



Universitat de les
Illes Balears



Treball Final de Carrera

ENGINYERIA INFORMÀTICA

Ponster, tu app de realidad aumentada

Álvaro Medina Ballester

Tutor

Ramón Mas Sansó

Escola Politècnica Superior
Universitat de les Illes Balears
Palma, September 11, 2014

CONTENTS

Contents	i
1 Abstract	1
2 Thanks	3
3 Intro	5
4 State of the art	7
4.1 Object recognition	7
4.2 Template matching	8
4.3 Feature detection	8
4.3.1 Speeded-Up Robust Features - SURF	9
4.3.2 Features from Accelerated Segment Test - FAST	10
4.3.3 Oriented FAST and Rotated BRIEF - ORB	10
4.4 Matching	10
4.4.1 Brute-Force Matcher	10
4.4.2 Fast Library for Approximate Nearest Neighbors - FLANN	11
4.5 Natural feature tracking	12
5 Technologies	13
5.1 OpenCV	13
5.2 FastCV	13
5.3 Vuforia	13
6 Development	15
7 Tracking	17
8 Conclusions	19

CHAPTER
1

ABSTRACT

CHAPTER
2

THANKS

CHAPTER 3

INTRO

Augmented reality has become a very popular topic in the last five years. With the introduction of mobile smartphones, developers started to have the chance to develop applications using the powerful CPUs and GPUs of those devices.

CHAPTER 4

STATE OF THE ART

The technique of mixing real world elements with virtual elements displayed on the screen of a device is what we call augmented reality. In the field of augmented reality, a lot of things have happened during the last years. The progress in the fields of computer vision and image processing have led to several new techniques of detection and tracking. This, combined with the increasing availability of powerful mobile devices, has enabled developers to build a plethora of high quality AR-based applications. Nowadays, modern mobile devices use integrated cameras, motion sensors and proximity sensors to make these AR-based experiences.

Augmented reality in mobile devices is slightly different from what can be seen in desktop environments. Although every year we have more powerful mobile smartphones (citation needed), processing and drawing into the device's screen is still an expensive operation in terms of computational cost. This is one of the reasons why cost-efficient computer vision algorithms and techniques have emerged in the past years.

Most of the augmented reality apps follow this behaviour:

- Get the input from the camera or a video.
- Search for an object of interest.
- Introduce our object into the scene, considering the camera or the input position.

In this chapter we are going to describe which are the different techniques that enables us to do them.

4.1 Object recognition

In order to provide an augmented reality experience, we have to know first which is the real world element that we are going to use as a reference to mix the real world input with our virtual elements. This reference can be from an image from the smartphone

camera to the user location. Ponster searches a particular image inside the camera input in order to draw the poster image above. In computer vision, this technique of searching an image and follow it along it's movement is usually referred as object tracking.

In computer vision there are a lot of object recognition techniques. In the development of Ponster, two of these techniques have been tested: template matching and feature-based detection. Detecting the image in a continuous input from the smart-phone camera, taking account of scale, rotation and perspective differences, becomes an object tracking technique.

4.2 Template matching

Template matching consists of finding areas of an image that are similar to a provided template image. We have to provide a template image (the image that we want to look for) and compare it with the source image (the image in which we want to search). Template matching is also called area-based approach.

OpenCV provides a method to perform template matching with several methods, such as SQDIFF or CCORR. With the latest, CCORR, we use a correlation formula to check if the template is inside the image. Instead of applying a yes/no approximation, we can bring a positive match with a certain threshold.

Performing a template matching operation using OpenCV on mobile devices is fast enough to deliver a smooth 25/30fps-like detection. However, match template does not take account of scale, rotation and perspective invariance by itself. There are several approaches to bring invariance to match template. For instance, image pyramids are used to make match template scale and rotation invariant, but it is not part of the OpenCV match template function, although it provides some methods to implement image pyramids.

Match template has been tested during the development of Ponster. Also, a basic image pyramid system has been developed for scale-invariance, but match template has been discarded in favor of feature detection algorithms because rotation and scale invariance and perspective warp are required features.

4.3 Feature detection

A feature-based approach can be presented as a three step method[2]. First of all, we have to detect keypoints[1] (also called interest points) in the image. Usually, interest points are corners, blobs or T-junctions[2]. A good keypoint is a *repeatable* keypoint; if we can find the same keypoint under different conditions such as light difference or rotation, it's considered as a quality keypoint. The second step is to compute *descriptors* or feature vectors. These descriptors are represented as neighbourhoods of interest points. Assuming that feature detection makes sense when we have two images to compare, these two steps have to be performed on both images. Once we've done that, we have a group of descriptors for each image, and we have to compare them in order to *match* features. If the features of the source image are present in the input image, we can assume that the object has been detected. The matching is based on the distance between the feature vectors.

Usually, source images with enough keypoints are easier to detect than more uniform images. This is why it's better to select a good source image with many features and good contrast. As we've said before, in order to deliver a good augmented reality experience, we need to make our detection algorithm scale, rotation and perspective invariant. Feature detection techniques can be scaled invariant by extracting features that are invariant to scale, such as feature vectors computed from interest point neighbourhoods. For rotation invariance, algorithms can estimate the orientation of the keypoint.

There are plenty of feature-detection based algorithms, many of them based on Scale-Invariant Feature Transform, or SIFT. Cost efficiency is one of the most important features of these algorithms, as every new technique introduced tries to maintain robustness and reducing computation time. Robustness it's also very important, but less robust algorithms are also been developed in favour of reducing computation time. One good example is FAST[orb paper] keypoint detector, which is not rotation invariant. Next, we are going to describe the feature-detection algorithms tested on Ponster: SURF, FREAK and ORB.

4.3.1 Speeded-Up Robust Features - SURF

Speeded-Up Robust Features, also know as SURF, is a group of detector and descriptor introduced by Herbert Bay et al. SURF is faster and more robust than other alternatives like SIFT[2]. It's descriptors are rotation and scale invariant. Perspective transformations are also considered, but in lower order.

The keypoint detection in SURF uses a Hessian-matrix approximation. This use of integral images reduces computational cost in comparison with another interest point detection techniques such as Harris corner detection. Scale invariance is achieved by calculating integral image pyramids, but instead of reducing the image size, integral images allow SURF to upscale and build the pyramids more efficiently.

The SURF descriptor calculation is slightly based on SIFT. SURF descriptor describes the distribution of the intensity content within the interest point neighbourhood, which is similar to the gradient information used by SIFT. It is done in two steps, fixing a reproducible orientation based on information from a circular region around the keypoint, and then building a square region aligned to the calculated orientation. Once we have the descriptors, the last step is to perform the matching. Descriptors are compared only if they have the same type of contrast, allowing to perform a faster matching. In Ponster, two different matching algorithms have been tested, Brute-force Matcher and FLANN-based matching.

SURF is as robust as other alternatives such as SIFT, but it's faster to compute due to the use of integral images. It is rotational and scale-invariant, which is better than the template matching technique described before, but it's performance running on the device (iPhone 5, iOS 7.1.2) is not good enough to deliver a decent user experience, taking between 0.7 and 1 second to compute each image. Also, SURF is a patented algorithm and it's not intended to use it in commercial applications.

4.3.2 Features from Accelerated Segment Test - FAST

FAST is a keypoint detector based on corner detection. It was introduced by Edward Rosten and Tom Drummond[3] and its primary purpose was to bring a real time interest point detector. FAST considers a circle of sixteen pixels around each corner candidate, and detects a candidate as a corner if there are n contiguous pixels in that circle with brighter intensity than the candidate pixel.

This technique is faster than others for corner detection, but FAST is not scale and rotation invariant. Also, it does not perform very well under high noise images. Many other techniques use FAST as a starting point of a detector, bringing scale and rotation invariance and a corresponding extractor. One example of this is ORB, tested in the development of Ponster.

4.3.3 Oriented FAST and Rotated BRIEF - ORB

As we have said before, real time performance has been a very popular topic in object detection during the last years. The main characteristic of ORB is to perform as good as SIFT, but doing it twice as fast. ORB uses a variant of FAST as the interest point detector, and BRIEF as the descriptor extractor.

FAST does not have rotation invariance. This is why ORB uses oFAST (FAST keypoint orientation), which is a variant of FAST that computes orientation by intensity centroid. This technique assumes that a corner's intensity is offset from its center, and this vector may be used to impute an orientation [paper]. To bring scale invariance, ORB employs a scale pyramid of the image and computes FAST on each level of the pyramid.

ORB also uses a variant of BRIEF called rBRIEF, or Rotation-aware BRIEF. The BRIEF descriptor, unlike SURF, is a binary descriptor. rBRIEF is based on steered BRIEF, which uses the keypoint orientation; in addition to steered BRIEF, a learning method for choosing good binary features is applied, resulting into rBRIEF.

In Ponster, ORB has been tested with better results than the other previous techniques, but again with poor performance in the device. Only a 15 fps processing have been achieved, with slightly worst detection than with SURF.

4.4 Matching

Once we have calculated the interest points and computed the descriptors in the two images that we want to compare, we have to perform a match between these two sources. Depending on the descriptor extraction method, one or another matcher must be used. ORB uses binary descriptors, but SURF does not, so the matching is performed in a different way.

We will discuss two of these methods, both used in the development of Ponster, Brute-force matcher and Fast Library for Approximate Nearest Neighbors.

4.4.1 Brute-Force Matcher

Brute-Force Matcher, as its name states, will compare each of the descriptors found in the images performing a linear search. Although it may seem that this approach is not very efficient, BF-matcher performs really well on binary descriptors like ORB.

The comparison is done by a distance function. There are many functions in BF-matcher:

- NORM_L1 better with SURF/SIFT.
- NORM_L2 better with SURF/SIFT.
- NORM_HAMMING better with ORB.
- NORM_HAMMING2 better with ORB.

Hamming distance can be computed with bit manipulation operations, which are very quickly. In Ponster, Hamming distance has been tested for ORB and L2 normalization with SURF.

4.4.2 Fast Library for Approximate Nearest Neighbors - FLANN

Instead of performing a linear search for matching descriptors, we can use a nearest neighbour matching technique. The nearest neighbour search tries to find, given a set of points P in a vector space X , all the points that are close to a given point q . FLANN is a library that enables us to perform this kind of searches with several algorithms. Two have been tested in Ponster, randomized KD-tree search for SURF and Locality-Sensitive Hashing for ORB.

Randomized KD-tree

Basic KD-tree search performs well for small datasets, but quickly degrades its performance when the dimensionality increases. Several KD-tree algorithm variants have been introduced, such as approximate nearest neighbour, in order to reduce the computational cost of KD-tree searches with large datasets.

This algorithm creates multiple randomized KD-trees, built by choosing the split dimension randomly from the first 5 dimensions on which data has the greatest variance. Then, a priority queue is created while searching all the trees, so the search can be ordered by increasing distance to each bin boundary.

Using this approximation techniques can boost performance by reducing the precision of the matching, although the loss is usually small enough to maintain a 95% precision.

Locality-Sensitive Hashing

Locality-Sensitive Hashing is a matching algorithm to solve the nearest neighbour search in high datasets. LSH is used with binary descriptors like the ones computed with ORB. The main idea of LSH is to hash the points with functions that ensure that close points will be more likely to key collision, thus allowing to get the nearest neighbours of each point querying the other points in its bucket.

The LSH parameters defines the hash functions *amplification*. This means that the hash functions must be *amplified* enough to ensure hash collision; otherwise, the algorithm would be useless. The effect of this parameters and a more in-depth explanation of LSH can be found in this[PDF] paper.

4.5 Natural feature tracking

Although the technologies explained in previous sections are robust and reliable, and have been successfully tested, they have not been used in the last version of Ponster due to its performance on mobile devices. Another technique has been proven to be successful to deliver both high performance in mobile devices and robust object tracking.

Natural feature tracking consists of computing the motion of a feature in the scene[cita aqui]. We call natural feature to any point or region candidate to be detected and selected as a reference. Usually this follows the next steps:

- Natural feature detection and selection.
- Motion estimation based on detected features.
- Evaluation feedback for stabilized detection and tracking.

The feature detection consists in selecting points and regions in the image with specific characteristics that are easy and robust to track. The motion estimation can be executed by several ways. For example, in [paper], optical flow is presented to estimate the camera movement. In our mobile environment, we can also use the gyroscope. The evaluation feedback provides a way to reject detected features that do not specifically match with the plane that should be representing, thus allowing the algorithm to avoid false positives and do not corrupt the tracking output.

The technology used in Ponster to perform the augmented reality is called Vuforia and uses Natural feature tracking to track the object in the scene.

CHAPTER 5

TECHNOLOGIES

Describe here Vuforia, OpenCV, FastCV, UIKit, etc.

5.1 OpenCV

5.2 FastCV

5.3 Vuforia

CHAPTER
6

DEVELOPMENT

CHAPTER 7

TRACKING

CHAPTER 8

CONCLUSIONS

