

"Stereo Vision for 3D Measurement in Robot Systems"

Asst. Prof. Dr. Muhyi. AL-Azawi
Electrical Eng. Dept.
AL-Mustansiriyah Univ.
Baghdad /Iraq.

Asst. Prof. Dr. Dhafer R. Zaghar
Electrical Eng. Dept.
AL-Mustansiriyah Univ.
Baghdad /Iraq.

M.Sc. Student Saba Ali Sadeq
at Electrical Eng. Dept.
AL-Mustansiriyah.
Baghdad /Iraq.

Abstract:

The major obstacle in the application of stereo vision to extract 3D information is the error in disparity map and highly computational cost at conventional computers. This paper is concerned with finding software solution to obtain tradeoff between disparity map accuracy and reductions in execution time through developing the stereo vision algorithm. This algorithm is based on block matching technique, in which an image is partitioning into blocks. An analysis of the essential parameter of this technique (the size of block) is performed to obtain the optimal solution. Adaptive block size with vertical edge detection has been adopted to implement the proposed algorithm. The significance of this work is reduction of the execution time of the stereo vision algorithm, where the execution time of the proposed algorithm is about (2%) of the required time to execute standard vision algorithm, when running at conventional computers.

Key words: Stereo vision, disparity mapping, block matching, 3D vision, adaptive block.

"الرؤية المجسمة لمقياس ثلاثي الأبعاد في أنظمة الإنسان الآلي"

أ.م.د. محي العزاوي
قسم الهندسة الكهربائية
كلية الهندسة/الجامعة المستنصرية

أ.م.د. ظافر رافع صغير
قسم الهندسة الكهربائية
كلية الهندسة/الجامعة المستنصرية

الطالبة صبا علي صادق
قسم الهندسة الكهربائية
كلية الهندسة/الجامعة المستنصرية

المستخلص:-

إن العقبة الرئيسية في تطبيق الرؤية المجسمة (Stereo Vision) لإستخلاص البعد الثالث هي الخطأ في مخطط التباين والتكلفة الحسابية العالية في الحواسيب التقليدية, لذا يتناول هذا البحث إيجاد حلول برمجية للحصول على موازنة بين دقة مخطط التباين (Disparity Map) وتقليل وقت تنفيذ الخوارزمية من خلال تطوير خوارزمية الرؤية

المجسمة. تستند هذه الخوارزمية على تقنية تطابق اجزاء الصور, حيث يتم فيها تجزئة الصور الى مقاطع (Blocks). وأجري تحليل للمعيار الأساسي لهذه التقنية (والذي يمثل حجم المقاطع block size) للحصول على نتائج مثالية. تم استخدام حجم مقطع تكيفي (متغير) (Adaptive block size) مع خاصية إظهار الحافات العمودية في الصور (Sobel Filter) لتنفيذ الخوارزمية. تتمثل أهمية هذا العمل في تقليل وقت تنفيذ خوارزمية الرؤية المجسمة, حيث استغرق وقت تنفيذ الخوارزمية المقترحة حوالي (٢%) من الوقت اللازم لتنفيذ الخوارزمية القياسية عند تنفيذها في الحواسيب التقليدية.

الكلمات المفتاحية: الرؤية المجسمة, خارطة التباين, تطابق المقاطع, رؤية ثلاثية الابعاد, مقطع تكيفي.

١. Introduction

For safe operation in uncontrolled environments of robots, dependable 3D perception modules are required. Commonly used 3D sensors for mobile robots are laser range finder and time-of-flight cameras. Although these types have advantage of delivering accurate 3D data, but suffer from low resolution, and cannot provide enough information about the objects [1]. A promising alternative for robot navigation and mapping is vision sensors/cameras. Which are low-cost, light and compact, easily available, have low power consumption, and provide rich information about the environment [2]. The stereoscopic camera (Binocular System) can provide the third dimension by grabbing left and right images simultaneously, like as human's binocular vision [3]. In this paper binocular vision hardware module has been implemented using two cameras to imitate human vision. The proposed vision algorithm has been implemented and tested with real stereo images to illustrate its performance, and a comparison with classical block matching algorithm will be introduced.

٢. Stereo Vision Principle

In stereo vision, the same seen is captured using two cameras, displaced horizontally from one another, to obtain two different views on a scene from two different angles. The two captured images have a lot of similarities and smaller number of differences. In human sensitivity, the brain combines the two captured images together to get a 3D model for the seen objects. In machine vision, the 3D model for the captured objects is obtained by finding the similarities between the stereo images and using projective geometry to process these matches. The difficulties of reconstruction using stereo is finding matching correspondences between the stereo pair.

The computer compares the images while shifting the two images together over top of each other to find the parts that match. The amount of the *shift* between two matched pixels is called “*disparity*”, which relates to the object distance. The set of displacements between matched pixels is usually indicated as “*disparity map*” [4]. The higher disparity of object pixel means that the object is closer to the cameras and appears *brighter* in disparity map, while the less disparity means the object is far from the cameras appears *darker*. In addition, if the

object is very far away, the disparity is zero that's means the object on the left images is the same pixel location on the right image [7]. Figure (1) depicts the geometrical basis for stereoscopic images by using two identical cameras.

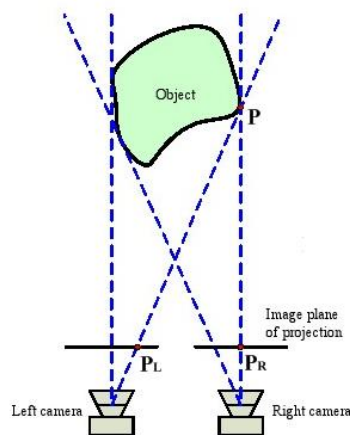


Figure (1): The position of two cameras and their image Planes of projection.

These cameras are set up on the same plane and turned in the same direction. The image planes are presented in front of the cameras for ease to model the projection (like as the model of pinhole camera) [7].

3. Epipolar Geometry

Epipolar geometry is a simplest stereo camera geometry that has optical axes that are *parallel* and *normal* to the *baseline*, the line connecting the lens centers of the cameras. This geometry depicted in Figure (2). This model shows two different perspective view of an object point (P) from two identical cameras center (O_L & O_R), which separate only in x direction by a *baseline* distance (b). The points P_L and P_R in the image plane are the perspective projections of (P) in left and right view, which are called a *conjugate pair*. The plane passing through the camera centers and the object point (p) in the scene is called the *epipolar plane*. The intersection of the epipolar plane with the image plane is called *epipolar line*. By referring the epipolar geometry, correspondences at points P_L and P_R must lie on the epipolar line. In this case, corresponding epipolar lines are horizontal and have the same y -

coordinate. This implies that a one-dimensional search is sufficient to find the correspondences.

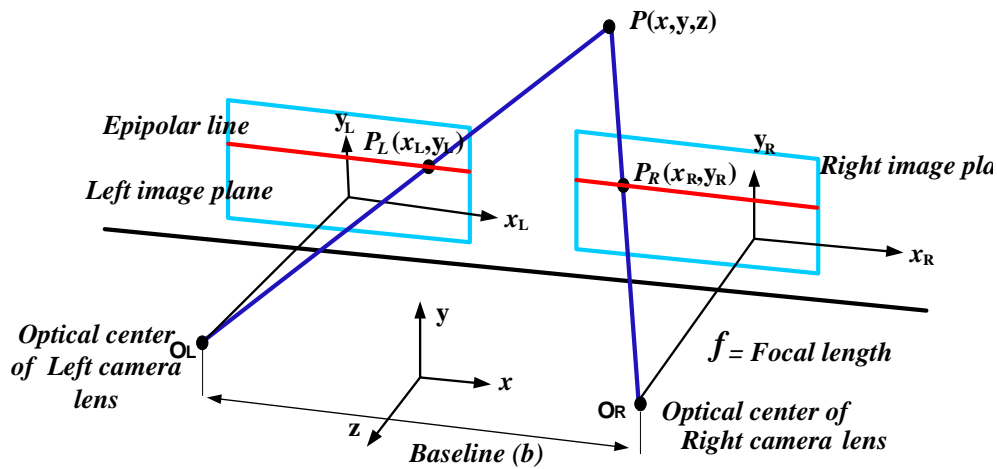


Figure (۲): The epipolar geometry of stereoscopic vision.

The third dimension or the depth (Z) can be determined by comparing the triangles which are similar:

$$\frac{x_L}{f} = \frac{x + \frac{b}{2}}{z} \quad (۱)$$

$$\frac{x_R}{f} = \frac{x - \frac{b}{2}}{z} \quad (۲)$$

$$z = \frac{b \cdot f}{x_L - x_R} = \frac{b \cdot f}{d} \quad (۳)$$

Where: the disparity is $d = (x_L - x_R)$. x_L and x_R are the x-axis coordinate of point P in the left and right image planes respectively, and x is x-axis of stereo-camera coordinate that considered to be the midway between the left and right camera coordinate systems. f is the focal length that represents the distance between the image plane and the optical center of camera lens. The baseline (b) is the distance between the lens centers of the cameras.

Measurement of depth from equation (۳) requires knowledge about the coordinates of corresponding points in the images [۶].

۴. Block Matching Correspondence

Extracting depth information from two or more cameras requires solving the correspondence problem. A block matching based on local (area-based) method is used, since it produced a dense disparity maps and can be implemented effectively. Indeed analyzing only pixels is not sufficient, because there may be several candidates that have the same gray value. Block matching method uses the intensity values of pixels within a neighborhood called block or (window) to find matching pixels in two images. For example, this happens

with blocks of sizes $[o \times o]$ or $[V \times V]$. The blocks are usually defined on epipolar line for matching ease. Each block from the left (reference) image is matched into a block over the searching area of pixels in right image as shown in Figure (3). The *shift* that gives a best result of the matching criteria is considered as the best match or correspondence [1][4].

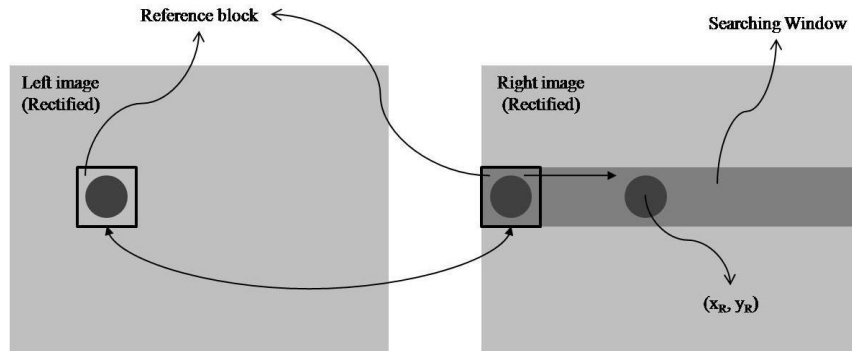


Figure (3): The block matching technique.

3. Criteria of Algorithm Performance

Matching algorithm is based on two main criteria, which would be used to support the analysis of matching algorithm results:

1. **Execution time:** represent the execution time of algorithm, and it is different from image to another.
2. **Reliability:** the information about objects should be accurate enough to recognize, and avoid obstacles.

In order to demonstrate the effectiveness of the implemented algorithms, it have been tested with many real stereo image pairs, which provided on [4] with ground-truth (true disparity map) that required to calculate the error percentage. The general approach used to compute error statistics with respect to ground truth data is:

- **The percentage of bad matching pixels (p);**

Is the error (measured in disparity units) between the computed depth map $d_C(x,y)$ and the ground truth map $d_T(x,y)$, as shown in Equation (4) [4].

$$p = \frac{1}{N} \sum_{(x,y)} |d_C(x,y) - d_T(x,y)| \quad (4)$$

where N is the total number of pixels.

4. Standard Matching Algorithm (SMA)

This algorithm represents the standard algorithm in stereo vision, which is based on block matching [4]. Using this algorithm as a first pass for determining a disparity map and to illustrate its result. Many statistical measurements for the intensity values are used to check the correspondence between two blocks such as Sum of Absolute Differences (SAD), Sum of Square Differences (SSD) and Normalized Cross Correlation (NCC). These functions give equivalent results in terms of accuracy, but (SAD) requires less computational time compared

with others; therefore, it will be adopted in this paper. Extraction disparity map is achieved by moving the reference window in left image *pixel-by-pixel* (full search) over the target windows in right image within a specific search range called *disparity range* and compute (SAD) at each moving step, the minimum SAD calculated between two windows representing that, these windows are *matched*. The amount of the *shift* between matched windows represents the *disparity value* (d) of the *central pixel* of the reference image. This procedure is repeated for whole pixels in the reference image to compute their disparity values. These values are sorted to produce *disparity map*. The SAD is calculated according to the Equation (9):

$$\text{SAD} = \sum_{(i,j) \in W(x,y)} |I_L(i,j) - I_R(i,j+d)| \quad (9)$$

where $W(x,y)$ is a window surrounding the position (x,y) , d is a disparity value, I_L and I_R are the intensity values in the left and right images respectively. This algorithm has been implemented with MATLAB (7,7) technical programming language, and executed with different stereo image pairs at different window sizes to illustrate its results. One of the stereo pairs called (Tsukuba) with size of 384×288 shown in Figure (4) would be chosen as an example to compare the results of the algorithms.

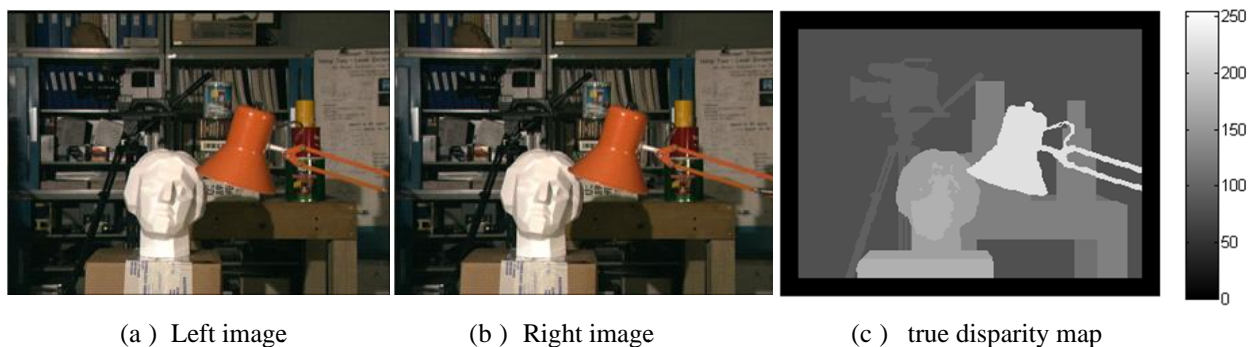


Figure (4): Tsukuba stereo image pair with true disparity map.

7.1. Summary of SMA

The matching algorithm SMA which implemented by MATLAB, can be summarized as the following steps:

Step 1. Convert left and right images from (RGB scale) color image to Gray scale image, to reduce the algorithm's complexity, which is multiplied with the number of color channels used.

Step 2. Prepare a 2D array (x, y, d) for disparity map, where (x, y) represent the coordinates of each pixel in the reference image (left), and d is for the disparity value of these pixels.

Step 3. Set the size of matching window, and Set initial match value (Initial SAD).

Step 4. Set the value of "Search Range" ($d = 1, \dots, d_{\max}$) that represents the search range within a loop. The value of (d_{\max}) is empirical and represents the maximum value of the shift in pixels between the left and right images.

Step 6. Compute SAD between reference window and candidate windows according to the Equation (6) at each shift.

Step 7. Iteratively update the initial SAD value with every new minimum SAD computed, until the match value converges.

Step 8. For each pixel (x, y) of the reference image, find the disparity element (x, y, d) . The set of disparity values represents the disparity map; the disparity is often treated as synonymous with inverse depth.

The result of matching algorithm is a disparity map that can be shown as image, which represents the $1/D$ of the $1/D$ images (left and right), as shown in Figure (6). The disparity is inversely proportional to depth, it can be seen that the closest object to the stereo cameras has large disparity value and appears in **bright gray color**. On the other hand, the farthest object has less disparity value and appears in **dark gray color**.

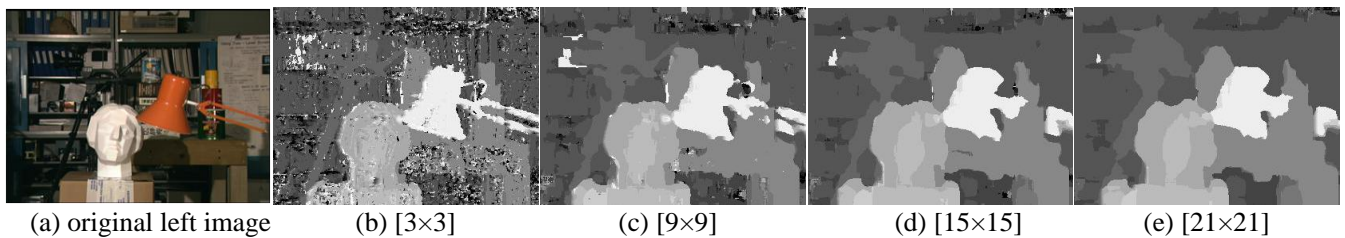


Figure (6): Shows (SMA) results, (a) represent Tsukuba left image, (b), (c), (d), and (e) are the disparity maps at different window sizes.

The results of the (SMA) in Figure (6) show that, at small window size the map seems noisy due to multi-matching problem, in which for each reference window in the left image there is **many matching windows** in the right image, because using too small window which does not cover enough intensity variation gives poor disparity estimation, especially at **texture-less region**. Additionally at too large window size, high error occurs due to ambiguity problem, this ambiguity arises from covering multi-objects by the window. Consequently, the matching result caused blurring at objects borders, as shown in Figure (6.e).

In general, this algorithm gives good disparity maps with less error as shown in (Table 1 column 3), because it is based on moving the reference window **pixel-by-pixel** over the target image (full search technique), and calculate the **disparity** for each pixel. But the drawback of this method represented in highly processing time. This time increases as the window size increases, due to increase the number of pixels within the window, thus more computational operation is required, as shown in (Table 1 column 4). This reason makes this algorithm unsuitable for robot application that required fast performance. Therefore, another algorithm has been proposed to overcome these drawbacks. Variable window size (adaptive window) is adopted to solve these problems with vertical edge detection to reduce the computation that results in decreasing the time.

4. The Proposed Matching Algorithm (SABMA)

In this method the search for the matched-blocks is achieved by moving the reference window *Block-by-Block* over the target windows. This method achieves a *large reduction* in execution time. In which the same disparity (depth) is given to all objects that exist in the same window, but this caused blurring at object borders when the window covers Multi-objects, as shown in Figure (6).

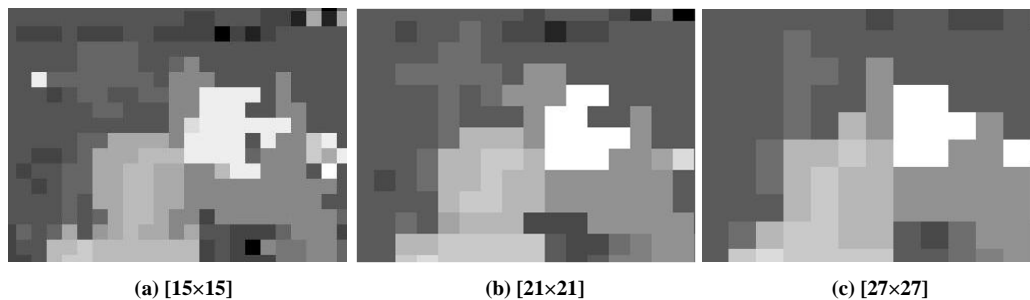


Figure (6): The blurring effect when applying Block-by-Block searching technique.

Therefore, to avoid this problem and the problems mentioned previously in SMA at objects borders, the selection of window size must be large enough to avoid unwanted noise, but small enough to avoid ambiguity and blurring. Therefore, *adaptive window size* is adopted, in which the window size changes adaptively according to local disparity characteristics.

The proposed algorithm is suggested to use *adaptive window* with *edge filter detection (Sobel filter)* to discriminate the borders of the objects in the scene in order to calculate the true disparity value for each object, as shown in Figure (7). The proposed algorithm will be called Sobel Adaptive Block Matching Algorithm (SABMA). The following three consecutive steps explain this method:

Step 1. Objects Boundary Detection

Vertical edge detection is performed for both stereo images by using *Sobel filter* with suitable threshold value (T). Where this value determines the amount of the detected edges.

Step 2. Select Window Size Adaptively

Two windows with different sizes are used:

1. Large window $[a \times a]$ is used to solve the problems of texture-less regions, and to accelerate the matching process.
2. Small sub-window $[b \times b]$ is used to solve the problem of ambiguity at objects boundaries. Where (b) is chosen as a part of (a) to increase the accuracy, e.g., $b \leq \frac{1}{3}a$.

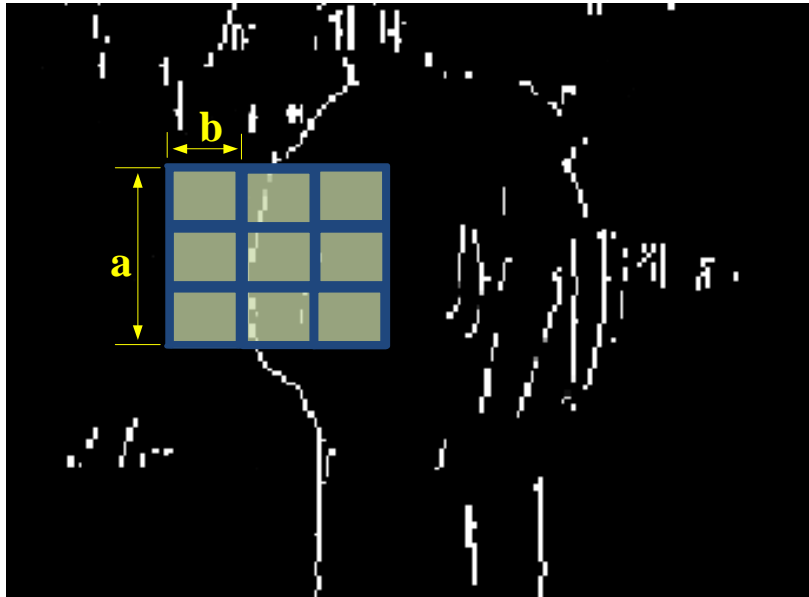


Figure (7): The division of large window [$a \times a$] into many Sub-Windows [$b \times b$] at object border under Sobel filter.

Step 5. Perform Matching Process

The search for the matching blocks is achieved by, moving the reference window *Block-by-Block* over the target windows in order to speed up the matching performance.

The division of window size is performed depending on the filtered image (*Sobel vertical edge detection*). Firstly large window size is used, and a conditional step is done on filtered image to check the presence of edges. If the window contained edges that means, there are many objects within the window that required dividing the window to obtain *accurate* depth result. Else, if the window doesn't contained edges which means that, the pixels within the window belong to same object, thus all these pixels have the same depth. Therefore, there is no need to divide the window. The effectiveness of this method can be seen when comparing Figure (7) that represents block-by-block searching technique without adaptive window, and Figure (9) that represents using the same searching technique with adaptive window size. Therefore, SABMA accomplishes compromise between the accuracy that achieved by adaptive window size, and reduction in execution time that achieved by block-by-block searching technique.

To illustrate the effectiveness of (SABMA), it has been executed with different images at different window sizes, as shown in Figures (8) and (9):

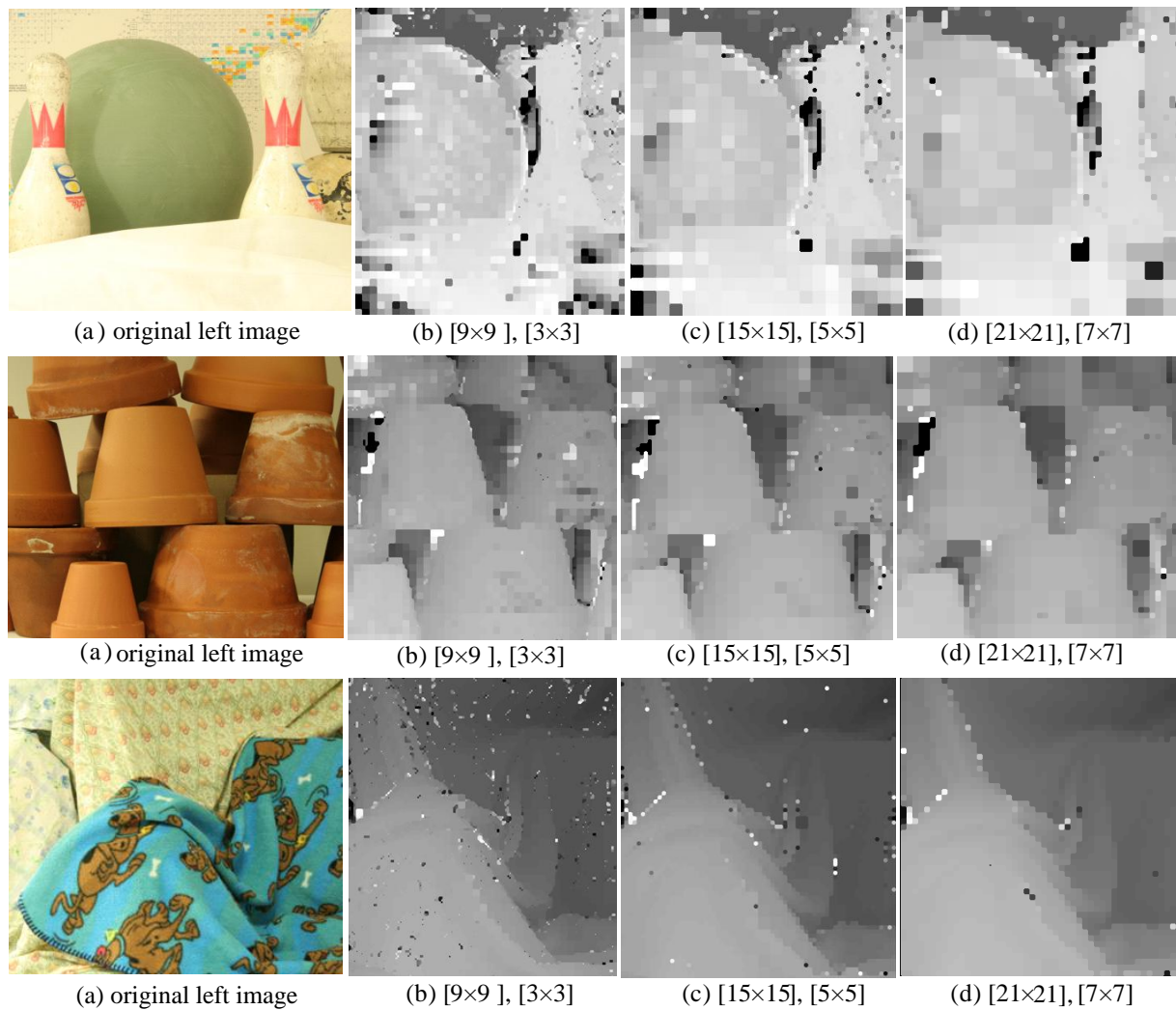


Figure (8): Disparity maps of (SABMA) algorithm at different images and different window sizes.

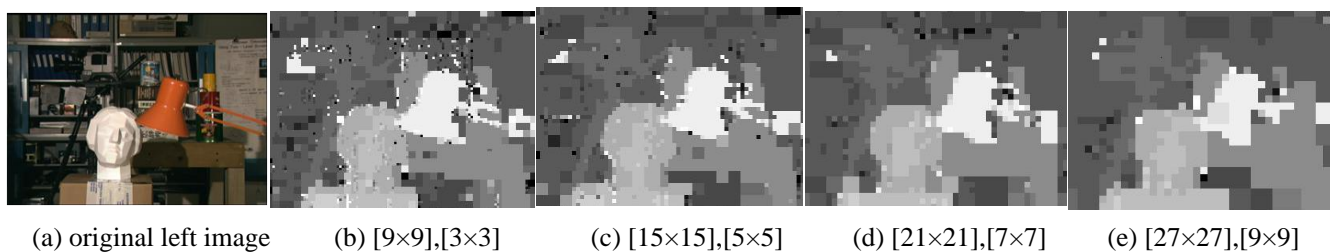


Figure (9): Tsukaba left image with disparity maps at different window sizes.

A comparison has been performed between SMA and SABMA in error percentage, and execution time for Tsukaba image pair, as shown in Table (1):

Large window size	Sub-window size	Error percentage(%) of SMA	Error percentage(%) of SABMA	Execution time (Sec.) of SMA	Execution time (Sec.) of SABMA	Execution time Percentage of SABMA with respect to SMA
-------------------	-----------------	----------------------------	------------------------------	------------------------------	--------------------------------	--

[9x9]	[3x3]	11,60 %	14,20 %	8,18	0,3	3.77 %
[10x10]	[5x5]	9,93 %	11,94 %	17.04	0,28	1.74 %
[21x21]	[7x7]	11,89 %	10,81 %	28.00	0,26	0.91 %
[27x27]	[9x9]	13,82 %	13,44 %	42.00	0,20	0.08 %

Table (1): The execution time and error calculation comparison of Tsukaba image pair between (SMA) and (SABMA) algorithms.

Table (1) shows the reduction in execution time of SABMA as compared with SMA and the same amount of reduction has been achieved with different image pairs. This reduction is attributing to the searching for the matched-blocks technique which based on moving the reference window *Block-by-Block* over the target windows instead of *Pixel-by-Pixel* searching technique that used in SMA. Thus a large amount of computational operation required to calculate SAD is reduced as shown:-

- **For Pixel-by-pixel searching technique**

Number of SAD computation required to extract disparity map for MxN image size = $M \times N \times \text{disparity range}$.

- **For Block-by-Block searching technique**

Number of SAD computation required to extract disparity map for MxN image size = $\frac{M \times N \times \text{disparity range}}{a^2}$.

Figure (10), shows a comparison in execution time of the two algorithms.

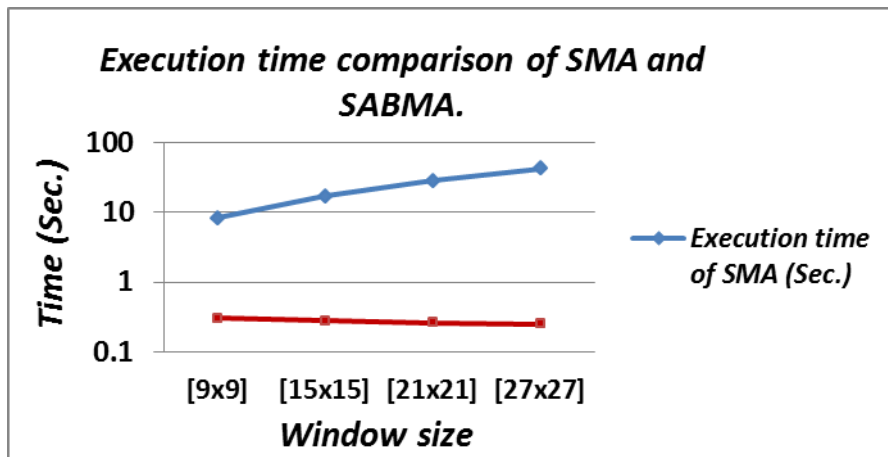


Figure (10): Comparison in Execution time between two algorithms.

Also Table (1) shows that, at each point the expected value of the error by the SABMA is approximately equal to that produced when any fixed-size window is used (at SMA). This

attribute to use adaptive window size that achieves an enhancement in extracting disparity values at object borders, as show in Figure (11).

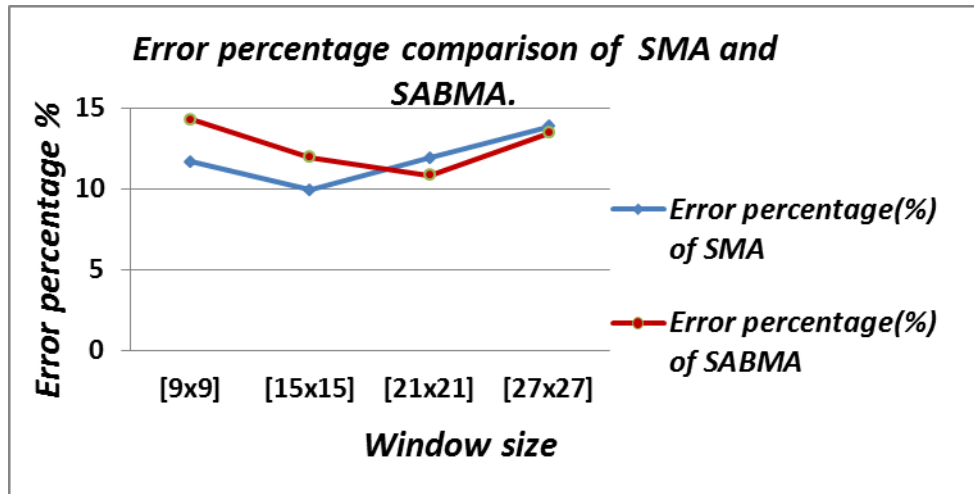


Figure (11): Comparison in Error percentage between two algorithms.

Λ. Implementation and Experimental Results

The stereo vision hardware used in this research, is implemented using two identical models USB web cameras, connected to PC (P4, core i2, 3.5 GHz processor and 3GB RAM) with USB (Universal Serial Bus) cable. As the baseline increases, there is a possibility to detect a disparity in object that farther away, since the baseline b is proportional to the range (depth) Z according to Equation (3). Therefore, in this research the baseline is 16cm selected, to compromise the increases of range detection and reduces non-overlap between two images, as shown in Figure (12).

In order to calculate the accuracy of the estimated depth Z , a comparison has been made between the estimated depth values, and the real depths of some selected points. Figure (13) shows stereo images (left and right) with size of 320×240 , captured from the implemented stereo cameras, and the disparity map which calculated with respect to left image.

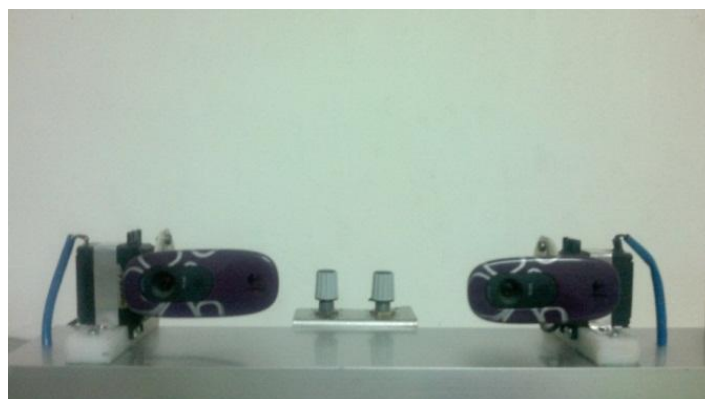


Figure (12): The implemented stereo cameras.

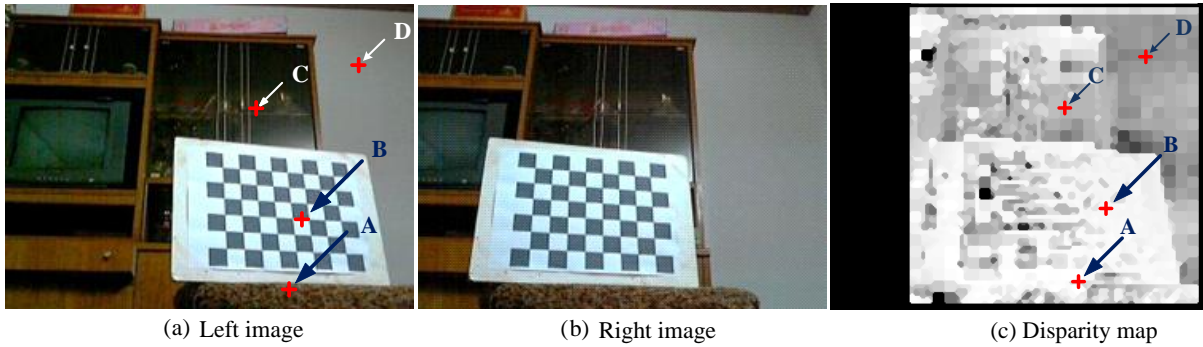


Figure (13): Selecting different points at different depths.

Table (2), shows the disparity values and depths for the selected points:

Table (2): Results of depth measurement and error percentage of the estimated depth.

Tested points	disparity value (pixels)	Real depth (cm)	Estimated depth (cm)	Error percentage (%)
A	102	68	67.81	0.2 %
B	83	83	83.34	0.41 %
C	28	240	247.00	0.83 %
D	24	290	288.22	0.61 %

Table (2) shows that, the disparity value decreases for the objects which are far from the camera, and it appears darker in disparity map and vice versa. Also the estimated depths that extracted from disparity map are approximated to the real objects depth with less error.

9. Conclusions

The following summarize the main important points that are concluded from this work:-

1. The execution time of the SABMA reaches about 2% of the required time to execute SMA at conventional computers, which attribute to the following proposed methods:
 - i. Sobel filter has a noticeable effect in reducing the time. It is found that the execution time of SMA with Sobel filter is about 92.2% of the time required to execute this algorithm without Sobel filter.
 - ii. **Block-by-Block** searching technique has a great influence on reducing the computational operation required to calculate SAD for disparity extraction.
2. Selection of window size does not only influence the execution time of the algorithm, but also the accuracy of disparity map. Therefore, the use of adaptive window size with Sobel filter helps in enhancing the reliability of the disparity map, by reducing the ambiguity at objects borders. This can be seen when less error is obtained at large window size, e.g., [21x21] and [27x27] of SABMA compared with SMA, as shown in Table (1).
3. When stereoscopic vision hardware has been implemented using two cameras, the experiment to estimate the depth of objects in real stereo images is performed. The

results show good outputs with less error as compared with the real objects depths, which is attributed to good result of disparity (d) values obtained by the SABMA.

10. Suggestions for Future Works

The following suggestions can be suggested for future works:

1. Extending the current work for motion estimation and velocity extracting of the moving objects from many successive pairs of frames captured at different times.
2. Extending the current work to increase the field of view of the implemented stereo vision system by using more than two cameras.
3. It is suggested to implement the SABMA with another programming language (e.g., C++) to approach real time performance.

References:

- [1] Ambrosch Kristian and Humenberger Martin, "*Parameter Optimization of the SAD-IGMCT for Stereo Vision in RGB and HSV Color Spaces*", Article, Vienna, Austria, 2010.
- [2] Barrera Alejandra, "*Mobile Robots Navigation*", India: In-Teh, pp. (50-52), 2010.
- [3] Chinapirom Teerapat, et al, "*Stereoscopic Camera for Autonomous Mini-Robots Applied in KheperaSot League*", paper, Heinz Nixdorf Institute, University of Paderborn, Germany, 2007.
- [4] Hamzah R A, et al, "*An Aligned Epipolar Line for Stereo Images with Multiple Size ROI in Depth Maps for Computer Vision Application*", International Journal of Information and Education Technology, Vol. 1, No. 1, April 2011.
- [5] Goshtasby A. Ardeshir, "*r-D and r-D Image Registration*", USA: A John Wiley & Sons, Inc., Publication, pp. (197-210), 2005.
- [6] Khaleghi Bahsdor, et al, "*A New Miniaturized Embedded Stereo-Vision System (MESVS-I)*", Intelligent Sensing System Laboratory, E.C.E. Department, University of Windsor, 2008.
- [7] Florczyk Stefan, "*Robot Vision-Video-based Indoor Exploration with Autonomous and Mobile Robots*", Germany: WILEY-VCH Verlag GmbH & Co. KGaA, Weinheim, 2005.
- [8] www.middlebury.edu/stereo.
- [9] Scharstein Daniel, "*A Taxonomy and Evaluation of Dense Two-Frame Stereo Correspondence Algorithms*", International Journal of Computer Vision, USA: Kluwer Academic Publishers, 2002.