality, 2009. ISMAR 2009. 8th IEEE International Symposium on*. [S.l.: s.n.], 2009. p. 145–151.
- LIMA, J. P.; ROBERTO, R.; SIMões, F.; ALMEIDA, M.; FIGUEIREDO, L.; TEIXEIRA, J. M.; TEICHRIEB, V. Markerless tracking system for augmented reality in the automotive industry. *Expert Systems with Applications*, v. 82, p. 100 – 114, 2017. ISSN 0957-4174. Disponível em: <<http://www.sciencedirect.com/science/article/pii/S0957417417302221>>.
- LIMA, J. P.; ROBERTO, R.; TEICHRIEB, V.; MARQUES, G. Study about natural feature tracking for augmented reality applications on mobile devices. In: *Virtual and Augmented Reality (SVR), 2015 XVII Symposium on*. [S.l.: s.n.], 2015. p. 7–14.
- LIMA, J. P.; ROBERTO, R.; TEIXEIRA, J. M.; TEICHRIEB, V. [poster] device vs. user perspective rendering in google glass ar applications. In: *Mixed and Augmented Reality (ISMAR), 2014 IEEE International Symposium on*. [S.l.: s.n.], 2014. p. 279–280.
- LINS, C.; ARRUDA, E.; NETO, E.; ROBERTO, R.; TEICHRIEB, V.; FREITAS, D.; TEIXEIRA, J. M. Animar: Augmenting the reality of storyboards and animations. In: *Virtual and Augmented Reality (SVR), 2014 XVI Symposium on*. [S.l.: s.n.], 2014. p. 106–109.

- LINS, C.; TEIXEIRA, J. M.; ROBERTO, R.; TEICHRIEB, V. Development of interactive applications for google glass. *Tendências e Técnicas em Realidade Virtual e Aumentada*, v. 4, p. 167 – 188, 2014. ISSN 2177-6776.
- LIU, H.; DARABI, H.; BANERJEE, P.; LIU, J. Survey of wireless indoor positioning techniques and systems. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, v. 37, n. 6, p. 1067–1080, Nov 2007. ISSN 1094-6977.
- LIU, Y. J.; ZHANG, J. B.; HOU, J. C.; REN, J. C.; TANG, W. Q. Cylinder detection in large-scale point cloud of pipeline plant. *IEEE Transactions on Visualization and Computer Graphics*, v. 19, n. 10, p. 1700–1707, Oct 2013. ISSN 1077-2626.
- LOPEZ-ESCOGIDO, D.; FRAGA, L. G. de la. Automatic extraction of geometric models from 3d point cloud datasets. In: *2014 11th International Conference on Electrical Engineering, Computing Science and Automatic Control (CCE)*. [S.l.: s.n.], 2014. p. 1–5.
- LOURAKIS, M. I. A.; ARGYROS, A. A. Sba: A software package for generic sparse bundle adjustment. *ACM Trans. Math. Softw.*, ACM, New York, NY, USA, v. 36, n. 1, p. 2:1–2:30, mar. 2009. ISSN 0098-3500. Disponível em: <<http://doi.acm.org/10.1145/1486525.1486527>>.
- LOWE, D. G. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, Kluwer Academic Publishers, Hingham, MA, USA, v. 60, n. 2, p. 91–110, nov. 2004. ISSN 0920-5691. Disponível em: <<http://dx.doi.org/10.1023/B:VISI.0000029664.99615.94>>.
- LUCAS, B. D.; KANADE, T. An iterative image registration technique with an application to stereo vision. In: *Proceedings of the 7th International Joint Conference on Artificial Intelligence - Volume 2*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1981. (IJCAI'81), p. 674–679. Disponível em: <<http://dl.acm.org/citation.cfm?id=1623264.1623280>>.
- LV, Z. Wearable smartphone: Wearable hybrid framework for hand and foot gesture interaction on smartphone. In: *Computer Vision Workshops (ICCVW), 2013 IEEE International Conference on*. [S.l.: s.n.], 2013. p. 436–443.
- MA, Y.; SOATTO, S.; KOSECKA, J.; SASTRY, S. *An Invitation to 3-D Vision: From Images to Geometric Models*. Springer New York, 2005. (Interdisciplinary Applied Mathematics). ISBN 9780387008936. Disponível em: <<https://books.google.com.br/books?id=6tUqQmwan4UC>>.
- MARTIN, P.; MARCHAND, E.; HOULIER, P.; MARCHAL, I. Decoupled mapping and localization for augmented reality on a mobile phone. In: *Virtual Reality (VR), 2014 IEEE*. [S.l.: s.n.], 2014. p. 97–98.
- MICHAEL, K.; CLARKE, R. Location and tracking of mobile devices: überveillance stalks the streets. *Computer Law & Security Review*, v. 29, n. 3, p. 216 – 228, 2013. ISSN 0267-3649.
- Microsoft. *Microsoft HoloLens / The leader in mixed reality technology*. 2016. [Online; last access: 29-Dec-2017]. Disponível em: <<https://www.microsoft.com/en-us/hololens>>.

- MORENCY, L.-P. 3d constrained local model for rigid and non-rigid facial tracking. In: *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Washington, DC, USA: IEEE Computer Society, 2012. (CVPR '12), p. 2610–2617. ISBN 978-1-4673-1226-4. Disponível em: <<http://dl.acm.org/citation.cfm?id=2354409.2354947>>.
- MOTA, R.; MARQUES, G.; SANTOS, A.; ARAUJO, J.; ROBERTO, R.; TEICHRIEB, V. Blockit: Tangible interfaces for music education. In: *Workshop em Realidade Virtual e Aumentada (WRVA)*. [S.l.: s.n.], 2014.
- MOTA, R. C.; ROBERTO, R. A.; TEICHRIEB, V. [poster] authoring tools in augmented reality: An analysis and classification of content design tools. In: *Mixed and Augmented Reality (ISMAR), 2015 IEEE International Symposium on*. [S.l.: s.n.], 2015. p. 164–167.
- MUR-ARTAL, R.; MONTIEL, J. M. M.; TARDOS, J. D. ORB-SLAM: A versatile and accurate monocular slam system. *IEEE Transactions on Robotics*, v. 31, n. 5, p. 1147–1163, Oct 2015. ISSN 1552-3098.
- NGUYEN, G.; ANDERSEN, H.; HOILUND, C. Street navigation using visual information on mobile phones. In: *Intelligent Systems Design and Applications (ISDA), 2010 10th International Conference on*. [S.l.: s.n.], 2010. p. 37–42.
- NGUYEN, T.; REITMAYR, G.; SCHMALSTIEG, D. Structural modeling from depth images. *IEEE Transactions on Visualization and Computer Graphics*, v. 21, n. 11, p. 1230–1240, Nov 2015. ISSN 1077-2626.
- NORMAND, J.-M.; MOREAU, G. Dof-based classification of augmented reality applications. In: *IEEE ISMAR workshop “Classifying the AR presentation space”*. [S.l.: s.n.], 2012. p. 1–8.
- OEHLER, B.; STUECKLER, J.; WELLE, J.; SCHULZ, D.; BEHNKE, S. Efficient multi-resolution plane segmentation of 3d point clouds. In: _____. *Intelligent Robotics and Applications: 4th International Conference, ICIRA 2011, Aachen, Germany, December 6-8, 2011, Proceedings, Part II*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011. p. 145–156. ISBN 978-3-642-25489-5. Disponível em: <https://doi.org/10.1007/978-3-642-25489-5_15>.
- OLIVIER, N.; UCHIYAMA, H.; MISHIMA, M.; THOMAS, D.; TANIGUCHI, R.-i.; ROBERTO, R.; LIMA, J. P.; TEICHRIEB, V. Live structural modeling using rgb-d slam. Unpublished paper accepted at the 2018 IEEE International Conference on Robotics and Automation (ICRA). 2018.
- OLSSON, T.; SALO, M. Online user survey on current mobile augmented reality applications. In: *Mixed and Augmented Reality (ISMAR), 2011 10th IEEE International Symposium on*. [S.l.: s.n.], 2011. p. 75–84.
- OpenSignal, Inc. *The State of LTE February 2016 - OpenSignal*. 2016. [Online; last access: 22-Jun-2016]. Disponível em: <<http://opensignal.com/reports/2016/02/state-of-lte-q4-2015/>>.
- OUI, W.; NG, E.; KHAN, R. An augmented reality’s framework for mobile. In: *Information Technology and Multimedia (ICIM), 2011 International Conference on*. [S.l.: s.n.], 2011. p. 1–4.

- PANG, G.; NEUMANN, U. Training-based object recognition in cluttered 3d point clouds. In: *2013 International Conference on 3D Vision - 3DV 2013*. [S.l.: s.n.], 2013. p. 87–94. ISSN 1550-6185.
- PANG, G.; QIU, R.; HUANG, J.; YOU, S.; NEUMANN, U. Automatic 3d industrial point cloud modeling and recognition. In: *2015 14th IAPR International Conference on Machine Vision Applications (MVA)*. [S.l.: s.n.], 2015. p. 22–25.
- PETERSEN, K.; FELDT, R.; MUJTABA, S.; MATTSSON, M. Systematic mapping studies in software engineering. In: *Proceedings of the 12th International Conference on Evaluation and Assessment in Software Engineering*. Swinton, UK, UK: British Computer Society, 2008. (EASE' 08), p. 68–77.
- PETIT, A.; CARON, G.; UCHIYAMA, H.; MARCHAND, E. Evaluation of Model based Tracking with TrakMark Dataset. In: *2nd Int. Workshop on AR/MR Registration, Tracking and Benchmarking*. Basel, Switzerland, Switzerland: [s.n.], 2011. Disponível em: <<https://hal.inria.fr/hal-00639700>>.
- Pew Research Center. *Smartphone Ownership and Internet Usage Continues to Climb in Emerging Economies* / Pew Research Center. 2016. [Online; last access: 15-May-2018]. Disponível em: <<http://www.pewglobal.org/2016/02/22/smartphone-ownership-and-internet-usage-continues-to-climb-in-emerging-economies/>>.
- PIAO, J.-C.; KIM, S.-D. Adaptive monocular visual-inertial slam for real-time augmented reality applications in mobile devices. *Sensors*, v. 17, n. 11, 2017. ISSN 1424-8220. Disponível em: <<http://www.mdpi.com/1424-8220/17/11/2567>>.
- PIRCHHEIM, C.; REITMAYR, G. Homography-based planar mapping and tracking for mobile phones. In: *Mixed and Augmented Reality (ISMAR), 2011 10th IEEE International Symposium on*. [S.l.: s.n.], 2011. p. 27–36.
- PIRCHHEIM, C.; SCHMALSTIEG, D.; REITMAYR, G. Handling pure camera rotation in keyframe-based slam. In: *Mixed and Augmented Reality (ISMAR), 2013 IEEE International Symposium on*. [S.l.: s.n.], 2013. p. 229–238.
- POLO, A.; VIANI, F.; GIAROLA, E.; OLIVERI, G.; ROCCA, P.; MASSA, A. Semantic wireless localization enabling advanced services in museums. In: *Antennas and Propagation (EuCAP), 2014 8th European Conference on*. [S.l.: s.n.], 2014. p. 443–446.
- QIU, R.; ZHOU, Q.-Y.; NEUMANN, U. Pipe-run extraction and reconstruction from point clouds. In: _____. *Computer Vision – ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part III*. Cham: Springer International Publishing, 2014. p. 17–30. ISBN 978-3-319-10578-9. Disponível em: <https://doi.org/10.1007/978-3-319-10578-9_2>.
- Qualcomm Technologies, Inc. *Mobile Multimedia Optimization - Mobile Technologies - Qualcomm Developer Network*. 2015. [Online; last access: 03-February-2016]. Disponível em: <<https://developer.qualcomm.com/software/hexagon-dsp-sdk>>.
- RAI, H.; DEEPAK, K.; SYED, S.; KRISHNA, P. A smart mobile application for identifying storage location of small industrial assets. In: *Mobile Data Management (MDM), 2012 IEEE 13th International Conference on*. [S.l.: s.n.], 2012. p. 332–335.

- RAJ, C.; TOLETY, S.; IMMACULATE, C. Qr code based navigation system for closed building using smart phones. In: *Automation, Computing, Communication, Control and Compressed Sensing (iMac4s), 2013 International Multi-Conference on*. [S.l.: s.n.], 2013. p. 641–644.
- RAMADASAN, D.; CHATEAU, T.; CHEVALDONNÉ, M. Dcslam: A dynamically constrained real-time slam. In: *2015 IEEE International Conference on Image Processing (ICIP)*. [S.l.: s.n.], 2015. p. 1130–1134.
- RAMANAN, D. Face detection, pose estimation, and landmark localization in the wild. In: *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Washington, DC, USA: IEEE Computer Society, 2012. (CVPR '12), p. 2879–2886. ISBN 978-1-4673-1226-4. Disponível em: <<http://dl.acm.org/citation.cfm?id=2354409.2355119>>.
- RAO, J.; QIAO, Y.; REN, F.; WANG, J.; DU, Q. A mobile outdoor augmented reality method combining deep learning object detection and spatial relationships for geovisualization. In: *Sensors*. [S.l.: s.n.], 2017.
- REINER, H. The paramountcy of reconfigurable computing. In: _____. *Energy-Efficient Distributed Computing Systems*. Wiley-Blackwell, 2012. cap. 18, p. 465–547. ISBN 9781118342015. Disponível em: <<https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118342015.ch18>>.
- REN, S.; CAO, X.; WEI, Y.; SUN, J. Face alignment at 3000 fps via regressing local binary features. In: *2014 IEEE Conference on Computer Vision and Pattern Recognition*. [S.l.: s.n.], 2014. p. 1685–1692. ISSN 1063-6919.
- ROBERTO, R.; FREITAS, D.; SIMOES, F.; TEICHRIEB, V. A dynamic blocks platform based on projective augmented reality and tangible interfaces for educational activities. *Journal on Interactive Systems, SBC*, v. 4, n. 2, p. 8–18, 2013. ISSN 2236-3297.
- ROBERTO, R.; LIMA, J. P.; ARAÚJO, T.; TEICHRIEB, V. Evaluation of motion tracking and depth sensing accuracy of the tango tablet. In: *2016 IEEE International Symposium on Mixed and Augmented Reality (ISMAR-Adjunct)*. [S.l.: s.n.], 2016. p. 231–234.
- ROBERTO, R.; LIMA, J. P.; MOTA, R.; TEICHRIEB, V. Authoring tools for augmented reality: An analysis and classification of content design tools. In: _____. *Design, User Experience, and Usability: Interactive Experience Design: 5th International Conference, DUXU 2016, Held as Part of HCI International 2016*. [S.l.]: Springer International Publishing, 2016.
- ROBERTO, R.; LIMA, J. P.; TEICHRIEB, V. Tracking for mobile devices: A systematic mapping study. *Computers & Graphics*, v. 56, p. 20 – 30, 2016. ISSN 0097-8493. Disponível em: <<http://www.sciencedirect.com/science/article/pii/S0097849316300115>>.
- ROBERTO, R.; LIMA, J. P.; UCHIYAMA, H.; ARTH, C.; TEICHRIEB, V.; TANIGUCHI, R. i.; SCHMALSTIEG, D. Incremental structural modeling based on geometric and statistical analyses. In: *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. [S.l.: s.n.], 2018. p. 955–963.

- ROBERTO, R.; SILVA, M. da; FREITAS, D.; LIMA, Y.; SILVA, V.; ARAUJO, C.; TEIXEIRA, J. M.; TEICHRIEB, V. Voxar puzzle motion: An innovative ar application proposed using design techniques. In: *K-12 Embodied Learning through Virtual & Augmented Reality (KELVAR), IEEE Virtual Reality 2016 Workshop on*. [S.l.: s.n.], 2016.
- ROBERTO, R. A.; UCHIYAMA, H.; LIMA, J. P. S. M.; NAGAHARA, H.; TANIGUCHI, R.-i.; TEICHRIEB, V. Incremental structural modeling on sparse visual slam. *IPSI Transactions on Computer Vision and Applications*, v. 9, n. 1, p. 5, Mar 2017. ISSN 1882-6695. Disponível em: <<https://doi.org/10.1186/s41074-017-0018-3>>.
- ROBINSON, S.; JONES, M.; VARTIAINEN, E.; MARSDEN, G. Picotales: Collaborative authoring of animated stories using handheld projectors. In: *Proceedings of the ACM 2012 Conference on Computer Supported Cooperative Work*. New York, NY, USA: ACM, 2012. (CSCW '12), p. 671–680. ISBN 978-1-4503-1086-4.
- ROTH, G.; LEVINE, M. D. Extracting geometric primitives. *CVGIP: Image Underst.*, Academic Press, Inc., Orlando, FL, USA, v. 58, n. 1, p. 1–22, jul. 1993. ISSN 1049-9660. Disponível em: <<http://dx.doi.org/10.1006/ciun.1993.1028>>.
- ROY, A.; ZHANG, X.; WOLLEB, N.; QUINTERO, C. P.; JÄGERSAND, M. Tracking benchmark and evaluation for manipulation tasks. In: *Robotics and Automation (ICRA), 2015 IEEE International Conference on*. [S.l.: s.n.], 2015. p. 2448–2453.
- RUBLEE, E.; RABAUD, V.; KONOLIGE, K.; BRADSKI, G. Orb: An efficient alternative to sift or surf. In: *2011 International Conference on Computer Vision*. [S.l.: s.n.], 2011. p. 2564–2571. ISSN 1550-5499.
- RUNCEANU, L.; BECKER, S.; HAALA, N.; FRITSCH, D. Indoor point cloud segmentation for automatic object interpretation. *Publikationen der Deutschen Gesellschaft für Photogrammetrie*, p. 147–159, 2017. Disponível em: <<https://www.semanticscholar.org/paper/Indoor-Point-Cloud-Segmentation-for-Automatic-Runceanu-Becker/181fd778e36b935479a3eb33abe2d68dbc876f12>>.
- SANDLER, M.; HOWARD, A.; ZHU, M.; ZHMOGINOV, A.; CHEN, L.-C. Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation. *arXiv preprint arXiv:1801.04381*, 2018.
- SANKAR, A.; SEITZ, S. M. Interactive room capture on 3d-aware mobile devices. In: *Proceedings of the 30th Annual ACM Symposium on User Interface Software and Technology*. New York, NY, USA: ACM, 2017. (UIST '17), p. 415–426. ISBN 978-1-4503-4981-9. Disponível em: <<http://doi.acm.org/10.1145/3126594.3126629>>.
- SANTOS, J.; RODRIGUES, F.; OLIVEIRA, L. A web & mobile city maintenance reporting solution. *Procedia Technology*, v. 9, n. 0, p. 226 – 235, 2013. ISSN 2212-0173. CENTERIS 2013 - Conference on ENTERprise Information Systems / ProjMAN 2013 - International Conference on Project MANagement/ HCIST 2013 - International Conference on Health and Social Care Information Systems and Technologies.
- SCHALL, G.; MULLONI, A.; REITMAYR, G. North-centred orientation tracking on mobile phones. In: *Mixed and Augmented Reality (ISMAR), 2010 9th IEEE International Symposium on*. [S.l.: s.n.], 2010. p. 267–268.

- SCHNABEL, R.; WAHL, R.; KLEIN, R. Efficient RANSAC for point-cloud shape detection. *Computer Graphics Forum*, Blackwell Publishing Ltd, v. 26, n. 2, p. 214–226, 2007. ISSN 1467-8659. Disponível em: <<http://dx.doi.org/10.1111/j.1467-8659.2007.01016.x>>.
- SCHÖPS, T.; ENGEL, J.; CREMERS, D. Semi-dense visual odometry for ar on a smartphone. In: *Mixed and Augmented Reality (ISMAR), 2014 IEEE International Symposium on*. [S.l.: s.n.], 2014. p. 145–150.
- SCHÖPS, T.; SATTTLER, T.; HÄNE, C.; POLLEFEYS, M. Large-scale outdoor 3d reconstruction on a mobile device. *Computer Vision and Image Understanding*, v. 157, p. 151 – 166, 2017. ISSN 1077-3142. Large-Scale 3D Modeling of Urban Indoor or Outdoor Scenes from Images and Range Scans. Disponível em: <<http://www.sciencedirect.com/science/article/pii/S1077314216301412>>.
- SHANKLIN, T. A.; LOULIER, B.; MATSON, E. T. Embedded sensors for indoor positioning. In: *Sensors Applications Symposium (SAS), 2011 IEEE*. [S.l.: s.n.], 2011. p. 149–154.
- SHEKIN, D. *Handbook of Parametric and Nonparametric Statistical Procedures, Fifth Edition*. Taylor & Francis, 2011. ISBN 9781439858011. Disponível em: <<https://books.google.com.br/books?id=YDd2cgAACAAJ>>.
- SHIBATA, F.; IKEDA, S.; KURATA, T.; UCHIYAMA, H. An intermediate report of trakmark wg international voluntary activities on establishing benchmark test schemes for ar/mr geometric registration and tracking methods. In: *Mixed and Augmented Reality (ISMAR), 2010 9th IEEE International Symposium on*. [S.l.: s.n.], 2010. p. 298–302.
- SHIN, B.; LEE, J. H.; LEE, H.; KIM, E.; KIM, J.; LEE, S.; CHO, Y. su; PARK, S.; LEE, T. Indoor 3d pedestrian tracking algorithm based on pdr using smarthphone. In: *Control, Automation and Systems (ICCAS), 2012 12th International Conference on*. [S.l.: s.n.], 2012. p. 1442–1445.
- SHIN, B.-J.; LEE, K.-W.; CHOI, S.-H.; KIM, J.-Y.; LEE, W. J.; KIM, H. S. Indoor wifi positioning system for android-based smartphone. In: *Information and Communication Technology Convergence (ICTC), 2010 International Conference on*. [S.l.: s.n.], 2010. p. 319–320.
- SHIN, H.; CHON, Y.; CHA, H. Unsupervised construction of an indoor floor plan using a smartphone. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, v. 42, n. 6, p. 889–898, Nov 2012. ISSN 1094-6977.
- SILVA, M.; ROBERTO, R.; TEICHRIEB, V. Evaluation of augmented reality technology in the english language field. In: *Informática na Educação (SBIE), Anais do XXVI Simpósio Brasileiro de*. [S.l.: s.n.], 2015. p. 577–586.
- SILVA, M. da; ROBERTO, R.; TEICHRIEB, V.; CAVALCANTE, P. Towards the development of guidelines for educational evaluation of augmented reality tools. In: *K-12 Embodied Learning through Virtual & Augmented Reality (KELVAR), IEEE Virtual Reality 2016 Workshop on*. [S.l.: s.n.], 2016.
- SILVA, V. E.; LINS, C.; SILVA, A.; ROBERTO, R.; ARAÚJO, C.; TEICHRIEB, V. Voxar puzzle: An innovative hardware/software computer vision game for children development.

- In: *Virtual and Augmented Reality (SVR), 2015 XVII Symposium on*. [S.l.: s.n.], 2015. p. 147–153.
- SIM, R.; ROY, N. Global a-optimal robot exploration in slam. In: *Robotics and Automation, 2005. ICRA 2005. Proceedings of the 2005 IEEE International Conference on*. [S.l.: s.n.], 2005. p. 661–666.
- SIMÕES, F. P. M. *Object Detection and Pose Estimation from Natural Features for Augmented Reality in Complex Scenes*. Tese (Doutorado) — Federal University of Pernambuco, March 2016.
- SINHA, S. N.; STEEDLY, D.; SZELISKI, R.; AGRAWALA, M.; POLLEFEYS, M. Interactive 3D architectural modeling from unordered photo collections. *ACM Trans. Graph.*, ACM, New York, NY, USA, v. 27, n. 5, p. 159:1–159:10, dez. 2008. ISSN 0730-0301. Disponível em: <<http://doi.acm.org/10.1145/1409060.1409112>>.
- SOLIN, A.; CORTES, S.; RAHTU, E.; KANNALA, J. Pivo: Probabilistic inertial-visual odometry for occlusion-robust navigation. In: *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. [S.l.: s.n.], 2018. p. 616–625.
- SONG, H.; LIU, H.; CHEN, D. An automatic gui adjustment method for mobile computing. In: *Computer Science and Automation Engineering (CSAE), 2011 IEEE International Conference on*. [S.l.: s.n.], 2011. v. 3, p. 206–210.
- STURM, J.; ENGELHARD, N.; ENDRES, F.; BURGARD, W.; CREMERS, D. A benchmark for the evaluation of rgb-d slam systems. In: *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. [S.l.: s.n.], 2012. p. 573–580. ISSN 2153-0858.
- TAKACS, G.; CHANDRASEKHAR, V.; TSAI, S.; CHEN, D.; GRZESZCZUK, R.; GIROD, B. Unified real-time tracking and recognition with rotation-invariant fast features. In: *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. [S.l.: s.n.], 2010. p. 934–941. ISSN 1063-6919.
- TAMURA, H.; KATO, H. Proposal of international voluntary activities on establishing benchmark test schemes for ar/mr geometric registration and tracking methods. In: *Mixed and Augmented Reality, 2009. ISMAR 2009. 8th IEEE International Symposium on*. [S.l.: s.n.], 2009. p. 233–236.
- TAN, D.; ILIC, S. Multi-forest tracker: A chameleon in tracking. In: *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*. [S.l.: s.n.], 2014. p. 1202–1209.
- TANSKANEN, P.; KOLEV, K.; MEIER, L.; CAMPOSECO, F.; SAURER, O.; POLLEFEYS, M. Live metric 3d reconstruction on mobile phones. In: *Computer Vision (ICCV), 2013 IEEE International Conference on*. [S.l.: s.n.], 2013. p. 65–72. ISSN 1550-5499.
- TATENO, K.; TOMBARI, F.; LAINA, I.; NAVAB, N. Cnn-slam: Real-time dense monocular slam with learned depth prediction. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. [S.l.: s.n.], 2017. p. 6565–6574. ISSN 1063-6919.
- TEIXEIRA, J. M. X. N. *Analysis and Evaluation of Optimization Techniques for Tracking in Augmented Reality Applications*. Tese (Doutorado) — Federal University of Pernambuco, March 2013.

- TERAURA, N.; SAKURAI, K. Preventing the access of fraudulent web sites by using a special two-dimensional code. In: *Innovative Mobile and Internet Services in Ubiquitous Computing (IMIS), 2012 Sixth International Conference on*. [S.l.: s.n.], 2012. p. 645–650.
- The Next Generation Mobile Networks Ltd. *NGMN - 5G Work Overview*. 2015. [Online; last access: 03-February-2016]. Disponível em: <<http://www.ngmn.org/work-programme/5g-work-overview.html>>.
- TOBIAS, L.; DUCOURNAU, A.; ROUSSEAU, F.; MERCIER, G.; FABLET, R. Convolutional neural networks for object recognition on mobile devices: A case study. In: *2016 23rd International Conference on Pattern Recognition (ICPR)*. [S.l.: s.n.], 2016. p. 3530–3535.
- TOMASI, C.; MANDUCHI, R. Bilateral filtering for gray and color images. In: *Computer Vision. Sixth International Conference on*. [S.l.: s.n.], 1998. p. 839–846.
- UCHIYAMA, H.; TAKETOMI, T.; IKEDA, S.; LIMA, J. P. S. d. M. [poster] abecedary tracking and mapping: A toolkit for tracking competitions. In: *ISMAR*. [S.l.: s.n.], 2015. p. 198–199.
- ULLAH, A.; ISLAM, M.; AKTAR, S.; HOSSAIN, S. Remote-touch: Augmented reality based marker tracking for smart home control. In: *Computer and Information Technology (ICCIT), 2012 15th International Conference on*. [S.l.: s.n.], 2012. p. 473–477.
- VALOEN, L.; SHOESMITH, M. I. The effect of phev and hev duty cycles on battery and battery pack performance. In: *Plugin Highway 2007 Conference*. [s.n.], 2007. p. 1–9. Disponível em: <http://umanitoba.ca/outreach/conferences/phev2007/PHEV2007/proceedings/PluginHwy_PHEV2007_PaperReviewed_Valoen.pdf>.
- VENTURA, J.; ARTH, C.; REITMAYR, G.; SCHMALSTIEG, D. Global localization from monocular slam on a mobile phone. *Visualization and Computer Graphics, IEEE Transactions on*, v. 20, n. 4, p. 531–539, April 2014. ISSN 1077-2626.
- VENTURA, J.; HOLLERER, T. Wide-area scene mapping for mobile visual tracking. In: *Mixed and Augmented Reality (ISMAR), 2012 IEEE International Symposium on*. [S.l.: s.n.], 2012. p. 3–12.
- VINEET, V.; MIKSIK, O.; LIDEGAARD, M.; NIEßNER, M.; GOLODETZ, S.; PRISACARIU, V. A.; KÄHLER, O.; MURRAY, D. W.; IZADI, S.; PÉREZ, P.; TORR, P. H. S. Incremental dense semantic stereo fusion for large-scale semantic scene reconstruction. In: *2015 IEEE International Conference on Robotics and Automation (ICRA)*. [S.l.: s.n.], 2015. p. 75–82. ISSN 1050-4729.
- WAGNER, D.; MULLONI, A.; LANGLLOTZ, T.; SCHMALSTIEG, D. Real-time panoramic mapping and tracking on mobile phones. In: *Virtual Reality Conference (VR), 2010 IEEE*. [S.l.: s.n.], 2010. p. 211–218. ISSN 1087-8270.
- WAGNER, D.; REITMAYR, G.; MULLONI, A.; DRUMMOND, T.; SCHMALSTIEG, D. Real-time detection and tracking for augmented reality on mobile phones. *Visualization and Computer Graphics, IEEE Transactions on*, v. 16, n. 3, p. 355–368, May 2010. ISSN 1077-2626.

- WAGNER, D.; SCHMALSTIEG, D.; BISCHOF, H. Multiple target detection and tracking with guaranteed framerates on mobile phones. In: *Mixed and Augmented Reality, 2009. ISMAR 2009. 8th IEEE International Symposium on*. [S.l.: s.n.], 2009. p. 57–64.
- WIERINGA, R.; MAIDEN, N.; MEAD, N.; ROLLAND, C. Requirements engineering paper classification and evaluation criteria: A proposal and a discussion. *Requir. Eng.*, Springer-Verlag New York, Inc., Secaucus, NJ, USA, v. 11, n. 1, p. 102–107, dez. 2005. ISSN 0947-3602. Disponível em: <<http://dx.doi.org/10.1007/s00766-005-0021-6>>.
- XIAO, J.; RUSSELL, B. C.; TORRALBA, A. Localizing 3d cuboids in single-view images. In: *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*. USA: Curran Associates Inc., 2012. (NIPS'12), p. 746–754. Disponível em: <<http://dl.acm.org/citation.cfm?id=2999134.2999218>>.
- XU, L.; LI, S.; BIAN, K.; ZHAO, T.; YAN, W. Sober-drive: A smartphone-assisted drowsy driving detection system. In: *Computing, Networking and Communications (ICNC), 2014 International Conference on*. [S.l.: s.n.], 2014. p. 398–402.
- YOVCHEVA, Z.; BUHALIS, D.; GATZIDIS, C. Smartphone augmented reality applications for tourism. *e-Review of Tourism Research (eRTR)*, v. 10, n. 2, p. 63–66, 2012. Disponível em: <<http://eprints.bournemouth.ac.uk/20219/>>.
- ZHANG, L.; LIU, J.; JIANG, H.; GUAN, Y. Senstrack: Energy-efficient location tracking with smartphone sensors. *Sensors Journal, IEEE*, v. 13, n. 10, p. 3775–3784, Oct 2013. ISSN 1530-437X.
- ZHOU, F.; DUH, H. B.-L.; BILLINGHURST, M. Trends in augmented reality tracking, interaction and display: A review of ten years of ismar. In: *Mixed and Augmented Reality, 2008. ISMAR 2008. 7th IEEE/ACM International Symposium on*. [S.l.: s.n.], 2008. p. 193–202.
- ZHU, L.; HYYPPÄ, J.; KUKKO, A.; KAARTINEN, H.; CHEN, R. Photorealistic building reconstruction from mobile laser scanning data. *Remote Sensing*, v. 3, n. 7, p. 1406–1426, 2011. ISSN 2072-4292. Disponível em: <<http://www.mdpi.com/2072-4292/3/7/1406>>.