Paper:

# Improvement of Human Tracking in Stereoscopic Environment Using Subtraction Stereo with Shadow Detection

**Kenji Terabayashi\*, Yuma Hoshikawa\*\*, Alessandro Moro\*, and Kazunori Umeda\***

\*Department of Precision Mechanics, Faculty of Science and Engineering, Chuo University / CREST, JST
1-13-27 Kasuga, Bunkyo-ku, Tokyo 112-8551, Japan
E-mail: terabayashi@mech.chuo-u.ac.jp
\*\*Toyota Motor Corporation
375-1 Imazato Susono Sizuoka 410-1104, Japan

The combination of subtraction stereo with shadow detection we propose improves people tracking in stereoscopic environments. Subtraction stereo is a stereo matching method which is fast and robust for the correspondence problem – one of the most serious issues in computer vision – restricting the search range of matching to foreground regions. Shadow detection gives adequate foreground regions of tracked people by removing cast shadows. This leads to accurate three-dimensional measurement of positions in stereoscopic environment tracking. By focusing on disparity images obtained by subtraction stereo, we can detect people easily based on standard labeling. Objects can also be measured directly in size by subtraction stereo without geometric information about environments for tracking. This is important for installing the tracking system easily. To track multiple passing people, we use the extended Kalman filter to address the occlusion problem usually encountered in crowded environments. The proposed method is verified by experiments using unknown stereoscopic environments.

**Keywords:** subtraction stereo, human tracking, stereoscopic environment, shadow detection

## 1. Introduction

Tracking people in stereoscopic environments by measuring their three-dimensional (3D) positions is one of the challenges of autonomous systems. This has a variety of promising applications, such as surveillance systems, museum and store marketing, user guidance, and city planning. Key requirements for actual use are real-timeliness, portability, and easy tracking systems installation.

Many studies have focused on people tracking using three sensors: (i) a single camera (Collins et al. [1], Kagami et al. [2], Haga et al. [3], Cui et al. [4], Benezeth et al. [5]), (ii) a laser range finder (Cui et al. [4], Mozos et al. [6], Panangadan et al. [7]), and (iii) a stereo camera (Bahadori et al. [8], Ess et al. [9]). These studies usually treat flat environments such as the ground.

In measuring 3D tracking people information in stereo-

scopic environments, these sensors feature certain aspects. The stereo camera can obtain depth information based on calibrated intrinsic parameters [10] and directly measure size information for an object. This is useful for distinguishing whether the target is, for example, a human, a car, or a small animal. This sensor is appropriate for easy installation since it can obtain depth information without requiring information about measuring environment structure, for example, the ground or stereoscopic site, and the position and angle of the installed camera. The single camera needs, however, certain assumptions to measure distance. These assumptions include size information on a measuring object, environment structure, and camera position. This camera is thus not suitable for scalable use and easy installation. This point about scaling sometimes makes it difficult to distinguish a bird near the camera from a person, which is avoidable based on size information by using a stereo camera. A laser range finder, however, can measure distances directly and highly accurately, although it cannot obtain a range image in a moment due to sweeping the measuring environment. This is the main reason why laser range finders are usually used to track people on a flat plane instead of in a stereoscopic environment.

Terabayashi et al. [11] focused on stereoscopic environments as places for people tracking in 3D space because buildings often include parts of the environments, for example, stairs, atriums, and slopes. In this paper, subtraction stereo, proposed by Umeda et al. [12], is employed to make the tracking system real-time and easy to install. Subtraction stereo is faster and more robust than standard stereo when it comes to correspondence problems [10], since it restricts stereo matching to foreground regions alone.

The issue remaining in this paper is the adverse effect of shadows cast by tracked people in the 3D measurement of their locations. Cast shadows are often detected as foreground regions, especially in environments with strong illumination, which adversely affects 3D measurement using subtraction stereo.

The purpose of our study is to improve the people tracking in stereoscopic environments by reducing the cast shadow effect. For accurate 3D measurement, this pa-
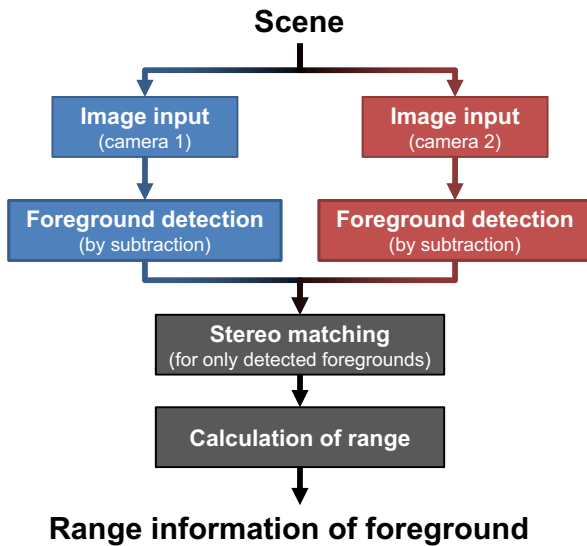
**Scene**

**Image input**
(camera 1)

**Image input**
(camera 2)

**Foreground detection**
(by subtraction)

**Foreground detection**
(by subtraction)

**Stereo matching**
(for only detected foregrounds)

**Calculation of range**

**Range information of foreground**

**Fig. 1.** Subtraction stereo flow.

per proposes combining subtraction stereo with shadow detection to remove cast shadows.

This paper is organized as follows: Section 2 proposes the subtraction stereo with shadow detection. Section 3 explains detection and 3D measurement used in people tracking. Section 4 focuses on people tracking with subtraction stereo. Section 5 details experimental results for people tracking in stereoscopic environments. Section 6 presents conclusions.

## 2. Subtraction Stereo with Shadow Detection

Subtraction stereo focuses on foreground information to increase the robustness of stereo matching and to make the matching process fast [12]. This method is appropriate for real-time people tracking in 3D environments. This section proposes a combination of subtraction stereo and shadow detection to improve the accuracy of 3D measurement in people tracking.

### 2.1. Subtraction Stereo Basics

**Figure 1** shows the basic subtraction stereo algorithm. In standard stereo vision, two images captured by left and right cameras are matched and disparities are obtained for each pixel. Subtraction stereo adds the step of extracting foreground regions from input images, then applies stereo matching to extracted foreground regions. One of the simplest ways to extract foreground regions is background subtraction.

Subtraction stereo is robust against the correspondence problem in stereo matching because search range of stereo matching is restricted to foreground regions. For the same reason, subtraction stereo can calculate disparity images more quickly than standard stereo.

**Figure 2** shows a comparative example of disparity images for subtraction stereo and standard stereo. **Fig. 2 (a)**

is a color image of input scene for the comparison. **Fig. 2 (b)** is an image subtracted from an image without people in it. **Figs. 2 (c)** and **(d)** are disparity images obtained by subtraction and standard stereo, in which pixel colors of represent disparities – e.g., green indicates a large disparity, i.e., a short distance. In contrast to the disparity image obtained by standard stereo matching (**Fig. 2 (d)**), disparities by subtraction stereo are only for foreground objects (**Fig. 2 (c)**).

### 2.2. Cast Shadow Removal

The cast shadow induced by tracked people in environments with strong illuminations adversely affects foreground regions for subtraction stereo. This reduces the accuracy of 3D position measurement in people tracking.

To improve people tracking, this paper proposes combining subtraction stereo and shadow detection to remove the cast shadow from foreground regions. Foreground regions refined by shadow detection are used for measuring the 3D positions of tracked people through subtraction stereo.

In shadow detection, color constancy criteria proposed by Moro et al. [14] are used to judge whether a pixel in an image is in shadow or not. When $I_{x,y,t}$ is the intensity value of an input image at a point with image coordinates $(x,y)$ at current time $t$ and $I'_{x,y}$ is the intensity value of a background image at the same point, the point can be determined for being in shadow by thresholding $\theta_{x,y,t}$ defined as

$$\theta_{x,y,t} = \begin{cases} \alpha\Psi_{x,y} + \beta\Lambda_{x,y} \\ \quad +(1-\alpha-\beta)\theta_{x,y,t-1} & \text{if } \dfrac{I_{x,y,t}}{\gamma} < I'_{x,y}; \\ \infty & \text{otherwise.} \end{cases} \quad (1)$$

Small $\theta_{x,y,t}$ corresponds to cast-shadow detection. In Eq. (1), $\Psi_{x,y}$ shows color constancy among nearby pixels, and $\Lambda_{x,y}$ shows color constancy within pixels located on image coordinates $(x,y)$. $\alpha$, $\beta$, and $\gamma$ are positive constant values determined empirically.

**Figure 3** shows the importance of shadow detection for segmenting people regions as foreground. **Figs. 3 (a)** and **(b)** are input scenes captured in indoor and outdoor environments. **Figs. 3 (c)** and **(d)** are shadow detection results overlapping input scenes with green pixels. In these figures, blue regions represent foreground refined through shadow detection, used for measuring 3D information by subtraction stereo. If shadow detection is not combined with subtraction stereo, the people's positions are adversely affected by the 3D positions of cast shadow shown in green in **Figs. 3 (c)** and **(d)**.

## 3. Human Detection and 3D Measurement

This section describes 3D measurement of people's positions for tracking in stereoscopic environments using subtraction stereo. In this paper, the calibrated parallel stereo camera is assumed to be fixed in place.
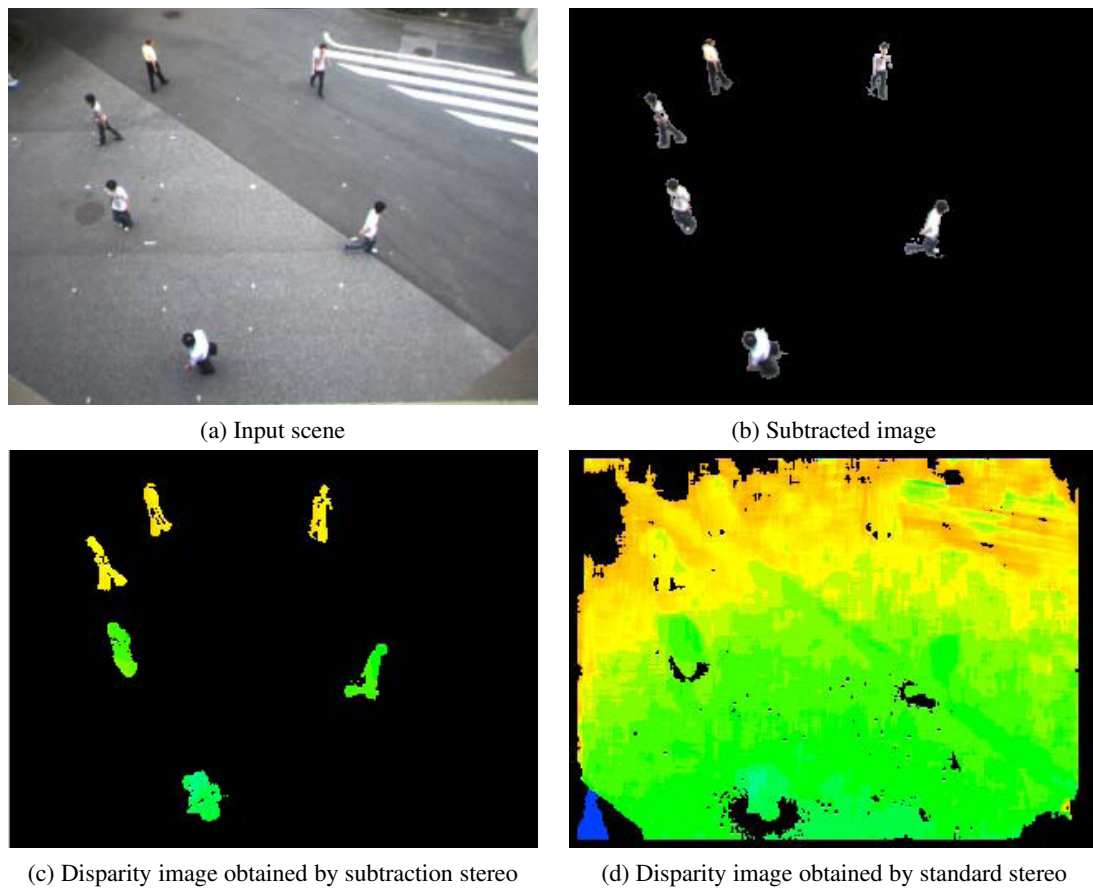
(a) Input scene

(b) Subtracted image

(c) Disparity image obtained by subtraction stereo

(d) Disparity image obtained by standard stereo

**Fig. 2.** Comparison of disparity images between subtraction and standard stereo.



(a) Input scene (indoor)

(b) Input scene (outdoor)

(c) Shadow detection (indoor)

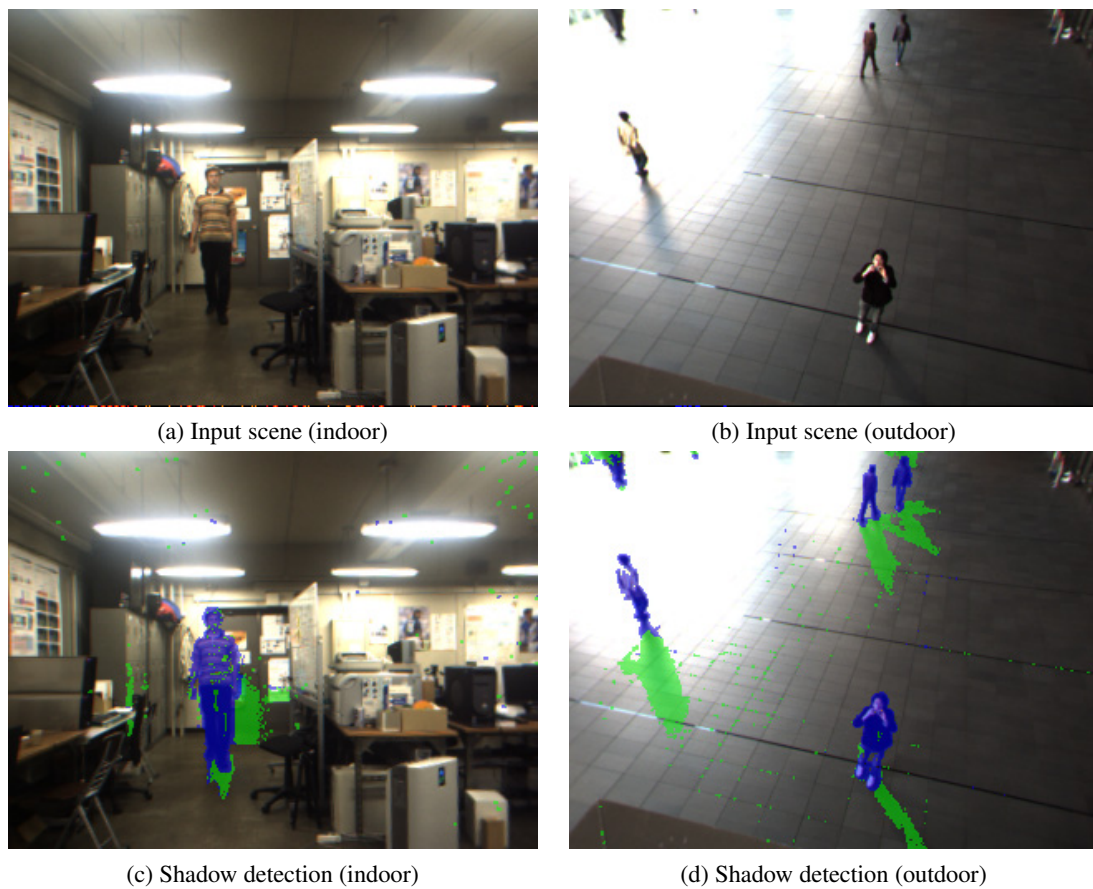(d) Shadow detection (outdoor)

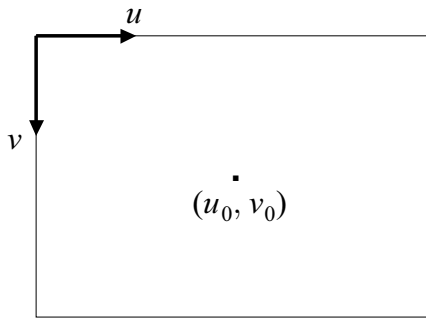**Fig. 3.** Importance of shadow detection for segmenting people regions in color images.

**Fig. 4.** Definition of image coordinates.

## 3.1. Detection by Labeling

Disparity images measured by subtraction stereo are originally restricted to foreground regions refined by shadow detection. This means that pedestrian group regions can be detected by standard labeling. To remove noise or to limit the size of objects, thresholding is applied to measured 3D size.

## 3.2. 3D Position Measurement

When the disparity of a point is given, the corresponding distance along the optical axis and the 3D position of the point are calculated. Hereafter, "distance" is used to mean the distance along the optical axis. Let disparity be $d$ and distance be $z$. Distance $z$ is calculated by

$$z = \frac{b f}{d p} \quad \ldots \ldots \ldots \ldots \ldots \ldots \quad (2)$$

where $b$ is baseline length, $f$ is the lens focal length, and $p$ is the width of each image pixel. Distance is inversely proportional to disparity. Measurement error of distance $z$ is proportional to the square of distance $z$.

The 3D position of a measured point is obtained from distance $z$ and image coordinates $(u, v)$ of a point in the image, as shown in **Fig. 4**. For the position of people, $z$ is used as the average distance of a labeled region, and $(u, v)$ is used as a point at the center of gravity (COG) of the region. Assuming no skew and that the aspect ratio of each pixel is 1, 3D position $\mathbf{x}$ is given as

$$\mathbf{x} = z \left[ \begin{array}{ccc} \frac{p}{f}(u - u_0) & \frac{p}{f}(v - v_0) & 1 \end{array} \right]^T \quad \ldots \quad (3)$$

where $(u_0, v_0)$ are the image coordinates of the image center.

## 4. People Tracking Using Measured 3D Positions

People tracking requires knowing the association of measured positions for each person frame by frame against measurement error and measurement blocked by occlusion. This paper uses the extended Kalman filter (EKF) to estimate the accurate position of each person.

State variable $\mathbf{x}$ of the EKF is defined as

$$\mathbf{x} = \left[ \begin{array}{cccccc} x & \dot{x} & y & \dot{y} & z & \dot{z} \end{array} \right]^T \quad \ldots \ldots \ldots \quad (4)$$

where $(x, y, z)$ and $(\dot{x}, \dot{y}, \dot{z})$ represent person position and velocity. Predicting $\mathbf{x}$ at the next time is calculated using current state $\mathbf{x}_t$ as

$$\mathbf{x}_{t+1} = \Phi \mathbf{x}_t + \omega \quad \ldots \ldots \ldots \ldots \ldots \quad (5)$$

where $\Phi$ is state transition matrix and $\omega$ is process noise. Based on the assumption that the velocity of a person is constant, the state transition matrix $\Phi$ is defined as

$$\Phi = \left[ \begin{array}{cccccc} 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right] . \quad \ldots \ldots \ldots \quad (6)$$

The assumption is adequate when the stereo camera frame rate is sufficiently high compared to walking speed. Measurement variable $\mathbf{z}$ of the EKF is given by

$$\mathbf{z} = \left[ \begin{array}{ccc} u & v & d \end{array} \right]^T \quad \ldots \ldots \ldots \ldots \ldots \quad (7)$$

where $u$ and $v$ are the image coordinates of a person in an image and $d$ is disparity. The relationship between state variable $\mathbf{x}$ and measurement variable $\mathbf{z}$ is as follows:

$$\mathbf{z}_t = f(\mathbf{x}_t) + \mathbf{v} \quad \ldots \ldots \ldots \ldots \ldots \ldots \quad (8)$$

$$f(\mathbf{x}_t) = \left[ \begin{array}{ccc} \frac{x_t f}{z_t p} + u_0 & \frac{y_t f}{z_t p} + v_0 & \frac{b f}{z_t p} \end{array} \right]^T \quad \ldots \quad (9)$$

where $f$ and $b$ are the camera focal length and baseline length, and $\mathbf{v}$ is measurement noise. As stated, the EKF estimates the position of each person in each frame using the following Jacobian matrix $\mathbf{J}$:

$$
\begin{aligned}
\mathbf{J} &= \frac{\partial f(\mathbf{x}_t)}{\partial \mathbf{x}_t} \\
&= \left[ \begin{array}{cccccc} \frac{f}{z_t p} & 0 & 0 & 0 & -\frac{x_t f}{z_t^2 p} & 0 \\ 0 & 0 & \frac{f}{z_t p} & 0 & -\frac{y_t f}{z_t^2 p} & 0 \\ -\frac{x_t f}{z_t^2 p} & 0 & -\frac{y_t f}{z_t^2 p} & 0 & -\frac{b f}{z_t^2 p} & 0 \end{array} \right] \quad (10)
\end{aligned}
$$

Each measured point is associated with a particular person when the distance between the estimated position and the measurement point is less than a threshold value. When there is no unique association with a person, the EKF simply predicts the person's position based on the Eq. (5).

## 5. Experiments

Tracking people in a stereoscopic environment is done using subtraction stereo with shadow detection without giving information about the measuring environment structure and stereo camera configuration. Use in unknown stereoscopic environments is important for easy

**Fig. 5.** Experimental scene in a stereoscopic atrium environment.

installation of the tracking system. Subtraction stereo is used with shadow detection.

### 5.1. Implementation of Subtraction Stereo with Shadow Detection

The subtraction stereo algorithm is implemented with a commercially available stereo camera (Point Grey Research Bumblebee2, color, $b = 120$ mm, $f = 3.8$ mm, $b = 14.8\ \mu$m) [13]. The image size is set to $320 \times 240$ pixels.

The shadow detection described in Section 2.2 is also implemented to improve the 3D measurement of people's positions in stereoscopic environments with strong illumination. To obtain a disparity image using subtraction stereo, stereo matching function in the Bumblebee2 library is applied to refined subtraction images from left and right cameras.

The computational speed for calculating disparity images by subtraction stereo is 30 fps, obtained using an Intel Core 2 Quad CPU (2.83 GHz) with 4 GB RAM. The accuracy of distance measured is discussed elsewhere (Umeda et al. [12]).

### 5.2. People Tracking in a Unknown Stereoscopic Atrium Environment

People passing on two floors in a unknown stereoscopic atrium environment were observed by a fixed stereo camera for tracking, e.g., as shown in **Fig. 5**. Note in the figure that the camera sees the two floors at the same time. In experiments, no layout information is given to measure 3D people positions.

Screenshots of tracking results are shown in **Fig. 6**. In this figure, colored rectangles are bounding boxes of people detected for tracking, and surrounding ID numbers are automatically assigned frame by frame based on EKF trackers. Note that people on different floors are tracked with each person is identical even after occlusion.

**Figure 7** shows the corresponding measured 3D positions of people in the stereoscopic environment. **Figs. 7 (a)** and **(b)** are the overhead and front views. In these

figures, colored markers indicate the locations of each tracked person. Results show that 3D positions of tracked people were measured correctly according to corridor structures in the atrium such as the corner and the two floors.

Without giving layout information, 3D positions of people tracked in an unknown stereoscopic environment were obtained by subtraction stereo with shadow detection. This is essential for ensuring that the tracking system is easy to install in stereoscopic environments. This is one advantage of using a stereo camera instead of single cameras and laser range finders.

The issue remaining for the tracking system is occlusion adversely affecting 3D measurements. An example is seen in the part of a person with ID_15 in **Fig. 6 (c)**. In this figure, the colored rectangle detected as the person is smaller than the actual size in the image. People detection with lacking parts makes 3D measurement inadequate, for example, the $Z$ coordinate variation shown in **Fig. 7 (b)**. To resolve this issue, the following two possibilities can be considered: (i) changing state variables of EKF trackers from the COG to the four corners of human regions; (ii) combining subtraction-based detection with appearance-based human detection (Zhe and Davis [15], Arie et al. [16]).

## 6. Conclusions

The major contribution of this paper is its 3D tracking in unknown stereoscopic environments using subtraction stereo with shadow detection. Shadow detection is applied to improve 3D measurement of people's positions by removing cast shadows from foreground regions for subtraction stereo. Experimental results show that people passing in an unknown 3D environment could be tracked and measured correctly without requiring environment layout information and camera configuration. This point is vital for easy system installation in various environments including stereoscopic structures.

**References:**
[1] R. Collins, A. Lipton, T. Kanade, H. Fujiyoshi, D. Duggins, Y. Tsin, D. Tolliver, N. Enomoto, and O. Hasegawa, "A system for video surveillance and monitoring," Technical Report CMU-RI-TR-00-12, Robotics Institute, Carnegie Mellon University, 2000.
[2] S. Kagami, K. Okada, M. Inaba, and H. Inoue, "Real-time 3D depth flow generation and its application to track to walking human being," Proceedings of 15th International Conference on Pattern Recognition (ICPR2000), Vol.4, pp. 197-200, 2000.
[3] T. Haga, K. Sumi, and Y. Yagi, "Human detection in outdoor scene using spatio-temporal motion analysis," Proceedings of 17th International Conference on Pattern Recognition (ICPR2004), Vol.4, pp. 331-334, 2004.
[4] J. Cui, H. Zha, H. Zhao, and R. Shibasaki, "Tracking multiple people using laser and vision," Proceedings of 2005 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS2005), pp. 2116-2121, 2005.
[5] Y. Benezeth, B. Emile, H. Laurent, and C. Rosenberger, "Vision-based system for human detection and tracking in indoor environment," International Journal of Social Robotics, Vol.2, No.1, pp. 41-52, 2010.
[6] O. M. Mozos, R. Kurazume, and T. Hasegawa, "Multi-part people detection using 2D range data," International Journal of Social Robotics, Vol.2, No.1, pp. 31-40, 2010.
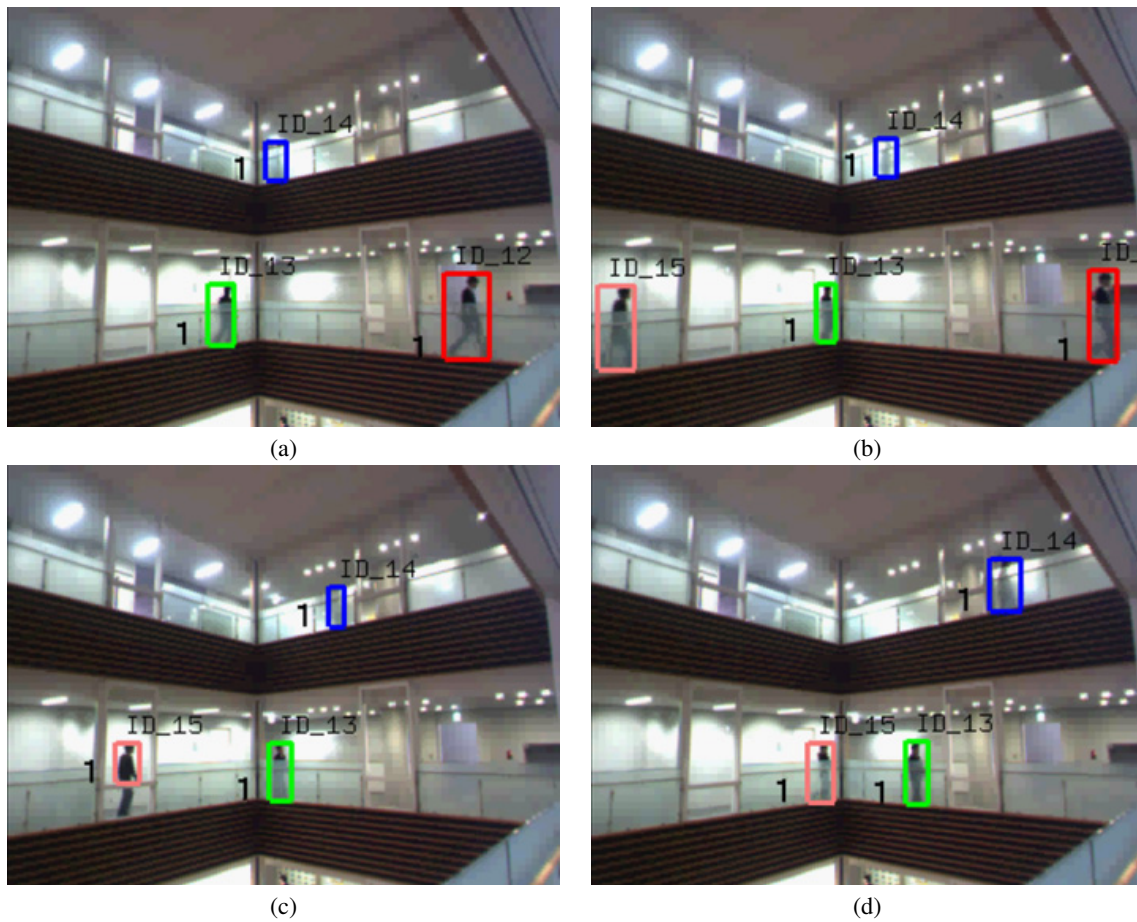
(a)

(b)

(c)

(d)

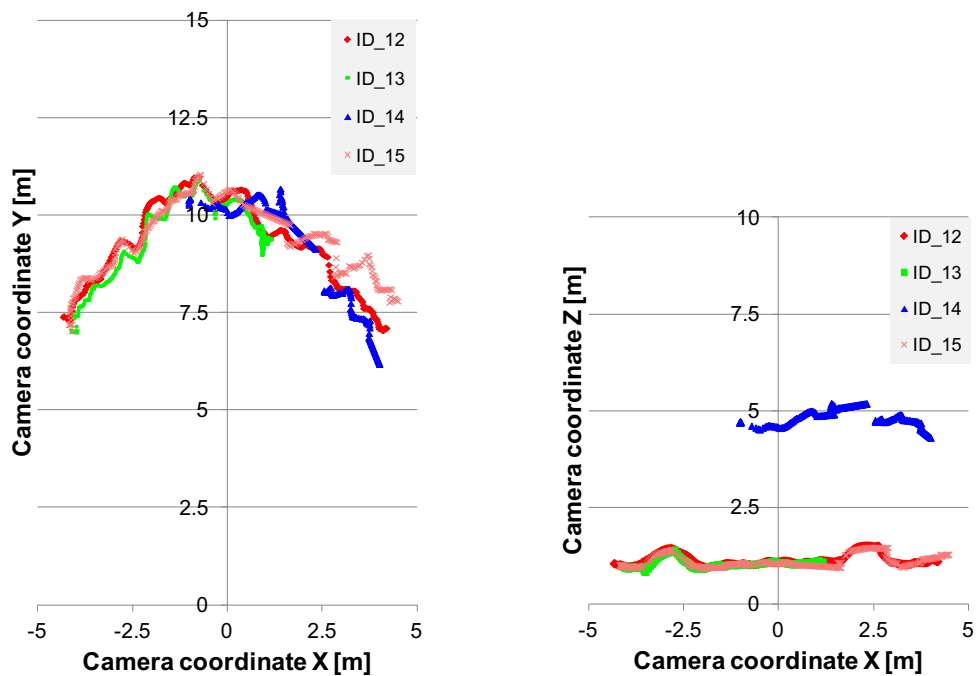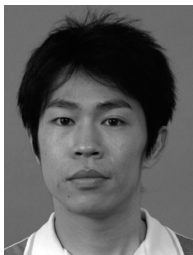**Fig. 6.** Experimental results of people tracking in a stereoscopic atrium environment.



(a) Measured positions in *X-Y* plane (overhead view)

(b) Measured positions in *X-Z* plane (front view)

**Fig. 7.** Experimental results of people tracking in a stereoscopic atrium environment.

Terabayashi, K. et al.

[7] A. Panangadan, M. Matarić, and G. S. Sukhatme, "Tracking and modeling of human activity using laser rangefinders," International Journal of Social Robotics, Vol.2, No.1, pp. 95-107, 2010.

[8] S. Bahadori, L. Iocchi, G. R. Leone, D. Nardi, and L. Scozzafava, "Real-time people localization and tracking through fixed stereo vision," Applied Intelligence, Vol.26, No.2, pp. 83-97, 2007.

[9] A. Ess, B. Leibe, K. Schindler, and L. van Gool, "Robust multiperson tracking from a mobile platform," IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol.31, No.10, pp. 1831-1846, 2009.

[10] R. Hartley and A. Zisserman, "Multiple View Geometry in Computer Vision," Cambridge Univ. Press, 2000.

[11] K. Terabayashi, Y. Hoshikawa, Y. Hashimoto, and K. Umeda, "Real-time human tracking in stereoscopic environments using subtraction Stereo," Proceedings of 4th International Asia Symposium on Mechatronics (AISM2010), pp. 97-104, 2010.

[12] K. Umeda, T. Nakanishi, Y. Hashimoto, K. Irie, and K. Terabayashi, "Subtraction stereo – a stereo camera system that focuses on moving regions –," Proceedings of SPIE 3D Imaging Metrology, Vol.7239, 2009.

[13] Point Grey Research, "http://www.ptgrey.com/"

[14] A. Moro, K. Terabayashi, K. Umeda, and E. Mumolo, "Auto-adaptive threshold and shadow detection approaches for pedestrians detection," Proceedings of Asian Workshop on Sensing and Visualization of City-Human Interaction (AWSVCI2009), pp. 9-12, 2009.

[15] L. Zhe and L. S. Davis, "Shape-Based Human Detection and Segmentation via Hierarchical Part-Template Matching," IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol.32, No.4, pp. 604-618, 2010.

[16] M. Arie, A. Moro, Y. Hoshikawa, T. Ubukata, K. Terabayashi, and K. Umeda, "Fast and Stable Human Detection Using Multiple Classifiers Based on Subtraction Stereo with HOG Features," Proceedings of 2011 IEEE International Conference on Robotics and Automation (ICRA2011), pp. 868-873, 2011.

**Name:**
Yuma Hoshikawa

**Affiliation:**
Toyota Motor Corporation

**Address:**
375-1 Imazato Susono Sizuoka 410-1104, Japan
**Brief Biographical History:**
● 2011 Received M. Eng. in precision engineering from Chuo University
● 2011- Joined Toyota Motor Corporation
**Main Works:**
● Y. Hoshikawa, K. Terabayashi, and K. Umeda, "Human Tracking Using Subtraction Stereo and Color Information," Proc. of Asian Workshop on Sensing and Visualization of City-Human Interaction, pp. 5-8, 2009.
● Y. Hoshikawa, Y. Hashimoto, A. Moro, K. Terabayashi, and K. Umeda, "Tracking of Human Groups Using Subtraction Stereo," SICE Journal of Control, Measurement, and System Integration, Vol.4, No.3, pp. 214-220, 2011.
**Membership in Academic Societies:**
● The Japan Society for Precision Engineering (JSPE)

**Name:**
Kenji Terabayashi

**Affiliation:**
Assistant Professor, Department of Precision Mechanics, Chuo University

**Address:**
1-13-27 Kasuga, Bunkyo-ku, Tokyo 112-8551, Japan
**Brief Biographical History:**
● 2004 Received M. Eng. degree in systems and information engineering from Hokkaido University
● 2008 Received Ph.D. degree in precision engineering from the University of Tokyo
● 2008- Assistant Professor at Chuo University
**Main Works:**
● K. Terabayashi, H. Mitsumoto, T. Morita, Y. Aragaki, N. Shimomura, and K. Umeda, "Measurement of Three Dimensional Environment with a Fish-eye Camera Based on Structure From Motion - Error Analysis," Journal of Robotics and Mechatronics, Vol.21, No.6, pp. 680-688, 2009.
● K. Terabayashi, N. Miyata, K. Umeda, and J. Ota, "Role of Pre-Operation in Experiencing Differently Sized Hands," Journal of Advanced Computational Intelligence and Intelligent Informatics, Vol.14, No.7, pp. 793-801, 2010.
**Membership in Academic Societies:**
● The Robotics Society of Japan (RSJ)
● The Virtual Reality Society of Japan (VRSJ)
● The Japan Society for Precision Engineering (JSPE)
● The Japan Society of Mechanical Engineers (JSME)
● The Institute of Electrical Engineers of Japan (IEEJ)
● The Institute of Image Electronics Engineers of Japan (IIEEJ)
● The Institute of Electrical and Electronics Engineers (IEEE)

**Name:**
Alessandro Moro

**Affiliation:**
Postdoctral Research Fellow, Chuo University

**Address:**
1-13-27 Kasuga, Bunkyo-ku, Tokyo 112-8551, Japan
**Brief Biographical History:**
● 2011 Received Ph.D. degree in computer science from University of Trieste
**Main Works:**
● A. Moro, K. Terabayashi, K. Umeda, and E. Mumolo, "A Framework for the Detection and Interaction with Pedestrian and Objects in an Unknown Environment," Proc. of Int. Conf. on Network Sensing Systems, pp. 257-260, 2010.
● A. Moro, K. Terabayashi, and K. Umeda, "Detection of Moving Objects with Removal of Cast Shadows and Periodic Changes Using Stereo Vision," Proc. of Int. Conf. on Pattern Recognition, pp. 328-331, 2010.

**Name:**
Kazunori Umeda

**Affiliation:**
Professor, Department of Precision Mechanics,
Chuo University

**Address:**
1-13-27 Kasuga, Bunkyo-ku, Tokyo 112-8551, Japan

**Brief Biographical History:**
● 1994 Received Ph.D. degree in precision machinery engineering from
the University of Tokyo
● 1994- Lecturer at Chuo University
● 2003-2004 Visiting Worker at National Research Council of Canada
● 2006- Professor at Chuo University

**Main Works:**
● N. Hikosaka, K. Watanabe, and K. Umeda, "Development of Obstacle
Recognition System of Humanoids Using Relative Disparity Maps from
Small Range Image Sensors," J. Robotics and Mechatronics, Vol.19, No.3,
pp. 290-297, 2007.
● M. Tateishi, H. Ishiyama, and K. Umeda, "A 200Hz Small Range Image
Sensor Using a Multi-Spot Laser Projector," Proc. of IEEE Int. Conf. on
Robotics and Automation, pp. 3022-3027, 2008.

**Membership in Academic Societies:**
● The Robotics Society of Japan (RSJ)
● The Japan Society for Precision Engineering (JSPE)
● The Japan Society of Mechanical Engineers (JSME)
● The Horological Institute of Japan (HIJ)
● The Institute of Electronics, Information and Communication Engineers
(IEICE)
● Information Processing Society of Japan (IPSJ)
● The Institute of Electrical and Electronics Engineers (IEEE)