



Fast Stereo 3D Imaging Based on Random Speckle Projection and Its FPGA Implementation

Yuhao Shang^{1,2}, Wei Yin^{1,2}, Shijie Feng^{1,2}, Tianyang Tao^{1,2},
Qian Chen², and Chao Zuo^{1,2}(✉)

¹ Smart Computational Imaging (SCI) Laboratory, Nanjing University of Science and Technology, Nanjing, Jiangsu 210094, China
zuochao@njjust.edu.cn

² Jiangsu Key Laboratory of Spectral Imaging and Intelligent Sense, Nanjing University of Science and Technology, Nanjing, Jiangsu 210094, China

Abstract. In this paper, we propose a fast stereo 3D imaging technique based on random speckle projection and its FPGA implementation. Stereo vision, as a classic passive method for 3D shape measurement based on the multi-view geometric constraints, can realize the 3D reconstruction of the tested scene using a pair of images captured through the binocular cameras. In addition, some complicated matching techniques, such as graph cut and block matching, are used to obtain a global disparity map but it leads to massive computing overhead. To solve this problem, we developed a fast stereo vision system based on FPGA. Benefiting from the full parallel architecture of FPGA, the complete computational framework is based on a full pipeline design, that is, the storage and calculation of data are performed under the system clock to implement different works of stereo vision (including stereo rectify and stereo matching) at the same time, promoting calculation speed and measurement efficiency. In order to further improve the accuracy of 3D measurement, by introducing structured light illumination into the existing system, a projection system based on random speckle is designed where fast speckle projection and synchronous acquisition are realized on the FPGA hardware. Experimental results verify that our method can achieve high-speed and robust 3D shape measurement.

Keywords: 3D imaging · Speckle · FPGA

1 Introduction

The acquisition of 3D information has extremely high significance in the fields of AR, VR, military, industrial inspection, robotics, and aerospace. Among plenty of state-of-the-art methods of achieving the 3D reconstruction of the tested scene (including binocular stereo vision (BSV), TOF, and structured light illumination [1–4]), BSV, which is based on the principle of triangulation, has been proven to be one of the most promising techniques due to its inherent advantages of non-contact, high efficiency, and low cost. In a conventional measurement system based on BSV only consisting of the binocular cameras, a series works of stereo vision (including stereo rectify, stereo

matching, and left-right consistency check) are performed sequentially to get a global disparity map with high quality, but the measurement efficiency and speed of the system based on BSV are limited by the inherent instruction cycle delay within traditional computers to bring massive computing overhead, which leads to the limits on the application of BSV. Meanwhile, different from the traditional computer and specified integrated circuit, FPGA offers high flexibility and programmability to meet the stringent requirements of parallelism and internal bandwidth [5], that makes many FPGA-based stereo vision systems have good performance like real-time. Dunn et al. [6] propose a stereo vision system based on multi-chip FPGA with the time required for a stereo matching of 256×256 resolution images only 34 ms. FPGA-based stereo vision systems on custom boards have been developed by Nishihara [7] using the Laplacian-of-Gaussian Sign-Correlation algorithm and a stereo vision system with the correlation-based algorithm is designed by Ding et al. [8] with a high processing speed. Jin et al. [9] propose a fully pipelined stereo vision system providing a dense disparity image with additional sub-pixel accuracy in real-time. Hariyama [10] proposes a processor architecture for high-speed and reliable stereo matching based on adaptive window-size control of SAD computation.

However, due to the limited accuracy obtained using BSV, we bring speckle pattern projection into the 3D measurement system based on BSV in order to improve the accuracy of the measurement. Meanwhile, researchers have also done a lot of work on structured light illumination measurement, such as Pan et al. [11] propose an improved DIC combining a 2D digital image correlation technique with the projection of a random speckle pattern using a conventional LCD projector, and Axel et al. [12] propose a fast and accurate method with a correlation technique which takes only the area of one pixel into account, used to locate the homologous points. Furthermore, the projection pattern design is also a critical step for the structured light illumination measurement. Hua et al. [13] study the quality of the speckle pattern used in image correlation technique using the mean subset fluctuation parameter. A method for designing a composite pattern, in which the speckle pattern is embedded into the conventional phase-shifting fringe pattern with a simple and effective evaluation criterion for the correlation quality of the designed speckle pattern in order to improve the matching accuracy significantly, is proposed by Yin et al. [14]. In addition, the 3D imaging method based on speckle projection in our system is similar to the temporal correlation method Schaffer et al. [15] have proposed in this paper, and the 3D imaging system is realized on FPGAs.

2 FPGA-Based Data Transmission Framework

2.1 Image Data Transferred from Cameras to SDRAM

In our system, the digital image sensor used in the binocular cameras is MT9V034, which is a 1/3-Inch wide-VGA CMOS chip and can capture 10-bit grayscale image with the resolution of 752×480 at 60 Hz. In order to process and transmit the image data quickly and conveniently, we fine-tune the register configuration of MT9V034 to make our binocular cameras capture images with the resolution of 640×480 at 75 Hz

under a camera clock with the frequency of 27 MHz. In the acquisition process of images, the whole image data are transmitted pixel by pixel and row by row into FPGA. On the basis of MT9V034's datasheet, it is obvious that 224 clock intervals will be needed in the transmission between two adjacent rows of an image, which means that the data transmission of images by cameras can be equivalent to the transmission without interruption under a pixel clock of 20 MHz. Meanwhile, the SDRAM in our FPGA development board, which consists of two IS42S16320B chips with the memory size of $32\text{ M} \times 16\text{ bit}$, makes the system has enough storage for a pair of images and is able to refresh the stored data of images in real time to guarantee the efficiency of the data processing.

In order to make full use of the SDRAM, the read/write clock frequency is set to 100 MHz. However, the transmission of image data, between binocular cameras and SDRAM, needs to transfer data from 20 MHz clock domain to 100 MHz clock domain, which may lead to image data overflow or loss. To avoid this problem, two FIFOs (each with the size of $213 \times 8\text{ bit}$) are set between the two different clock domains, achieving the transmission of image data normally without losing or overflowing data, as shown in Fig. 1(a). In addition, according to the storage characteristics of SDRAM, the data stored in SDRAM will not disappear unless the power is turned off, SDRAM is cleared, or new data is read in and overwrites the old data, so the image data will be stored and read out by Ping-Pong switching.

2.2 Data Transferred from SDRAM to Post-processing Module

The clock signal used in the post-processing modules is set to 20 MHz, which aims at making the image post-process operation and the image acquisition equivalent to real-time handling of the original pixel data from the camera each clock. There is also a mismatch between the clock domain of SDRAM and the clock domain of the post-processing module, so two more FIFOs need to be set between the two modules to deal with the problems mentioned above, as shown in Fig. 1(a).

Therefore, the data transmission between the two modules is that SDRAM will receive a read request from the back-end FIFOs after storing the data of a pair of images, and then will write the image data to these FIFOs. The post-processing module deals with the data from back-end FIFOs when those back-end FIFOs have stored enough data. In addition, the important point is that when the post-processing module starts reading data from the two back-end FIFOs these FIFOs will be not cleared until the data of the entire image has been read, ensuring that the order of the transferred data is consistent with the order in which the image data is entered from cameras.

2.3 Ethernet Transfers Data to Computers

In this system, due to its wide transmission bandwidth that enables Ethernet to realize efficient data transmission, the processed data is transferred by Gigabit Ethernet using TCP/IP protocol to send UDP packets to computers under the clock signal with a frequency of 125 MHz. There into, except the preamble code, each UDP packet contains 1328 bytes, including a 42-byte header, a 2-byte image row number flag, 1280-byte image data that are the data of the same row in the pair of images, and a

check code of 4 bytes, as shown in Table 1. Similar to the four front-end and back-end FIFOs, two FIFOs, of which each at least stores data with the amount M (M could be obtained by Eq. (1)), are set between the post-processing and Ethernet modules.

Through the FPGA-based system, the computer handles the data packets received, and saves the data in the form of image, as shown in Fig. 1(b)–(c).

$$(44 + 640)/125 = (640 - M)/20 \tag{1}$$

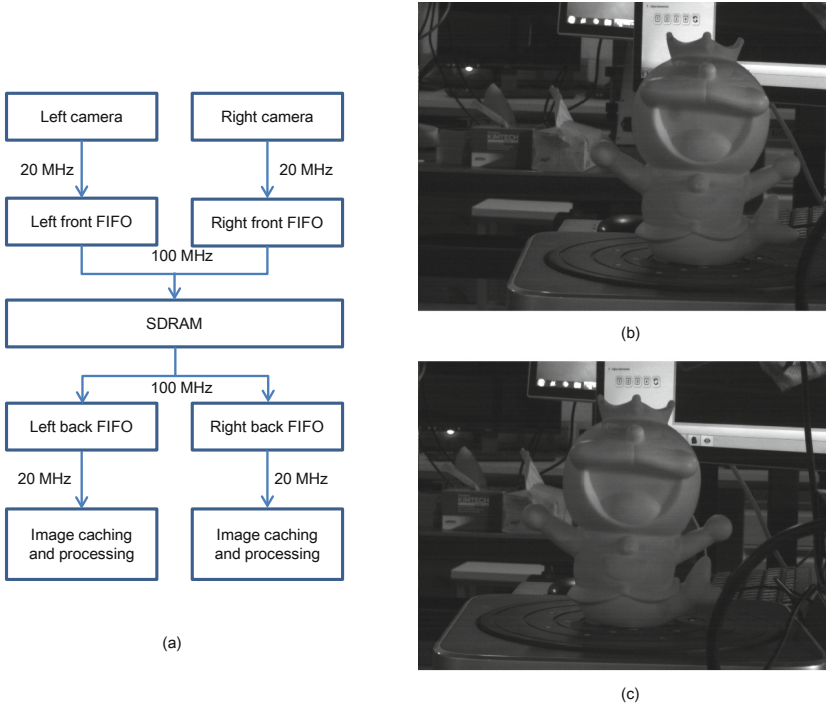


Fig. 1. (a) The data transmission between each module in this system. (b)–(c) There are two images captured by cameras. (b) The image is captured by left camera. (c) The image is captured by right camera.

Table 1. UDP packet contents (in byte).

Name	Preamble code	Ethernet header	IP header	UDP header
Length	8	14	20	8
Name	Image row number	Left image data	Right image data	Check code
Length	2	640	640	4

The key of the system implemented on FPGAs will be clearly stated in Sect. 4.

3 Image Processing

3.1 Speckle Pattern Design

In order to further improve the accuracy of 3D measurement, by introducing structured light illumination into the existing system, a projection system based on random speckle is designed where fast speckle projection and synchronous acquisition are realized on the FPGA hardware. Three-dimensional imaging technology based on optical structure is applied in more and more fields such as biomechanics, intelligent monitoring, robot navigation, industrial quality control, and human-computer interaction. In order to improve the three-dimensional imaging technology, researchers have carried out research on many factors that may have influence on the quality of the 3D imaging results, including the coding of projected structured light pattern. In order to accurately match the two images, a random encoding method will be proposed below.

Zhou et al. [16] propose a novel design method of color binary speckle pattern, and we will use this idea to design speckle patterns according to the relative position of the camera, tested object and the projector (see Fig. 2), as well as camera and projector parameters.

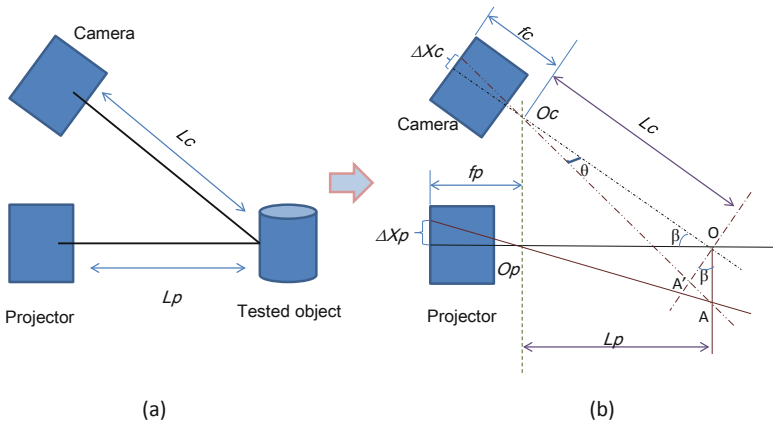


Fig. 2. Relationship between BWBSP design and system parameters. (a) Schematic binocular stereo measurement system. (b) Geometric relationship of the system.

According to the parameters of the projector and cameras, the relationship of the size between the speckle projected by the projector and the speckle captured by the camera can be known through Eq. (2), where $\Delta\delta_p$ is the pixel size of the projector (in mm), m_p is the width of the projected speckle, L_p is the projection distance of the projector, and f_p is the focal length of the projector, and $\Delta\delta_c$, m_c , and L_c are related parameters of the camera similarly.

$$m_p \Delta \delta_p = \frac{L_c^2 f_p}{L_p^2 f_c} (m_c \Delta \delta_c) \quad (2)$$

In addition, in order to satisfy the sampling theorem and obtain good image contrast, the allowable range of m_c value should belong to 3 to 5 [17]. In this design, the ratio of the camera distance L_c to the projection distance L_p is about 6/5; the size of a pixel of the camera sensor is 6.0 μm ; the focal length f_c is 12 mm; the size of the pixel of the projector is 7.6 μm ; and the focal length f_p of the projector is about 15 mm. If $m_c = 3$, then m_p is calculated as 4.

Speckle generation is generated by using MATLAB R2017a. The method is generated by referring to the method of Pan et al. [18]. The generated speckle pattern is shown in Fig. 3.

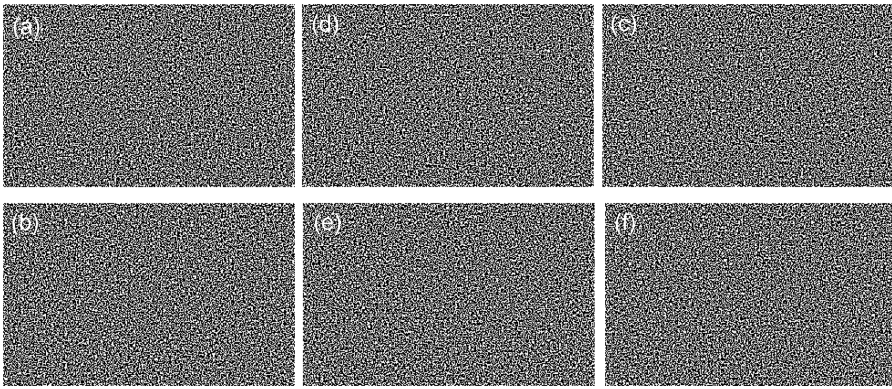


Fig. 3. Speckle patterns generated by MATLAB.

3.2 Image Rectification

Between binocular image acquisition and stereo matching, stereo rectification of binocular images needs to be implemented firstly which can narrow the matching range, thereby reducing the amount of computation generated by the matching algorithm, and indirectly speeding up the matching speed. Before image rectification, the work of camera calibration, which determines the parameters of the binocular cameras that facilitate the subsequent association of the world coordinate with the camera coordinate based on the relationship in the geometric models, would be done by Zhang's camera calibration method [19]. And then the polar line rectification of the images can be performed.

Camera Calibration and Inverse Mapping Pixel Calculation. The camera coordinate system and the world coordinate system can be converted to each other by rotation and translation transformations. In addition to the parameters required for the rotation and translation transformations, determining the relationship between the world coordinate system and the camera coordinate system also requires the precise focal length

and focus of the camera. The rotation matrix R and the translation matrix T are the external parameters of the camera, each having three parameters, while the focal length f and the focus c are the internal parameters of the camera, each of which consists of two parameters, so the calibration of a camera is to determine 10 variables of the camera. The conversion relationship between the world coordinate M_{world} and the camera coordinate M_{camera} of the object could be defined by Eq. (3). Then we can get the relationship between left and right cameras, as shown in Eq. (4), where R_{lr} is the rotation matrix and T_{lr} is the translation matrix between the binocular cameras, M_{right} and M_{left} are the coordinates of the object in the right and left camera coordinate systems, respectively.

$$M_{world} = R(M_{camera} - T) \tag{3}$$

$$M_{right} = R_{lr}(M_{left} - T_{lr}) \tag{4}$$

The rectification of the images is the calculation of the inverse mapping pixel. Before the value of pixels in the rectified images is calculated, the inverse mapping coefficient must be calculated before the subsequent calculation. However, the inverse mapping coefficient is the decimal part of the pixel coordinate, and the inverse pixel coordinates p are calculated by Eqs. (5) to (8).

$$p = [p_x \ p_y]^T \tag{5}$$

$$[PX \ PY \ PZ]^T = H'[x \ y \ 1]^T \tag{6}$$

$$H' = [H_x \ H_y \ H_z]^T = R^T K K^{-1} \tag{7}$$

$$p_x = f_x \frac{PX}{PZ} + c_x, \quad p_y = f_y \frac{PY}{PZ} + c_y \tag{8}$$

(f_x, f_y) and (c_x, c_y) are the focal length and principal point in x-axis and y-axis respectively, KK is the inherent parameters matrix, and R is the rotation matrix between right and left cameras.

It is easy to find that eight multiplications, two indispensable divisions, and eight additions are required in the process of calculating the inverse coordinates. However, the time required for calculating a multiplication in an FPGA is long, with the used resources large. In order to simplify the calculation of the inverse mapping pixel on the FPGA, the calculation process is reorganized, which makes Eqs. (4) to (7) changing into Eqs. (9) to (11).

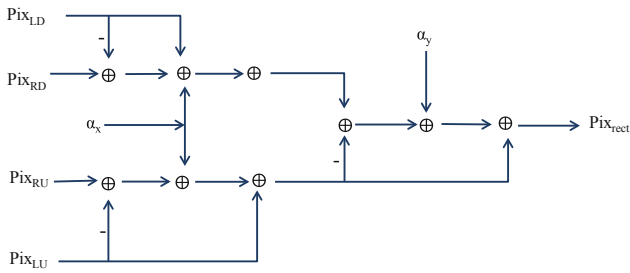
$$H'' = \begin{bmatrix} f_x & f_x & f_x \\ f_y & f_y & f_y \\ 1 & 1 & 1 \end{bmatrix} \cdot * H' \tag{9}$$

$$[X' \ Y' \ Z']^T = H''[x \ y \ 1]^T \tag{10}$$

$$p_x = \frac{X'}{Z'} + c_x, \quad p_y = \frac{Y'}{Z'} + c_y \tag{11}$$

Bilinear Interpolation Module. The equation for bilinear interpolation is given (see Eq. (12)). According to the principle of linear interpolation method, the integer part of the calculation result is the inverse mapping pixel coordinate, and the decimal part is the inverse mapping coefficient α_x and α_y , and the design of bilinear interpolation module [20] is shown in Fig. 4(a). The image correction result is shown in Fig. 4(b) and (c).

$$Pix_{rec} = (1 - \alpha_x)(1 - \alpha_y)Pix_{LU} + \alpha_x(1 - \alpha_y)Pix_{RU} + (1 - \alpha_x)\alpha_y Pix_{LD} + \alpha_x\alpha_y Pix_{RD} \tag{12}$$



(a)



(b)

(c)

Fig. 4. (a) The design of bilinear interpolation module; (b) Simulation results of the rectification algorithm on MATLAB; (c) The rectified images obtained from FPGA.

3.3 Image Matching

Statistical Pattern. Since the two images captured by cameras have been rectified in the previous work, the matching area of the images could be reduced from the whole frame to one line which means that potential wrong matching points are reduced, and it can be found that the disparity between the two images is within a certain range (about

100 or so) in Fig. 1 (b)–(c). Therefore, we can set the disparity value in the range from 80 to 144.

The temporal correlation method for statistical pattern is to project multiple speckle patterns on the tested object, and then the correlation between the two gray value sequences of the reference point and the point to be matched is obtained, among which the point with the largest correlation value is the matching point. According to the paper [15], the image related equation is shown in Eq. (13), where $g_i(p_j, t)$ denotes the gray value of the pixel p_j in the t -th image by camera i , and $\bar{g}_i(p_j)$ stands for the mean gray value of the pixel p_j in camera i over all the N images.

$$\rho(p_1, p_2) = \frac{\sum_{t=1}^N [g_1(p_1, t) - \bar{g}_1(p_1)] \cdot [g_2(p_2, t) - \bar{g}_2(p_2)]}{\sqrt{\sum_{t=1}^N [g_1(p_1, t) - \bar{g}_1(p_1)]^2 \cdot \sum_{t=1}^N [g_2(p_2, t) - \bar{g}_2(p_2)]^2}} \quad (13)$$

Algorithm Optimization. The number of speckle patterns projected in this experiment is small (only 6 different speckle patterns). Therefore, the quality of the result matched by Eq. (13) is very poor, the results have many wrong matching points, and the environmental interference can be found very serious, as shown in Fig. 5(a).

Therefore, the images are filtered by a simple filter that does not damage the image data before the matching is performed, and the background interference is filtered out (shown in Fig. 5(b)). After removing the background and then using Eq. (13) to correlate with 24 speckle patterns, it is found that the quality of the matching result is still poor as shown in Fig. 5(c), and by the method of spatial correlation we get the result shown in Fig. 5(d), of which the accuracy is not well. Therefore, we add the gray value of the neighboring pixels around this pixel for high quality matching, that is, the temporal correlation method is optimized by combining with spatial correlation method, with Eq. (13) changing into Eq. (14), where the radius of the added window centered on the pixel p_j to be matched is set to r , and $\bar{g}_i(p_j)$ means the average gray value of all the pixels in this window on all N images captured by camera i . And we can obtain disparity maps, as shown in Fig. 5(e)–(f). The result shown in Fig. 5(g) is the disparity map after left-right consistency check and occlusion area filling processing, and the result shown in Fig. 5(h) is the disparity map obtained by the median filter.

$$\rho(p_1, p_2) = \frac{\sum_{j=-r}^r \sum_{i=-r}^r \sum_{t=1}^N [g_1(p_1(i, j), t) - \bar{g}_1(p_1)] \cdot [g_2(p_2(i, j), t) - \bar{g}_2(p_2)]}{\sqrt{\sum_{j=-r}^r \sum_{i=-r}^r \sum_{t=1}^N [g_1(p_1(i, j), t) - \bar{g}_1(p_1)]^2 \cdot \sum_{j=-r}^r \sum_{i=-r}^r \sum_{t=1}^N [g_2(p_2(i, j), t) - \bar{g}_2(p_2)]^2}} \quad (14)$$

This optimized algorithm is realized on MATLAB with all experiments conducted on a 2.5 GHz Intel Core i7-6500U CPU, 8 GB of RAM and no GPU optimization, and the running time of processing single image with the resolution of 640×480 is about 8 min, that means it could not meet the requirement of real-time. Meanwhile, the operation using this proposed algorithm needs to get the data of all images with different speckle patterns at the same time and make correlation operation. So, on our FPGA development board, DE2-115, it is extremely difficult for this optimized

algorithm to implement due to the limited storage capacity, limited data bandwidth and the limited internal resources of FPGA chip. In our follow-up work, the problems mentioned above will be studied.

4 FPGA-Based Image Matching

Since temporal correlation algorithm is extremely difficult to implement on FPGAs, the stereo vision system is based on the Census transform matching method, by which the image data could be transformed into Census binary vectors to achieve the image matching, and the Census vectors are suit for the parallel computing mode of FPGA. Due to the space limitations, the details of the principle about the Census transform are omitted here. In order to better achieve the image matching module on FPGA, the pipeline design is introduced into the system, seen in Fig. 6(a), where the Hamming distances are stored in the registers with the corresponding disparity value after the exclusive OR operation is implemented on the Census binary vectors stored in shift registers, and the disparity value with the smallest Hamming distance is the desired.

The disparity maps obtained from the 3D imaging system built on the FPGA are shown in Fig. 6(b)–(e), where the result shown in Fig. 6(b) is obtained directly from FPGA board, the result shown in Fig. 6(c) is obtained after the left-right consistency check, and the result shown in Fig. 6(d) is the disparity map after the occlusion area filling processing with the result shown in Fig. 6(e) obtained after median filtering.

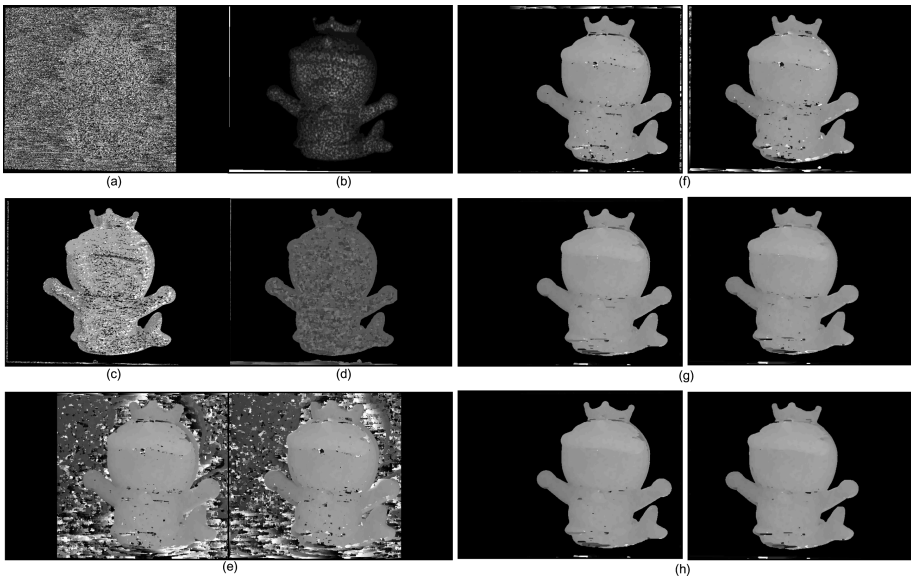


Fig. 5. Matching results. (a) The result is obtained by temporal correlation method; (b) The result is obtained by filtering out the background interference. (c) The result is obtained by temporal correlation method. (d) The result is obtained by spatial correlation method. (e)–(f) The results are obtained by the optimized method. (g)–(h) The results are obtained by left-right consistency check and the median filter.

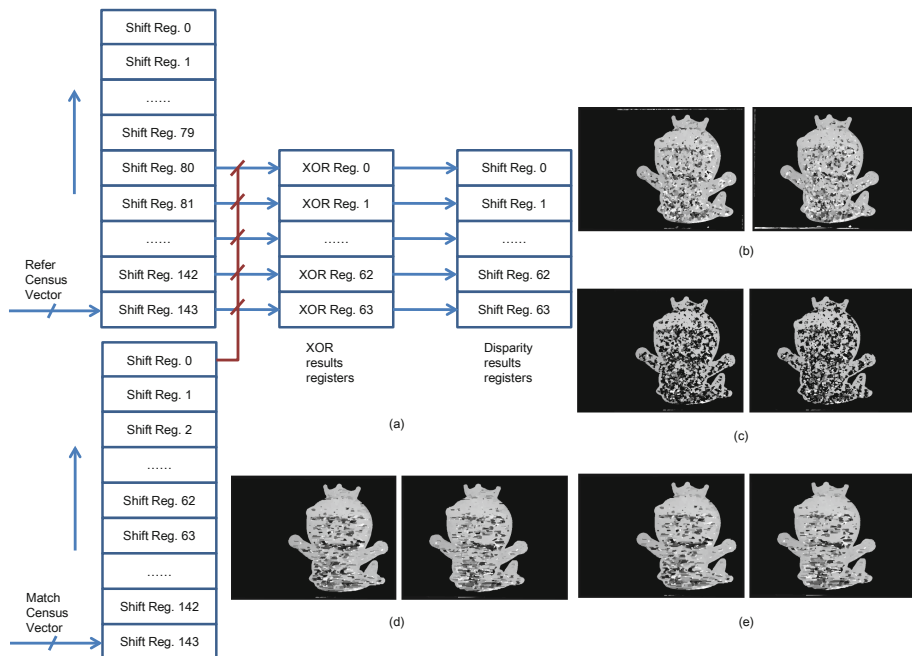


Fig. 6. (a) The pipeline design introduced on FPGA. (b)–(e) Disparity maps.

5 Conclusion

In this work, the 3D imaging technology based on random speckle projection is studied, and the hardware implementation of 3D imaging system is realized on FPGAs. The system can rectify and match the images acquired by the binocular cameras and use Gigabit Ethernet to transmit data. Then the correlation method is optimized by combining the spatial correlation method with temporal correlation method, by which the accuracy of the image matching results obtained with 6 speckle patterns is much higher than the results only by spatial correlation method or temporal correlation method (even more speckle patterns). Finally, the image process module of 3D imaging system realized by Census transform on FPGA could achieve real-time 3D measurement at 75 frames per second for the images with a resolution of 640×480 under a global 27 MHz clock signal.

References

1. Gorthi, S.S., Rastogi, P.: Fringe projection techniques: whither we are? *Opt. Laser Eng.* **48** (2), 133–140 (2010)
2. Feng, S., Zhang, L., Zuo, C., Tao, T., Chen, Q., Gu, G.: High dynamic range 3-D measurements with fringe projection profilometry: a review. *Meas. Sci. Technol.* **29**(12), 122001 (2018)

3. Zuo, C., Feng, S., Huang, L., Tao, T., Yin, W., Chen, Q.: Phase shifting algorithms for fringe projection profilometry: a review. *Opt. Laser Eng.* **109**, 23–59 (2018)
4. Yin, W., et al.: High-speed three-dimensional shape measurement using geometry-constraint-based number-theoretical phase unwrapping. *Opt. Laser Eng.* **115**, 21–31 (2019)
5. Zhang, L., Zhang, K., Chang, T.: Real-time high-definition stereo matching on FPGA. In: Proceedings of the ACM/SIGDA 19th International Symposium on Field Programmable Gate Arrays (2011)
6. Dunn, P., Corke, P.: Real-time stereopsis using FPGAs. In: Luk, W., Cheung, P.Y.K., Glesner, M. (eds.) *FPL 1997*. LNCS, vol. 1304, pp. 400–409. Springer, Heidelberg (1997). https://doi.org/10.1007/3-540-63465-7_245
7. Nishihara, H.K.: Real-time stereo- and motion-based figure ground discrimination and tracking using LOG sign correlation. In: Conference on Signals, Systems & Computers. IEEE (2002)
8. Ding, J., Du, X., Wang, X.: Improved real-time correlation-based FPGA stereo vision system. In: International Conference on Mechatronics & Automation. IEEE (2010)
9. Jin, S., Cho, J., Pham, X.D.: FPGA design and implementation of a real-time stereo vision system. *IEEE Trans. Circuits Syst. Video Technol.* **20**(1), 15–26 (2010)
10. Hariyama, M.: FPGA implementation of a stereo matching processor based on window-parallel-and-pixel-parallel architecture. In: Symposium on Circuits & Systems. IEEE (2005)
11. Pan, B., Xie, H., Gao, J.: Improved speckle projection profilometry for out-of-plane shape measurement. *Appl. Opt.* **47**(29), 5527–5533 (2008)
12. Axel, W., Holger, W., Richard, K.: Human face measurement by projecting bandlimited random patterns. *Opt. Express* **14**(17), 7692–7698 (2006)
13. Hua, T., Xie, H., Wang, S.: Evaluation of the quality of a speckle pattern in the digital image correlation method by mean subset fluctuation. *Opt. Laser Technol.* **43**(1), 9–13 (2011)
14. Yin, W., et al.: High-speed 3D shape measurement using the optimized composite fringe patterns and stereo-assisted structured light system. *Opt. Express* **27**, 2411–2431 (2019)
15. Schaffer, M., Marcus, G., Harendt, B.: Statistical patterns: an approach for high-speed and high-accuracy shape measurements. *Opt. Eng.* **53**(11), 112205 (2014)
16. Zhou, P., Zhu, J., Jing, H.: Optical 3-D surface reconstruction with color binary speckle pattern encoding. *Opt. Express* **26**(3), 3452 (2018)
17. Lionello, G., Cristofolini, L.: A practical approach to optimizing the preparation of speckle patterns for digital-image correlation. *Meas. Sci. Technol.* **25**(10), 107001 (2014)
18. Pan, B., Lu, Z., Xie, H.: Mean intensity gradient: an effective global parameter for quality assessment of the speckle patterns used in digital image correlation. *Opt. Lasers Eng.* **48**(4), 469–477 (2010)
19. Zhang, Z.: A flexible new technique for camera calibration. *IEEE Trans. Pattern Anal. Mach. Intell.* **22**(11), 1330–1334 (2000)
20. Ma, J., Yin, W., Zuo, C., Feng, S., Chen, Q.: Real-time binocular stereo vision system based on FPGA. In: ICOPEN 2018, p. 108271U (2018)