

**Federal University of Bahia
State University of Feira de Santana**

MASTER THESIS

A multi-view environment for markerless augmented reality

Caio Sacramento de Britto Almeida

Master in Computer Science – MMCC

Salvador
December 15th, 2014

MMCC-Msc-0007

CAIO SACRAMENTO DE BRITTO ALMEIDA

**A MULTI-VIEW ENVIRONMENT FOR MARKERLESS
AUGMENTED REALITY**

Thesis submitted to the Master Program in Computer Science from Federal University of Bahia and State University of Feira de Santana, in partial fulfillment of the requirements for the degree of Master in Computer Science.

Advisor: Antônio Lopes Apolinário Júnior

Salvador
December 15th, 2014

Index card.

Sacramento de Britto Almeida, Caio

A multi-view environment for markerless augmented reality/ Caio Sacramento de Britto Almeida– Salvador, December 15th, 2014.

27p.: il.

Advisor: Antônio Lopes Apolinário Júnior.

Thesis (master)– Federal University of Bahia, Institute of Mathematics, December 15th, 2014.

1. Multi-view environment. 2. Augmented reality. 3. Computer graphics.

I. Apolinário Jr., Antônio Lopes. II. Federal University of Bahia. Institute of Mathematics. III A multi-view environment for markerless augmented reality.

CDD 20.ed. 123.45

APPROVAL SHEET**CAIO SACRAMENTO DE BRITTO ALMEIDA****A MULTI-VIEW ENVIRONMENT FOR
MARKERLESS AUGMENTED REALITY**

This thesis was considered worthy of acceptance of the master's degree in Computer Science and approved on its final form by the UFBA-UEFS Master Program in Computer Science.

Salvador, December 15th, 2014

Prof. Dr. Antônio Lopes Apolinário Júnior 1
Federal University of Bahia - Brazil

Prof. Dr. Michelle Ângelo 2
State University of Feira de Santana - Brazil

Prof. Dr. Rodrigo Silva 3
Federal University of Juiz de Fora - Brazil

ACKNOWLEDGEMENTS

I would like to thank my family for all the support during all those years, not only the time spent on the master program, but also during all my graduation. Namely, first I would like to thank my brother, Rodrigo, for being my inspiration on following this career; my twin sister, Thalita, for being my support of all times; my parents, for the great education that was given to me; my grandmother Lourdes; my girlfriend Jéssica; my friends; and specially my advisor Antônio Apolinário for relying on me to make this work happen.

ABSTRACT

Augmented reality is a technology which allows 2D and 3D computer graphics to be aligned or registered with scenes of the real-world in real-time. This projection of virtual images requires a reference in the captured real image, which is often achieved by using one or more markers. But, there are situations where using markers can be unsuitable, like medical applications, for example. In this work, we present a multi-view environment, composed by augmented reality glasses and two Kinect devices, which doesn't use fiducial markers in order to run augmented reality applications. All devices are calibrated according to a common reference system, and then the virtual models are transformed accordingly too. In order to achieve that, two approaches were specified and implemented: one based on one Kinect plus optical flow and accelerometer data from augmented reality glasses, and another one based purely on two Kinect devices. The results regarding quality and performance achieved by these two approaches are presented and discussed, as well as a comparison between them.

Keywords: augmented reality, augmented reality glasses, kinect, transformation, optical flow, markerless

CONTENTS

Chapter 1—Introduction	1
1.1 Motivation	1
1.2 Objectives	2
1.2.1 Features	3
1.3 Thesis overview	4
Chapter 2—Conceptual primer	5
2.1 Augmented reality	5
2.1.1 Direct or indirect vision	6
2.1.2 Markers	6
2.1.2.1 Fiducial markers	6
2.1.2.2 Markerless	8
2.2 Cameras	8
2.2.1 Calibration	8
2.2.2 Parameters	9
2.2.2.1 Intrinsic parameters	9
2.2.2.2 Extrinsic parameters	9
2.2.3 Calibration approaches	11
2.2.4 Multi-view	12
2.3 Sensor-based registration approach	13
2.3.1 Kinect	13
2.3.2 Registration	13
2.3.3 Transformation	13
2.4 Vision-based registration approach	13
2.4.1 Optical flow	13
2.4.2 Lucas-Kanade algorithm	14
Chapter 3—Related works	17
3.1 Markerless augmented reality	17
3.2 Multiple Kinects	17
3.3 Multi-view environment	17

Chapter 4—Solution architecture	19
4.1 Environment	19
4.2 Calibration	19
4.2.1 Augmented reality glasses calibration	19
4.2.2 Initial calibration between Kinects	19
4.2.3 Initial calibration between Kinect and glasses	19
4.3 Communication	19
4.4 Transformations	19
4.5 Method 1: Glasses accelerometer and one Kinect	19
4.6 Method 2: Two Kinects	19
4.7 Hybrid approach	19
Chapter 5—Result	21
5.1 Scope	21
5.2 Analysis	21
5.3 Comparison	21
Chapter 6—Conclusions	23
6.1 Future work	23

LIST OF FIGURES

1.1	Fiducial marker used to represent a tridimensional model over it (??)	2
1.2	Global vision	3
2.1	The augmented reality is localized between the extremes of the reality-virtuality continuum	5
2.2	Immersive augmented reality	6
2.3	Non-immersive augmented reality	7
2.4	Example of fiducial marker used by medical applications	7
2.5	Example of fiducial marker for motion capture	8
2.6	Images of a chessboard being held at various orientations (left) provide enough information to completely solve for the locations of those images in global coordinates (relative to the camera) and the camera intrinsics (Bradski; Kaehler, 2008)	12
2.7	Pyramid Lucas-Kanade optical flow: running optical flow at the top of the pyramid first mitigates the problems caused by violating our assumptions of small and coherent motion; the motion estimate from the preceding level is taken as the starting point for estimating motion at the next layer down (Bradski; Kaehler, 2008)	16
2.8	The same input image being used for features identification by Shi-Tomasi algorithm (red dots) and by Harris Corner Detector (green dots)	16

LIST OF TABLES

Chapter

1

In this chapter I present the motivation, objectives and overview of this work.

INTRODUCTION

1.1 MOTIVATION

Augmented reality has taken advantage from the progresses on the fields of multimedia and virtual reality, making feasible new forms of interaction between humans and machines. Differently from virtual reality, that takes the user to a virtual environment, the augmented reality keeps the user on his physical environment and takes the virtual environment to the user's space, allowing the interaction with the virtual world, in a more natural manner and without needing training or adaptation (??). Many times, this interaction means merging virtual images with images captured from a real environment.

One of the greatest challenges on the field of augmented reality is to determine, in real time, which virtual image to be displayed, in which position and how it should be represented. In order to obtain an integration illusion between real objects and virtual objects, the generated object should be aligned with the tridimensional position and orientation of the real objects (??). This can be achieved by estimating the camera position.

On many situations, fiducial markers¹ are used (often this is due to the real time requirements of the augmented reality applications) (??) and are drawn in a way that they can be easily identified. Those markers need to be placed on the target scene and can achieve great results using just a few computational resources. Figure 1.1 shows the usage of a fiducial marker and a tridimensional object being projected over it.

However, besides requiring human interference on the scene, there situations where the usage of fiducial markers is not possible, feasible or comfortable for the target model. That is the case, for example, of medical applications on which this model is a patient. It's also possible to cite other limitations of fiducial markers, like, for example, occlusion (a virtual image could be not projected if the marker is not completely visible) and illumination (the intensity of light reflected by the marker could make it hard to be

¹A fiducial marker or fiducial is an object placed in the field of view of an imaging system which appears in the image produced, for use as a point of reference or a measure.

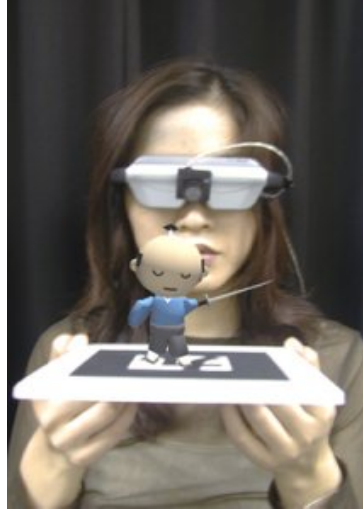


Figure 1.1 Fiducial marker used to represent a tridimensional model over it (??)

identified). Less common, there are approaches that replace fiducial markers (??) (??) by GPS, gyroscopes, accelerometers, cameras, among others (??) (??). These approaches have the advantage of not requiring human interference on the scene (to put a marker or to move it around).

Depending on the way that a user sees the mixed world, augmented reality can be classified on two ways. When the user sees the mixed world pointing his eyes straight to the real positions with optical scene or video, this augmented reality is called *immersive* or of *direct vision*. On the other hand, when the user sees the mixed world by some device, like a monitor screen or projector, not aligned with the real positions, this augmented reality is *non-immersive* or of *indirect vision* (??).

This work proposes a multi-view environment for augmented reality, of direct vision, composed by two Kinects (??) and augmented reality glasses, that allows a watcher visualize, in real time, virtual images merged with real images from the target model. In this approach, it's not intended to use any fiducial marker. Instead, it will be used a geometric approach based on the data captured by each Kinect. This proposed environment can be used, for example, on the medical field (real situations, education and training) or in other situation where a multi-view environment for markerless augmented reality is applicable.

1.2 OBJECTIVES

The environment proposed on this work aims to contribute to augmented reality applications where virtual images need to be merged with real images in real time, without using fiducial markers, and considering the viewing angle of the observer and the position of the target object.

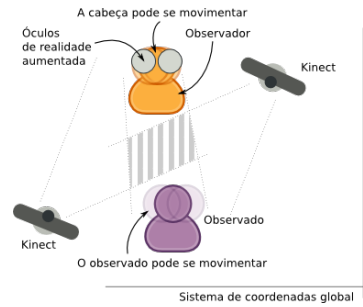


Figure 1.2 Global vision

1.2.1 Features

Based on the study of related works, it was defined the following features that define the scope of the first version of the augmented reality environment, graphically represented on Figure (??):

- There are two main elements on the environment, the observer and the target model;
- Observer and target model are positioned one in front of the other, with some restriction of minimum and maximum distance;
- The observer sees the combination of a real image with a pre-defined virtual image over the target model, through the augmented reality glasses;
- The observer and target model can move their heads, in a way that the virtual image adapts itself in real time to fit the new viewing angles;
- No fiducial marker is used.

The elements presented on this environment are the following:

- Observer - He is the user of the system that wears a pair of augmented reality glasses and is positioned in front of the target model (which can be a human being, for example, or a static object). On a medical context, the observer would be a medical specialist, responsible for observing the patient and responsible for analyzing the combination of the real image (part of the patient body) with a virtual image (from magnetic resonance imaging or computed tomography). The observer is able to move his head.
- Augmented reality glasses - The augmented reality glasses are wore by the observer and have two cameras. The images captured from the target model will be merged with virtual images and displayed on the two lenses of these glasses. On one of the approaches implemented by this work, the observer movimentation is determined from sensors present on the glasses: acceletometer, magnetometer and gyroscopes. Based on sensor data, it's possible to determine the variation on the orientation of the glasses, and so it's possible to know how much the observer has moved his

head. This calculation returns values that define movements on longitudinal axis, vertical axis and lateral axis. The virtual image should be reprojected in real time according to those movements.

- **Target model** - The target model (a human being, for example), is placed in front of the observer and doesn't use any kind of fiducial marker. The main goal is that the virtual image is placed over the target model. In order to calculate where and how the virtual image should be displayed, it's necessary to identify the position and orientation of the real object relative to the observer. This is done based on two sensors placed on the environment, where the first one captures data from the observer and the second one captures data from the target model.
- **Sensors** - Two sensors are placed on the environment and capture data from the observer and the target model (one for each). The information captured by the target model's sensor contains its model, that will be merged with the virtual image.

Each device presented on this multi-view (glasses and sensors) environment has its own coordinate system, but all information must be converted to a global coordinate system.

Since there are two sensors and one pair of augmented reality glasses, another objective is to implement two different approaches: one that uses the augmented reality glasses to determine the observer's pose (based on data from accelerometer and magnetometer) and another one that uses a second Kinect device to determine the observer's pose based on a reconstruction of his model.

1.3 THESIS OVERVIEW

The next chapter "Conceptual primer" describes some basic theory needed. It also describes the hardware that are used by this work and how to calibrate them. After that, the "Related work" chapter presents some works related to this one, divided by subject. The third chapter, "Solution architecture", after the theory and related works were presented, explains the steps performed in order to implement the objectives of this work. The results of this implementation are presented on the chapter later, "Results", and finally the conclusions about those results are presented on the last chapter, "Conclusions", where possible future works are also listed.

In this chapter I present the main concepts behind this work.

CONCEPTUAL PRIMER

2.1 AUGMENTED REALITY

Augmented reality was born on the decade of 1990, to merge a virtual image or virtual environment with a real image or real environment. But only from the 2000s it became more popular, due to lower costs of hardware and software devices, and ready to be used on tangible and multimodal (voice, touch, gesture, etc.) applications (??). It can be considered the mixing of real and virtual worlds at some point of the continuum reality-virtuality, that connects completely virtual environments to completely real environments (??), like shown on Figure 2.1. It can also be considered a system that completes the real world with virtual objects, in a way that they seem to exist on the same space, respecting the following features:

- Real objects are mixed with virtual objects;
- Execution is interactive and on real time;
- Virtual and real objects are aligned;
- Applicable to all human sensory systems, including auditory, olfactory and somatosensory (??).

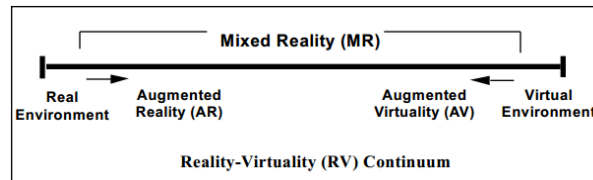


Figure 2.1 The augmented reality is localized between the extremes of the reality-virtuality continuum

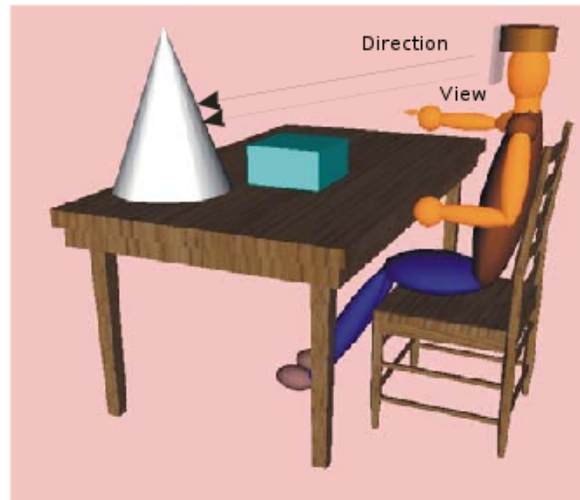


Figure 2.2 Immersive augmented reality

From a human-computer interaction perspective, the augmented reality can be considered a new way of interaction between humans and computers, and, on this aspect, it can be classified as direct vision or indirect vision.

2.1.1 Direct or indirect vision

Augmented reality is classified as *direct vision* or *immersive* when the user sees the mixed world pointing his eyes straight to the real position of the objects of interest, like shown on Figure 2.1.1.

On the direct vision, images from the real world can be seen with the naked eye or brought by video, while the generated virtual images can be projected on the eyes, on the real scenario or mixed with the real world video. Immersive augmented reality can be implemented with optical helmets, for example (??).

On the other hand, non-immersive augmented reality (indirect vision), happens when the user sees the mixed world through other output device, like a monitor screen, for example, not aligned with the real positions, as shown on Figure 2.1.1. On this kind of vision, real and virtual images are merged and displayed as video to the user. It can be achieved by using cameras and projectors (??).

2.1.2 Markers

In order to identify the position where a virtual image should be rendered, two main approaches can be applied by augmented reality applications: one is to use fiducial markers, the other is not use them.

2.1.2.1 Fiducial markers Fiducial markers are often implemented as square white cards with a black symbol printed (or drawn) on it, easy to be recognized, working like a barcode or QR code. Computer vision techniques are used to calculate the position of

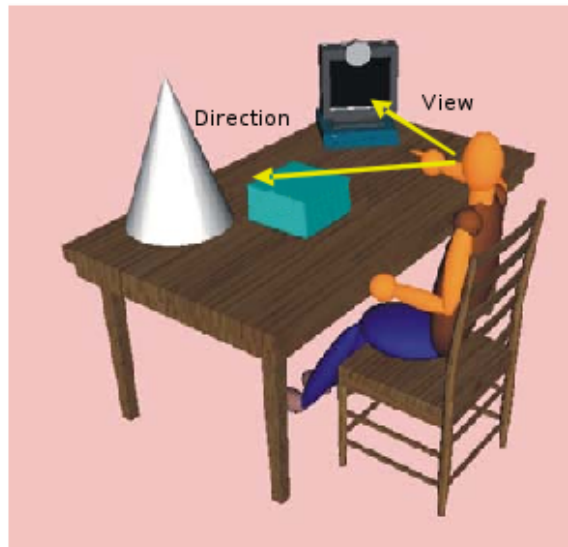


Figure 2.3 Non-immersive augmented reality

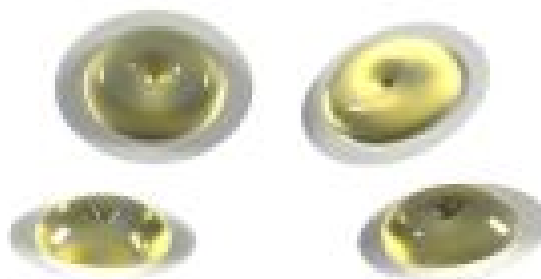


Figure 2.4 Example of fiducial marker used by medical applications

the real camera and its orientation relative to the markers, in a way that virtual objects can be projected over them.

Fiducial markers can assume other shapes besides a square card. The Figure 2.1.2.1 shows an example of a fiducial marker used on medical applications, similar to a sticker, which is fixed on the patient's skin. On these applications, augmented reality can be used for visualization and training on surgeries. It's possible to collect patient's data in real time, by using non-invasive sensors as the ones used for magnetic resonance imaging and computed tomography. This dataset can be merged in real time with the real image of the patient (??).

On motion capture systems, other kinds of markers can also be used, like, for example, the one described on (??) and represented on Figure 2.1.2.1. In order to estimate the pose, it uses a single camera and three spherical fiducial markers, which are also reflexive. Two of those markers are placed at the left and right sides of the actor, perpendicular to his ears. The third marker is placed on the same height as the other two, but in front of



Figure 2.5 Example of fiducial marker for motion capture

the actor. The markers have a color close to orange for high reflectance (in order to be easily recognized) and are supported by a structure made of carbon fiber (which is light).

Disadvantages of fiducial markers are listed on (??), for example, they are invasive (need to be placed on the scene or fixed on the target object), they have limited interaction (can't be moved around so much because it needs to still be visible and recognizable) and need to be printed or drawn before using and stored for future usage.

2.1.2.2 Markerless Markerless augmented reality means that fiducial markers are not used on the scene. Instead, features from the target scene need to be used in order to estimate the camera pose and objects orientations.

Target detection based on computer vision has been extensively studied e successfully applied on augmented reality applications. On related computer vision literature, geometric primitives can be used for pose estimation, on most cases, points, segments, lines, edges of edge points, cones, cylinders, or a combination of two or more of those features.

Markerless augmented reality has the advantage of using parts of the real environment as targets and can even extract information from the environment to be used by the augmented reality system (??).

Avoiding markers leads to a much more effective augmented reality experience, but requires the implementation of several image processing or sensor function techniques, resulting in more complex algorithms and in higher computational resources (Shumaker, 2011).

2.2 CAMERAS

Many tasks on augmented reality deal with an imaging device. Usually this imaging device is a camera, which performs a mapping from a 3D world to a 2D image (Hanning, 2011). The problem called *camera calibration* is the one that handles the determination of the parameters for this mapping.

2.2.1 Calibration

Camera calibration is a necessary step in 3D computer vision in order to extract metric information from 2D images. Much work has been done, by both the photogrammetry community and the computer vision community. Those techniques can be divided roughly into two categories: photogrammetric calibration (performed by observing a calibration object whose geometry in a 3-D space is known with very good precision) and self-calibration (does not use a calibration object, just move the camera in a scene and get its parameters) (Zhang, 2000).

Camera calibration means to determine the camera model parameters which fit best to the observed behavior of the actual camera (Hanning, 2011). So, it's necessary to measure the the distance of an observation to a given camera model. The determination of the optimal camera mapping concerning each of these distance functions defines a non-linear optimization problem, and as such, it depends on the initial value.

Ordinary cameras are very often modelled as pinhole cameras. A pinhole is an imaginary wall with a tiny hole in the center that blocks all rays except the ones that pass through this tiny center hole (Bradski; Kaehler, 2008). Unfortunately, a real pinhole is not a good way to make images because it does not gather enough light for rapid exposure. This is why human eyes and cameras use lenses to gather more light than what would be available at a single point. This way, the simple geometry of the pinhole camera model is not enough and it also introduces distortion of the lens itself.

The camera coordinate system (CCS) is a Cartesian coordinate system defined by the principal plane: The x-axis and y-axis of the CCS determine the principal plane, the z-axis is given by the optical axis. The optical center of the lens determines the origin (0, 0, 0) of the CCS. Thus, in the camera coordinate system, the principal plane becomes $z = 0$ (Hanning, 2011).

Camera calibration is usually split up into two distinct parts, the intrinsic and extrinsic parameters, which will be covered on the following section.

2.2.2 Parameters

The process of camera calibration gives us both a model of the camera's geometry and a distortion model of the lens. Both compose the intrinsic parameters of the camera.

2.2.2.1 Intrinsic parameters The intrinsic parameters are those that describe the internal workings of the camera and consist of the focal length, the principal point, the skew coefficients and radial and tangential distortions (Tillapaugh; Engineering, 2008). These values are needed to help to describe imperfections in the lens of the camera and give a mapping from camera reference frame to the image plane. The intrinsic parameters depend only on the camera itself, and so they just need to be found once, regardless the environment changes or not.

2.2.2.2 Extrinsic parameters The extrinsic parameters, on the other hand, represent the viewpoint of the camera by a rigid transformation, which describes its position

and orientation (Bajramovic, 2010). This part of the model is independent of the camera itself. The according parameters are called extrinsic, as they describe the relation between the camera and the world.

According to (Tillapaugh; Engineering, 2008), the extrinsic parameters are those that are dependent on the environment. To relate an object's coordinate system to the world's coordinate system, a translation and a rotation matrix are needed. Therefore, the extrinsic parameters consist of these two matrices, so that a mapping from the world coordinate system to the camera reference frame can be found. Since the extrinsic parameters for the camera explain how the camera relates to the environment, if the camera changes position, the parameters have to be recalculated (differently from the intrinsic parameters as explained on the previous section).

The relation between world coordinate system and the camera reference frame can be described by the equation:

$$X_c = R_c \times X + T_c$$

Where X is a 3x1 vector that represents a point on the world coordinate space, R_c is the extrinsic rotation matrix, T_c is the extrinsic translation matrix and X_c is a 3x1 vector in the camera reference frame. The rotation matrix is build from the combination of three single-axis rotation matrices:

$$R_x(\theta) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos(\theta) & \sin(\theta) \\ 0 & -\sin(\theta) & \cos(\theta) \end{bmatrix}$$

$$R_y(\phi) = \begin{bmatrix} \cos(\phi) & 0 & -\sin(\phi) \\ 0 & 1 & 0 \\ \sin(\phi) & 0 & \cos(\phi) \end{bmatrix}$$

$$R_z(\omega) = \begin{bmatrix} \cos(\omega) & \sin(\omega) & 0 \\ -\sin(\omega) & \cos(\omega) & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

The rotation matrix R has the property that its inverse is its transpose, hence $R^T R = R R^T = I$, where I is the identity matrix consisting of 1s along the diagonal and 0s everywhere else.

The translation vector T is how it's possible to represent a shift from one coordinate system to another system whose origin is displaced to another location; in other words, the translation vector is just the offset from the origin of the first coordinate system to the origin of the second coordinate system (Bradski; Kaehler, 2008). So, in order to shift from a coordinate system centered on an object to one centered at the camera, the appropriate translation vector is simply $T = origin_{object} - origin_{camera}$. Thus, a point in

the object (or world) coordinate frame P_o has coordinates P_c in the camera coordinate frame: $P_c = R(P_o - T)$.

The extrinsic parameters can be represented by a final matrix called *homogenous matrix*. It looks as follows:

$$\begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix}$$

So, in order to convert from a world origin coordinate to a camera orientation coordinate system, only a single homogenous matrix needs to be used.

2.2.3 Calibration approaches

There are multiple ways to calibrate cameras. The most common way to perform the calibration is to have multiple points that have known relationship to each other in the world's coordinate system captured in an image from the camera (Tillapaugh; Engineering, 2008). As explained previously, for the intrinsic parameters it does not matter how the points are related to the world coordinates, it's only important how they relate to each other. Such parameters consist of the horizontal and vertical focal lengths and the principal point of the camera and the skew. The skew is generally assumed to be zero (Furht, 2011) for many cameras.

Various techniques were created to obtain accurate mappings, including vanishing points for orthogonal directions and calibration from rotation purely (Medioni; Kang, 2004), features extraction which can be lines or points, single-image or multiple-images configurations, different camera models, linear and non-linear algorithms, among many others (Armangué; Salvi; Batlle, 2000) (Clarke; Fryer, 1998). Self-calibration is indeed feasible with simple image correspondences among frames, yet the participation of control points produces more robust calibration results, in closer agreement with object space constraints (Douskos; Kalisperakis; Karras, 2007).

The works by (Tsai, 1986) and (Zhang, 2000) are extensively referenced and propose closed form solutions for the estimation of intrinsic and extrinsic parameters using 3D and 2D calibration patterns respectively. Camera calibration is a much discussed topic but the lack of robust algorithms for features detection difficulties the construction of automatic calibration process (Laureano; Paiva; Silva, 2013). Calibration pattern recognition is a hard task, where the lighting problems and high level of ambiguities are the principal challenges. For this reason, the algorithms often require user intervention for a reliable detection of the calibration points.

There are some tools available for automatic camera calibration. Two of the most popular ones are The Bouguet MatLab Toolbox (Bouguet, 2008) and the OpenCV library (Bradski, 2000). The former is an application that asks the user to define four extreme points that represent the area where an algorithm searches for the calibration points, given the number of rows and columns of the pattern. The latter is a very popular computer vision library that offers an automatic way to detect chessboard patterns in images by the *findChessboardCorners()* function. The method performs successive morphological operators until a number of black and white contours are identified, subsequently the

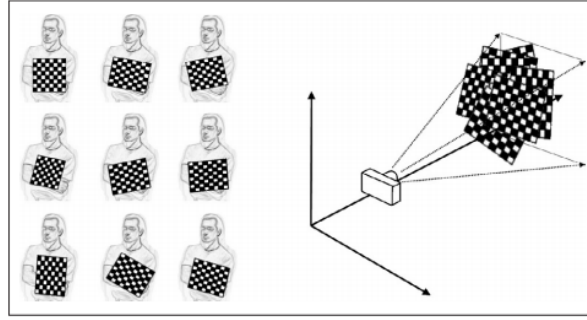


Figure 2.6 Images of a chessboard being held at various orientations (left) provide enough information to completely solve for the locations of those images in global coordinates (relative to the camera) and the camera intrinsics (Bradski; Kaehler, 2008)

corners of the contours make up the calibration point set. The pattern is recognized only if all rectangles are identified. Figure 2.2.3 shows how a person can use a chessboard in order to calibrate a camera using the OpenCV library.

As stated by (Arca; Casiraghi; Lombardi, 2005), the well-known algorithms for chessboard detection proposed by (Tsai, 1986) and (Zhang, 2000) can achieve good calibration results since both of them are based on an initialization procedure that requires the precise corners positions of a calibration pattern (a chessboard, for example). Although the pattern to be detected is generally a simple object, the detection of its corners at subpixel precision is a very difficult task to be solved under uncontrolled acquisition conditions. At the state of art, most of the presented methods determine the positions of the chessboard corners by means of a two step processing scheme. At first the corners are detected with pixel precision, by means of edge detection methods; to increase the detection accuracy, the second step modifies their positions by the interpolation of the found borders. Those methods rarely reach a high accuracy, since the interpolation provides just a guess about the precise corners positions. Some authors try to improve the accuracy in the later steps of camera calibration, which increases computational cost, while others, like (Arca; Casiraghi; Lombardi, 2005), does not make any assumption on the orientation and scale of the chessboard and try to first find the chessboard, then determine the size of the squares and later on find the corners by using a simple statistical model.

2.2.4 Multi-view

Calibrating a multi-camera system accordingly means estimating the intrinsic and extrinsic parameters of all cameras (Bajramovic, 2010). Actually it is mainly concerned with estimating extrinsic parameters, since most approaches first compute intrinsic parameters individually for each camera. Extrinsic parameters are subsequently estimated given the intrinsic parameters.

The calibration task is usually formulated such that the world coordinate system may be chosen arbitrarily, which amounts to calibrate up to a rigid transformation of the whole system. Furthermore, if no metric measurements are used, the scale of the multi-camera system cannot be determined. This is the case if only correspondences between

the cameras are used as input (of the extrinsic calibration) (Bajramovic, 2010).

According to (Hanning, 2011), the calibration of a stereo camera system is more than calibrating two camera separately: An additional constraint for the stereo camera system can be applied, since a calibration target is observed by both cameras.

So, in theory, calibration of a multi-camera system is possible by individually calibrating intrinsic and extrinsic parameters of all cameras involved with respect to the same world coordinate system using classical calibration methods. This just requires that all cameras observe a common calibration object or that a calibration object is moved to precisely known positions such that every camera can observe it. The first approach imposes very strict limitations on the allowable relative camera poses (relative orientation and position) or requires a very large calibration object, which is generally complicated (Chen; Davis, 2000).

Once cameras are calibrated, they are ready to register objects from a mixed reality scene. In general, traditional methods can be classified into three types: sensor-based methods, image-based methods and hybrid methods (Shumaker, 2011). The first two will be covered in the next sections; the third, which is also applied by this work, won't be covered here since it's basically a combination of the first two methods.

2.3 SENSOR-BASED REGISTRATION APPROACH

Three-dimensional (3D) reconstruction of an environment is an important problem that has received much attention in past decades. It can be defined as the process of capturing the shape and appearance of real objects (Pati, 2012). Recovering 3D surface structure of an object has been a central issue in computer vision. On the last years, with the launch of inexpensive RGBD sensors (notably the Kinect), applications based on sensor registration became cheaper and more feasible.

2.3.1 Kinect

2.3.2 Registration

2.3.3 Transformation

2.4 VISION-BASED REGISTRATION APPROACH

Registration based on sensors relies solely upon geometric information, for example, the spatial coordinates of the points in two clouds. As a result, environments that don't provide high geometric texture (for example, plane surfaces like a wall) can cause sensor tracking to fail. For such cases, a vision-based registration approach can be applied, by using the RGB information from the images (Peasley; Birchfield, 2013).

2.4.1 Optical flow

Optical flow is a pattern that describes the apparent motion of objects, edges or surfaces on a visual scene caused by the relative motion between that object and an observer (a camera or the human eye) (Burton; Radford, 1978). This concept was first introduced by

the psychology in order to describe the visual stimulus provided to animals moving on the world (Gibson, 1950).

A fundamental problem in the processing of image sequences is the measure of the optical flow (or image velocity). The main objective is to compute an approximation to the 2D motion field - a projection of the 3D velocities of surface points onto the imaging surface - from spatiotemporal patterns of image density. Once computed, the measurements of image velocity can be used for many purposes, including object tracking (Barron; Fleet; Beauchemin, 1994).

The field of optical flow estimation is making big progresses as evidenced by the increasing accuracy of current methods on the Middlebury optical flow benchmark (Scharstein; Szeliski, 2002). After almost 30 years of research, these methods have obtained an impressive level of reliability and accuracy, by combining a data term that assumes constancy of some image property with a spatial term that models how the flow is expected to vary across the image (Sun; Roth; Black, 2014).

Differential methods belong to the most widely used techniques for optic flow computation in image sequences. They can be classified into local methods such as Lucas-Kanade technique or Bigun's structure tensor method, and into global methods such as the Horn/Schunck approach and its extensions. Often local methods are more robust under noise, while global techniques yield dense flow fields (Bruhn; Weickert; Schnörr, 2005).

Despite their differences, many of these techniques can be viewed conceptually in terms of three stages of processing (Barron; Fleet; Beauchemin, 1994):

- Prefiltering or smoothing with low-pass/band-pass filter in order to extract signal structure of interest and to enhance the signal-to-noise ratio;
- The extraction of basic measurements, such as spatiotemporal derivatives (to measure normal components of velocity) or local correlation surfaces;
- The integration of these measurements to produce a 2D flow field, which often involves assumptions about the smoothness of the underlying flow field.

One of the most popular methods for computing the optical flow of an image is the Lucas-Kanade algorithm, which will be explained on the following section.

2.4.2 Lucas-Kanade algorithm

The Lucas-Kanade (Lucas; Kanade, 1981) algorithm is a differential method for computing the optical flow of an image (Peasley; Birchfield, 2013). The goal of Lucas-Kanade is to align a template image $T(x)$ to an input image $I(x)$ where $x = (x, y)^T$ is a column vector containing the pixel coordinates. If the Lucas-Kanade algorithm is being used to compute optical flow or to track an image from time $t = 1$ to time $t = 2$, the template $T(x)$ is an extracted sub-region (a window) of the image at $t = 1$ and $I(x)$ is the image at $t = 2$. Applications of the Lucas-Kanade algorithm range from optical flow and tracking to layered motion, mosaic construction and face coding. Numerous variations have been made to the original algorithm (Baker; Matthews, 2004).

The basic idea of the Lucas-Kanade algorithm is based on three assumptions (Bradski; Kaehler, 2008):

- Brightness constancy: a pixel from the image of an object in the scene does not change in appearance as it (possibly) moves from frame to frame. For grayscale images (although Lucas-Kanade can also be done in color) this means that it's assumed that the brightness of a pixel does not change as it is tracked from frame to frame;
- Temporal persistence: the image motion of a surface patch changes slowly in time. In practice, this means the temporal increments are fast enough relative to the scale of motion in the image that the object does not move much from frame to frame;
- Spatial coherence: neighboring points in a scene belong to the same surface, have similar motion, and project to nearby points on the image plane.

Mathematically, the goal of the Lucas-Kanade is to minimize the sum of squared error between the two images, the template T and the image I warped back onto the coordinate frame of the template (Baker; Matthews, 2004):

$$\sum_x [I(W(x; p)) - T(x)]^2$$

Where $W(x; p)$ denote the parametrized set of allowed warps, where $p = (p_1, \dots, p_n)^T$ is a vector of parameters. The warp $W(w; p)$ takes the pixel x in the coordinate frame of the template T and maps it to the sub-pixel location $W(x; p)$ in the coordinate frame of the image I . In this context of optical flow computation, the warps $W(w; p)$ are the translations:

$$W(w; p) = \begin{pmatrix} x + p_1 \\ y + p_2 \end{pmatrix}$$

Where the vector of parameters $p = (p_1, p_2)^T$ is then the optical flow.

Warping I back to compute $I(W(x; p))$ requires interpolating the image I at the sub-pixel locations $W(x; p)$. The minimization is performed with respect to p and the sum is performed over all of the pixels x in the template image $T(x)$. Minimizing the expression is a non-linear optimization task even if $W(w; p)$ is linear in p because the pixel values $I(x)$ are, in general, non-linear in x . In fact, the pixel values $I(x)$ are essentially un-related to the pixel coordinates x . To optimize the expression, the Lucas-Kanade algorithm assumes that a current estimate of p is known and then iteratively solves for increments to the parameters Δp , so the expression is (approximately) minimized:

$$\sum_x [I(W(x; p + \Delta p)) - T(x)]^2$$

With respect to Δp , and then the parameters are updated:

$$p \leftarrow p + \Delta p$$

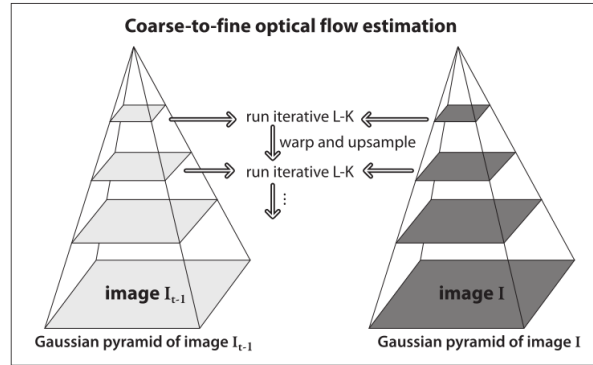


Figure 2.7 Pyramid Lucas-Kanade optical flow: running optical flow at the top of the pyramid first mitigates the problems caused by violating our assumptions of small and coherent motion; the motion estimate from the preceding level is taken as the starting point for estimating motion at the next layer down (Bradski; Kaehler, 2008)

These steps are iterated until the estimates for the parameters p converge.

The result of the algorithm is a set of optical flow vectors distributed over the image which give an estimation idea of the movement of objects in the scene, although some of those vectors will be erroneous (Rojas, 2014).

The Lucas-Kanade method can be applied in a sparse context because it relies only on local information that is derived from some small window surrounding each of the points of interest. The disadvantage of using small local windows in Lucas-Kanade is that large motions can move points outside of the local window and thus become impossible for the algorithm to find, which led to the development of the "pyramidal" Lucas-Kanade algorithm (Bouguet, 2000), one of many approaches that have been used to improve the convergence rate and reduce the likelihood of falling into a local minimum.

One component in many algorithms is a coarse-to-fine strategy. The most common approach is to build image pyramids by repeated blurring and downsampling. So, optical flow is first computed on the top level (fewest pixels) and then upsampled and used to initialize the estimate at the next level (Baker et al., 2011), as shown on Figure 2.4.2. Computation at the higher levels in the pyramid involves far fewer unknowns and so is far faster. The initialization at each level from the previous level also means that far fewer iterations are required at each level. Due to this, pyramid algorithms are usually faster than a single solution at the bottom level.

The features to be tracked by the Lucas-Kanade algorithm can be manually defined or identified from an algorithm, like the popular one proposed by Shi and Tomasi (Shi; Tomasi, 1994), called "good features to track", and implemented by the OpenCV library (Bradski, 2000). The idea of this algorithm is a modification of the Harris Corner Detection (Harris; Stephens, 1988), which basically replaces Harris' scoring function. Figure 2.4.2 shows the feature points on the same input image as identified by both algorithms.

So, the Lucas-Kanade algorithm, applied over a good set of feature points, is an efficient method for obtaining optical flow information at interesting points in an image and works for moderate object speeds.

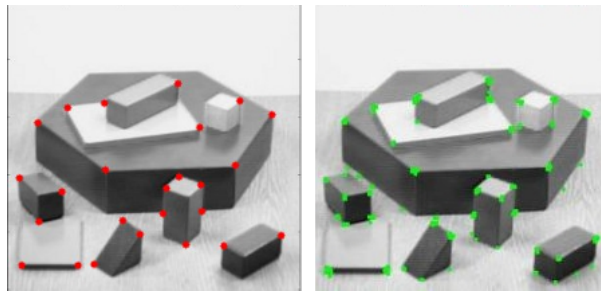


Figure 2.8 The same input image being used for features identification by Shi-Tomasi algorithm (red dots) and by Harris Corner Detector (green dots)

Chapter

3

In this chapter I present some related works.

RELATED WORKS

Cover all the related works, with multiple Kinects, optical flow, markerless augmented reality, medical applications, multi-view environment, reconstruction, etc.

3.1 MARKERLESS AUGMENTED REALITY

3.2 MULTIPLE KINECTS

3.3 MULTI-VIEW ENVIRONMENT

In this chapter I explain in details the steps performed in order to implement the objective of this work.

SOLUTION ARCHITECTURE

4.1 ENVIRONMENT

Technologies, machines, SOs, etc.

4.2 CALIBRATION

4.2.1 Augmented reality glasses calibration

4.2.2 Initial calibration between Kinects

4.2.3 Initial calibration between Kinect and glasses

4.3 COMMUNICATION

Network, sockets, etc.

4.4 TRANSFORMATIONS

4.5 METHOD 1: GLASSES ACCELEROMETER AND ONE KINECT

Cover Lucas-Kanade algorithm, etc.

4.6 METHOD 2: TWO KINECTS

Talk about performance of two Kinects fighting for a single GPU.

4.7 HYBRID APPROACH

When optical flow has just a few feature points, we switch to the second Kinect.

Chapter

5

In this chapter I present the results of the procedure explained in the previous chapter.

RESULT

5.1 SCOPE

Talk about error propagation.

5.2 ANALYSIS

Talk about performance and alignment results.

5.3 COMPARISON

Compare methods 1 and 2 with regards to performance and quality.

Chapter

6

In this chapter I discuss the conclusions of this work and list some possibilities of future works.

CONCLUSIONS

6.1 FUTURE WORK

BIBLIOGRAPHY

Arca, S.; Casiraghi, E.; Lombardi, G. *CORNER LOCALIZATION IN CHESSBOARDS FOR CAMERA CALIBRATION*. 2005.

Armangué, X.; Salvi, J.; Batlle, J. A comparative review of camera calibrating methods with accuracy evaluation. *Pattern Recognition*, v. 35, p. 1617–1635, 2000.

Bajramovic, F. *Self-Calibration of Multi-Camera Systems*. [S.l.]: Logos Verlag Berlin, 2010. ISBN 9783832527365.

Baker, S.; Matthews, I. Lucas-kanade 20 years on: A unifying framework. *International Journal of Computer Vision*, Kluwer Academic Publishers, v. 56, n. 3, p. 221–255, 2004. ISSN 0920-5691.

Baker, S. et al. A database and evaluation methodology for optical flow. *International Journal of Computer Vision*, Springer US, v. 92, n. 1, p. 1–31, 2011. ISSN 0920-5691. Available at: <<http://dx.doi.org/10.1007/s11263-010-0390-2>>.

Barron, J.; Fleet, D.; Beauchemin, S. Performance of optical flow techniques. *International Journal of Computer Vision*, Kluwer Academic Publishers, v. 12, n. 1, p. 43–77, 1994. ISSN 0920-5691.

Bouguet, J. Y. *Camera calibration toolbox for Matlab*. 2008. Available at: <http://www.vision.caltech.edu/bouguetj/calib_doc/>.

Bouguet, J. yves. Pyramidal implementation of the lucas kanade feature tracker. *Intel*, 2000.

Bradski, G. Opencv. *Dr. Dobb's Journal of Software Tools*, 2000.

Bradski, G.; Kaehler, A. *Learning OpenCV: Computer Vision with the OpenCV Library*. [S.l.]: O'Reilly Media, 2008. ISBN 9780596554040.

Bruhn, A.; Weickert, J.; Schnörr, C. Lucas/kanade meets horn/schunck: Combining local and global optic flow methods. *International Journal of Computer Vision*, Kluwer Academic Publishers, v. 61, n. 3, p. 211–231, 2005. ISSN 0920-5691.

Burton, A.; Radford, J. *Thinking in Perspective: Critical Essays in the Study of Thought Processes*. [S.l.]: Methuen, 1978. (Psychology in progress). ISBN 9780416858402.

Chen, X.; Davis, J. Wide area camera calibration using virtual calibration objects. In: *In Proceedings of CVPR*. [S.l.: s.n.], 2000. p. 520–527.

Clarke, T.; Fryer, J. The development of camera calibration methods and models. *The Photogrammetric Record*, Blackwell Publishers Ltd., v. 16, n. 91, p. 51–66, 1998. Available at: <<http://dx.doi.org/10.1111/0031-868x.00113>>.

Douskos, V.; Kalisperakis, I.; Karras, G. *AUTOMATIC CALIBRATION OF DIGITAL CAMERAS USING PLANAR CHESS-BOARD PATTERNS*. 2007.

Furht, B. *Handbook of Augmented Reality*. [S.l.]: Springer, 2011. (SpringerLink : Bücher). ISBN 9781461400646.

Gibson, J. *The perception of the visual world*. [S.l.]: Houghton Mifflin, 1950.

Hanning, T. *High Precision Camera Calibration*. [S.l.]: Vieweg+Teubner Verlag / Springer Fachmedien Wiesbaden GmbH, Wiesbaden, 2011. (Vieweg + Teubner research). ISBN 9783834898302.

Harris, C.; Stephens, M. A combined corner and edge detector. In: *In Proc. of Fourth Alvey Vision Conference*. [S.l.: s.n.], 1988. p. 147–151.

Laureano, G. T.; Paiva, M. S. V. de; Silva, A. S. da. *Topological Detection of Chessboard Pattern for Camera Calibration*. 2013.

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. Available at: <<http://dl.acm.org/citation.cfm?id=1623264.1623280>>.

Medioni, G.; Kang, S. B. *Emerging Topics in Computer Vision*. Upper Saddle River, NJ, USA: Prentice Hall PTR, 2004. ISBN 0131013661.

Pati, U. C. *3-D Surface Geometry and Reconstruction: Developing Concepts and Applications: Developing Concepts and Applications*. [S.l.]: Information Science Reference, 2012. (Premier reference source). ISBN 9781466601147.

Peasley, B.; Birchfield, S. Replacing projective data association with lucas-kanade for kinectfusion. In: *Robotics and Automation (ICRA), 2013 IEEE International Conference on*. [S.l.: s.n.], 2013. p. 638–645. ISSN 1050-4729.

Rojas, R. *Lucas-Kanade in a nutsheGll*. 2014.

Scharstein, D.; Szeliski, R. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *Int. J. Comput. Vision*, Kluwer Academic Publishers, Hingham, MA, USA, v. 47, n. 1-3, p. 7–42, abr. 2002. ISSN 0920-5691. Available at: <<http://dx.doi.org/10.1023/A:1014573219977>>.

Shi, J.; Tomasi, C. Good features to track. In: . [S.l.: s.n.], 1994. p. 593–600.

Shumaker, R. *Virtual and Mixed Reality - New Trends, Part I: International Conference, Virtual and Mixed Reality 2011, Held as Part of HCI International 2011, Orlando, FL, USA, July 9-14, 2011, Proceedings*. [S.l.]: Springer, 2011. (Information Systems and Applications, incl. Internet/Web, and HCI). ISBN 9783642220203.

Sun, D.; Roth, S.; Black, M. A quantitative analysis of current practices in optical flow estimation and the principles behind them. *International Journal of Computer Vision*, Springer US, v. 106, n. 2, p. 115–137, 2014. ISSN 0920-5691. Available at: <<http://dx.doi.org/10.1007/s11263-013-0644-x>>.

Tillapaugh, B.; Engineering, R. I. of T. C. *Indirect Camera Calibration in a Medical Environment*. [S.l.]: Rochester Institute of Technology, 2008. ISBN 9780549934479.

Tsai, R. Y. An efficient and accurate camera calibration technique for 3D machine vision. In: *Proc. Conf. Computer Vision and Pattern Recognition*. Miami: [s.n.], 1986. p. 364–374.

Zhang, Z. A flexible new technique for camera calibration. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, IEEE, v. 22, n. 11, p. 1330–1334, 2000.