

Fusion of stereo and optical flow data using occupancy grids

Christophe Braillon¹, Kane Usher², Cédric Pradalier², James L. Crowley¹ and Christian Laugier¹

¹Laboratoire GRAVIR
INRIA Rhône-Alpes
655 avenue de l'Europe
38334 Saint Ismier Cedex, France

²CSIRO ICT Centre
Autonomous Systems Lab
1 Technology court
Pullenvale QLD 4069, Australia

Email: firstname.lastname@inrialpes.fr

Email: firstname.lastname@csiro.au

Abstract—In this paper, we propose a real-time method to detect obstacles using theoretical models of the ground plane, first in a 3D point cloud given by a stereo camera, and then in an optical flow field given by one of the stereo pair's camera.

The idea of our method is to combine two partial occupancy grids from both sensor modalities with an occupancy grid framework. The two methods do not have the same range, precision and resolution. For example, the stereo method is precise for close objects but cannot see further than 7 m (with our lenses), while the optical flow method can see considerably further but has lower accuracy.

Experiments that have been carried on the CyCab mobile robot and on a tractor demonstrate that we can combine the advantages of both algorithms to build local occupancy grids from incomplete data (optical flow from a monocular camera cannot give depth information without time integration).

I. INTRODUCTION

This work takes place in the general context of mobile robots navigating in open and dynamic environments. Computer vision for ITS (Intelligent Transport Systems) is an active research area [6]. One of the key issues of ITS is the ability to avoid obstacles. This requires a method to perceive them.

In this article we address the problem of obstacle sensing through their motion (optical flow) in an image sequence. The perceived motion can be caused either by the obstacle itself or by the motion of the camera (which is the motion of the robot in the case of a camera fixed on it). We also use a stereo camera to improve the range (for short sight) of the resulting sensor.

Many methods have been developed to find moving objects in an image sequence. Most of them use a fixed camera and use background subtraction (for example in [13] and [5]).

Recently, in [15], a new approach to obstacle avoidance has been developed, based on ground detection by finding planes in images. The weak point of this method is that the robot must be in a static environment.

Model based approaches using ego-motion have been demonstrated in [9], [14]. The first one detects the ground plane by virtually rotating the camera and visually estimating the ego-motion. The second one uses dense stereo and optical flow to find moving objects and robot ego-motion. These two

methods have a large computational cost as several successive calculations (stereo, optical flow, ego-motion, ...) are required.

In this paper, we demonstrate that without knowing the motion of the camera, we can model the motion of the ground plane and determine the location of the obstacles. Moreover we show that the stereo data improves the accuracy for short range obstacle detection.

One key point in this method is that we do not compute explicitly the optical flow of the image at any time. Optical flow computation is very expensive in terms of CPU time, is inaccurate and sensitive to noise. In general we can see in the survey led by Barron et al. ([1]) that the accuracy of optical flow computation is linked to the computational cost. We model the expected optical flow (which is easy and quick to compute) to get rid of the inherent noise and the time consuming optical flow step. As a consequence we are able to demonstrate robust and real-time obstacle detection.

II. MODEL-BASED OBSTACLE DETECTION (STEREO AND OPTICAL FLOW)

The two algorithms we will use in the next two parts are model-based approaches. In a first step, we model our expected observation in the sensor space (in this article we focus on the ground plane model).

Let $n \in \mathbb{N}$ be the dimension of the observation space and $m \in \mathbb{N}$ a number of parameters. The model $\mathcal{F}_{n,m,p}$ with m parameters is a function $(\mathbb{R}^n \times \mathbb{R}^m) \rightarrow \mathbb{R}^p$ and is defined by the relation:

$$\mathcal{F}_{n,m,p}(Z, P) = 0, P \in \mathbb{R}^m \text{ and } Z \in \mathbb{R}^n$$

Given a model $\mathcal{F}_{n,m,p}$ of our observations we try to extract the parameter set $P \in \mathbb{R}^m$ from a set of observed data $(Z_i)_{i=1 \dots k}$ by minimising the error in the system:

$$\mathcal{F}_{n,m,p}(Z_i, P) = 0, \forall i \in \{1 \dots k\}$$

This minimisation can be done by performing a Least Mean Squares (LMS) minimisation. However, for better outlier rejection, here we use a Least Median Squares technique with the minimisation performed by the Nelder-Mead Simplex search [12].

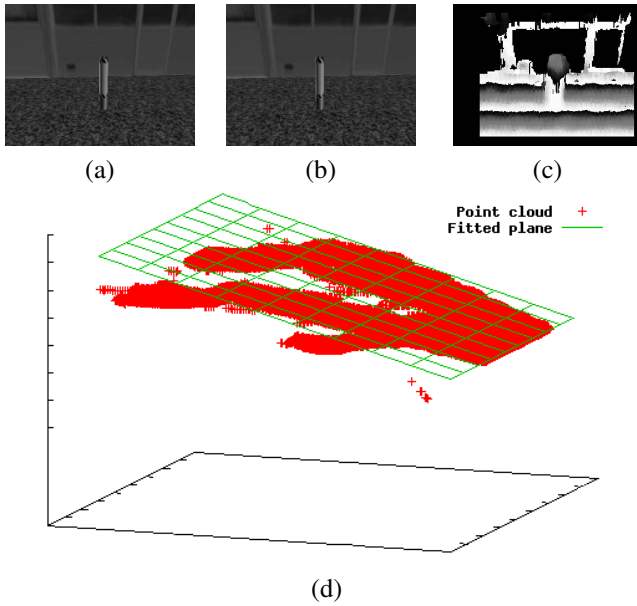


Fig. 1. (a) and (b) are respectively the left and right images from a synthetic image sequence (c) is the disparity map computed with SVS software, and (d) is the corresponding stereo point cloud, also shown is the plane fitted to the data.

Once the parameters are retrieved, we can cluster the observation in two sets (observations that match the model and observations that do not).

In the next part, we will use a ground plane model in both optical flow and 3D world spaces. In that case, the parameter sets will give us the position of the camera (with respect to the ground plane) and its ego-motion. The observation will be clustered in two sets: the ground plane and the obstacles.

A. Stereo obstacle detection

For obstacle detection using the stereo camera, we take the point cloud data generated from the stereo images, find the dominant plane in the point cloud (which should be the ground-plane) using the Least Median Squares method, and then build an occupancy grid based upon the distance of each point from the dominant plane.

The observation space using the point cloud data is \mathbb{R}^3 . Working in Cartesian space, the ground plane can be described by a set of four parameters (\mathbb{R}^4). The model is then expressed as:

$$\mathcal{F}_{3,4,1}(Z, P) = p_1x + p_2y + p_3z + p_4$$

where $Z = (x, y, z)$ and $P = (p_1, p_2, p_3, p_4)$. Figure II-A shows an example of the stereo point cloud data together with the plane fitted to the data.

Having found the ground plane, we can now estimate the camera height, roll and tilt, and compare these to the 'expected' values (from knowledge of the camera mounting position). We can also populate the occupancy grid using the idea that points not on the ground plane (at least within a tolerance) must belong to an obstacle. That is, we can calculate

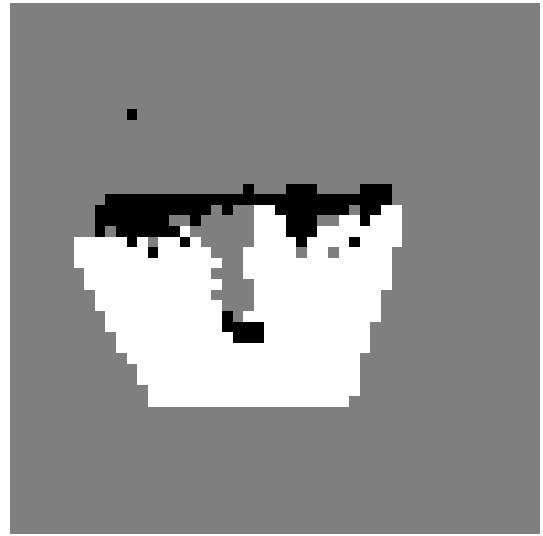


Fig. 2. Example of an occupancy grid generated from the stereo data and plane fitting process. The viewpoint is from the bottom centre of the image.

the distance of each point from the ground plane using the following equation:

$$e_i = p_1x_i + p_2y_i + p_3z_i + p_4$$

where the subscript i denotes the i th point in the cloud. If this distance exceeds a threshold, then the point contributes to the evidence that the ground plane cell to which the point belongs to contains an obstacle. Otherwise, if the distance is below the threshold, the point contributes to the evidence that the cell is free. Figure 2 illustrates an example occupancy grid generated using this method. In this figure, an 'empty' cell is represented by white, an 'occupied' cell is represented by black, and 'unknown' cells are represented by grey. The viewpoint is from the bottom of the image. Note the difference to, for example, a scanning laser generated occupancy grid, which physically can't 'see' behind objects.

B. Optical flow obstacle detection

By definition, an optical flow field is a vector field that describes the velocity of pixels in an image sequence. The first step of our method is the modelling of the optical flow field for our camera.

This model is based on the classical pinhole camera model, that is to say, we neglect the distortion due to the lens. We also assume that there is no skew factor. We will see in the experimental results that these two assumptions are valid.

1) *Naive first approach*: The parametrisation of our model can be found in [11], it says that only 8 parameters (parameter space is \mathbb{R}^8 and $P = (p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8)$) are needed to fully describe the visual motion of a plane. We call $Z = (u, v, f_u, f_v)$ an observation of an optical flow vector $\vec{f} = (f_u, f_v)$ at pixel (u, v) . Therefore the observation space is \mathbb{R}^4 . We can write the new model as follows:

$$\mathcal{G}_{4,8,2}(Z, P) = \begin{pmatrix} (p_1 - f_u) + p_2u + p_3v + p_4u^2 + p_5uv \\ (p_6 - f_v) + p_7u + p_8v + p_4uv + p_5v^2 \end{pmatrix}$$

Using this method for optical flow requires a good accuracy on the ground plane to evaluate the parameters of the ground plane. Indeed, the part of the flow field we want to model is the ground plane. Therefore we need a good accuracy on its optical flow. Moreover we want to respect our real-time constraint. Thus we need a method that perform accurate optical flow computation in real-time.

We studied all the characteristics of various optical flow method described in [1] but no method was really appropriate (either inaccurate or slow). We used brand new optical flow computation methods developed by Bruhn, Weickert *et al.* (in [2], [3], [4], ...). They are the most accurate real-time method we found. They give good information on uniform surfaces (they use a global constraint like in Horn and Schunck technique [8] to compute the flow field on uniform surfaces).

Even with state-of-the-art techniques, the optical flow we compute on the ground plane is inaccurate. Indeed, the ground plane is often poorly-textured (asphalt on the road) and is a large part of the image. We can see on figure 3 that the optical flow of the ground plane is inaccurate.

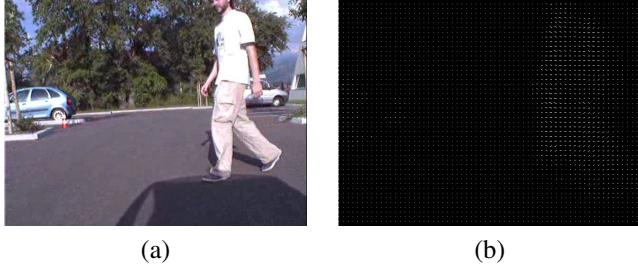


Fig. 3. Figure (a) is an image from a video sequence where the camera is translating and the pedestrian is moving in front of the robot. Figure (b) is the corresponding optical flow computed with [3]. Note the incorrect optical flow vectors on the ground plane.

2) *Odometry based optical flow model*: Having observed that the optical flow computation does not give results good enough for our optimisation, we proposed a reverse method which tries to match an optical flow model given by the odometry data to the image.

To describe our model, we will use the projective geometry formalism. We will call (u, v, w) the homogeneous coordinates of a pixel in the image, \mathcal{H} the homography matrix that projects one point on the ground plane in the image and $\dot{\mathcal{H}}$ its derivative.

The projection equation is:

$$\begin{pmatrix} \dot{u} \\ \dot{v} \\ \dot{w} \end{pmatrix}_{Image} = \dot{\mathcal{H}} \begin{pmatrix} X \\ Y \\ 1 \end{pmatrix}_{Image} \quad (1)$$

We can then infer the optical flow equation by differentiating equation 1 and we obtain:

$$\begin{pmatrix} \dot{u} \\ \dot{v} \\ \dot{w} \end{pmatrix}_{Image} = \dot{\mathcal{H}} \mathcal{H}^{-1} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix}_{Image} \quad (2)$$

Now we can obtain the optical flow vector \vec{f} for the pixel at Euclidean coordinates (u, v) by the formula:

$$\vec{f}(u, v) = \begin{pmatrix} \dot{u} - u\dot{w} \\ \dot{v} - v\dot{w} \end{pmatrix} \quad (3)$$

Finally from equations (2) and (3) we can express the theoretical optical flow vector for each pixel in the image (with the assumption that each pixel is in the ground plane).

The homography matrix and its derivative are evaluated using the position of the camera (c_x, c_y, c_z) , its orientation ϕ , the motion of the camera (given by the odometry of the robot) v (linear velocity) and ω (angular velocity).

$$\mathcal{H} = \begin{pmatrix} h_{1,1} & h_{1,2} & h_{1,3} \\ h_{2,1} & h_{2,2} & h_{2,3} \\ h_{3,1} & h_{3,2} & h_{3,3} \end{pmatrix}$$

with:

$$\begin{aligned} h_{1,1} &= u_0 \cos \phi \\ h_{1,2} &= -\alpha_u \\ h_{1,3} &= \alpha_u + u_0 (-\cos \phi + c_z \sin \phi) \\ h_{2,1} &= -\alpha_u \sin \phi + v_0 \cos \phi \\ h_{2,2} &= 0 \\ h_{2,3} &= \alpha_v (\sin \phi + c_z \cos \phi) + v_0 (-\cos \phi + c_z \sin \phi) \\ h_{3,1} &= \cos \phi \\ h_{3,2} &= 0 \\ h_{3,3} &= -\cos \phi + c_z \sin \phi \end{aligned}$$

We have also:

$$\dot{\mathcal{H}} = \begin{pmatrix} \dot{h}_{1,1} & \dot{h}_{1,2} & \dot{h}_{1,3} \\ \dot{h}_{2,1} & \dot{h}_{2,2} & \dot{h}_{2,3} \\ \dot{h}_{3,1} & \dot{h}_{3,2} & \dot{h}_{3,3} \end{pmatrix}$$

with:

$$\begin{aligned} \dot{h}_{1,1} &= \alpha_u \omega_t \\ \dot{h}_{1,2} &= u_0 \omega_t \cos \phi \\ \dot{h}_{1,3} &= -u_0 v_t \cos \phi \\ \dot{h}_{2,1} &= 0 \\ \dot{h}_{2,2} &= (-\alpha_v \sin \phi + v_0 \cos \phi) \omega_t \\ \dot{h}_{2,3} &= (\alpha_v \sin \phi - v_0 \cos \phi) v_t \\ \dot{h}_{3,1} &= 0 \\ \dot{h}_{3,2} &= \omega_t \cos \phi \\ \dot{h}_{3,3} &= -v_t \cos \phi \end{aligned}$$

Figure 4 shows the result of our model for a camera at position $c_x = 1.74$ m, $c_y = 0$ m, $c_z = 0.83$ m and $\phi = 0$ rad. The model is valid only below the horizon line whose equation is: $y = v_0 - \alpha_v \tan \phi$. Therefore there is no flow vector above the horizon line.

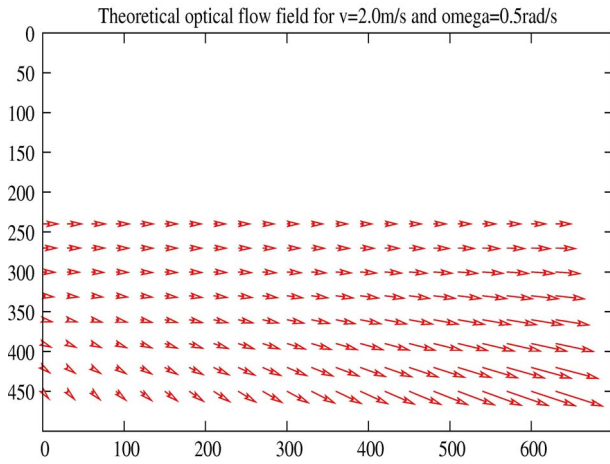


Fig. 4. Example of theoretical optical flow field for a moving robot with a velocity of $2 \text{ m}\cdot\text{s}^{-1}$ and a rotation speed of $0.5 \text{ rad}\cdot\text{s}^{-1}$

Once theoretical optical flow field is computed, we can match it to the observed data. We use two consecutive images and try to match one pixel in the previous image to the corresponding theoretical pixel in the current image. The matching is done by computing an SSD (Sum of Squared Differences) measure. An example of the SSD matching can be seen on figure 5

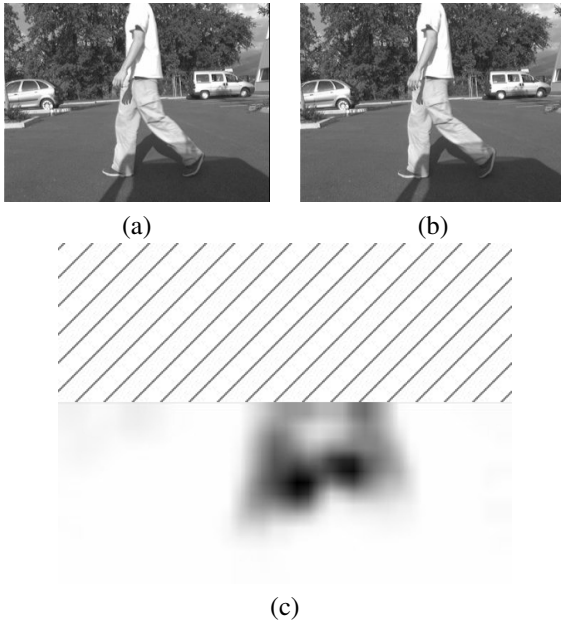


Fig. 5. Result of the optical flow matching. Image (a) and (b) are two consecutive frame of a video sequence. Subfigure (c) is the result of the SSD matching. The dashed area is the one where the model does not apply. The lighter the area is, the better the match is

III. VISION-ORIENTED DATA FUSION

In this section we propose a method to cope with the 3D occupancy grid memory and its fusion with 2D occupancy

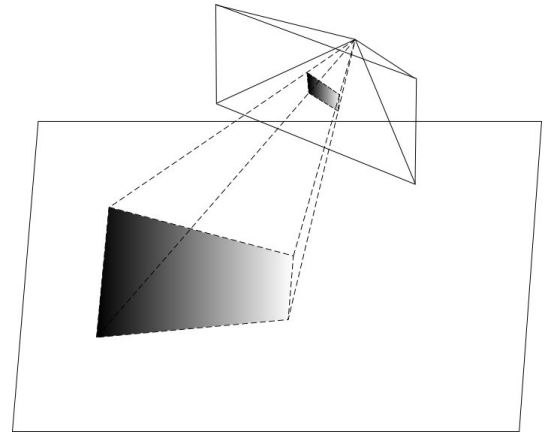


Fig. 6. Camera model, the shape in the image is projected on the ground plane. The dashed pyramid corresponds to the potentially occupied space.

grids. This model is based on the projective description of camera sensors. Using a particular sensor configuration we will be able to extract a 2D model of occupancy grid for a monocular camera.

A. 3D camera sensor model

It is difficult to deal with cameras in an occupancy grid framework. Indeed, one pixel of the image can correspond to an infinite set of 3D world points. This set of points is known as projective line.

The set of projective lines corresponding to all the image pixels is a pyramid (dimension 3). Therefore we need to express the occupancy grids in a 3D space.

Figure 6 shows the camera model we use. The shape on the image can correspond to the dashed pyramid. Saying that a pixel in the image belongs to an obstacle, means that the projective line is potentially occupied.

B. 2D projected model

The occupancy grid framework is very expensive when it comes to a 3D space. The idea is to project the 3D occupancy grid on the ground plane. This projection respects the semantics of occupancy grids and translates the incompleteness of the monocular camera into uncertainty.

The projection of a 3D occupancy grid $C_{i,j,k}^{3D}$ into a 2D occupancy grid $C_{i,j}$ is defined in our work as:

$$C_{i,j} = \max_k \left(p_k \times C_{i,j,k}^{3D} + (1 - p_k) \frac{1}{2} \right) \quad (4)$$

In equation (4) we use the maximum operator to have an occupancy grid as safe as possible. Indeed, we impose that the probability for a cell to be occupied is the maximum of the probabilities for all the 3D cells on top of the 2D ground cell, which means that we do not risk saying that a 2D cell is occupied if a 3D cell over it is occupied.

The term p_k is *a priori* knowledge. It gives *a priori* knowledge on the vertical distribution of the obstacles. The

value 1 means a strong confidence and the value 0 means that no obstacle can be at the given height.

We imposed the following function for p_k (the function is represented in figure 7):

$$p_k = \begin{cases} 1 & \text{if } k \leq z_0 \\ 2 \left(\frac{z-z_0}{\Delta z}\right)^3 - 3 \left(\frac{z-z_0}{\Delta z}\right)^2 + 1 & \text{if } k \in]z_0, z_0 + \Delta z] \\ 0 & \text{elsewhere} \end{cases} \quad (5)$$

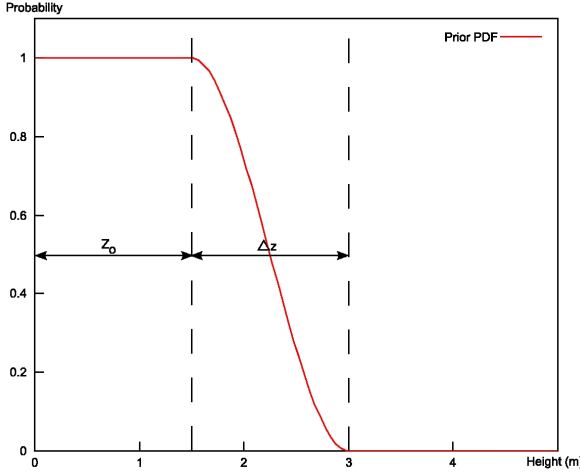


Fig. 7. Graph of the function p_k defined in equation (5)

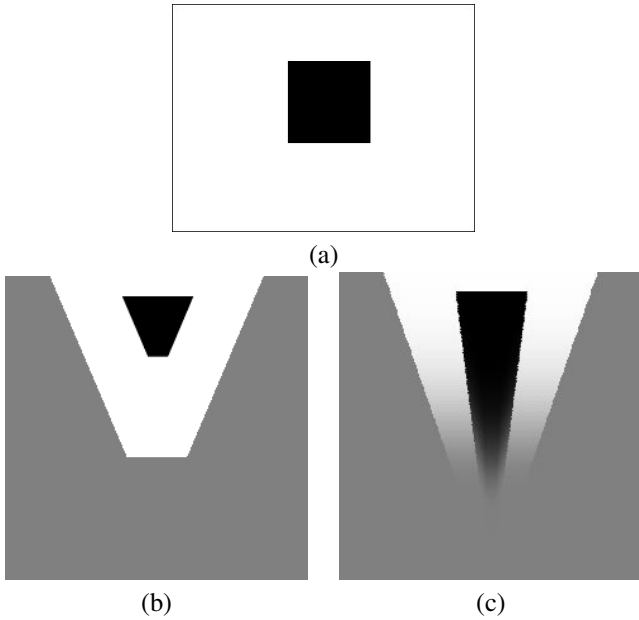


Fig. 8. (a) is a square on the image, (b) is its basic projection on the ground plane, (c) is the 3D model projected on the ground plane.

Figure 8 shows the different steps of the pyramid projection. On subfigure (c) we can see that the shape is fading out at the

bottom, this is because the sensor model we used gives more probability to the obstacles close to the ground.

IV. CAMERA MODALITIES FUSION

To fuse all the sensor modalities we use an occupancy grid framework [7], [10], to express the fusion in formal probabilistic terms. The probability for cell (i, j) , to be occupied given the sensor observations $(Z_k)_{k=1 \dots n}$ for this cell can be written as:

$$P(Occ_{i,j} | Z_1 \dots Z_n) = \frac{P(Occ_{i,j})}{P(Z_1 \dots Z_n)} \prod_{k=1}^n P(Z_k | Occ_{i,j})$$

For a given set of observation we can write:

$$P(Occ_{i,j} | Z_1 = z_1 \dots Z_n = z_n) \propto \prod_{k=1}^n P(Z_k = z_k | Occ_{i,j})$$

To add the concept of sensor confidence, we can improve the sensor model by taking into account the failure possibility. We introduce a binary random variable $H \in \{Right, Wrong\}$ that indicates if the sensor failed or not. The new sensor model is then:

$$P(Z_k = z_k | Occ_{i,j}) = \sum_H P(H) P(Z_k = z_k | Occ_{i,j}, H)$$

We consider that in the case of a failure, the sensor model is a uniform law. Therefore we write the following final sensor model:

$$P(Z_k = z_k | Occ_{i,j}) = \alpha P(Z_k = z_k | Occ_{i,j}, H = Right) + (1 - \alpha) \frac{1}{2}$$

V. EXPERIMENTAL RESULT

On figure 9 we can see the result of the whole process.

Subfigure (a) is the current frame where we analyse the results of our algorithm. Subfigure (b) is the result of the optical flow detection algorithm. We can clearly see the pedestrian moving in front of the camera. We can also see the blue car in the back. The relative importance of the car and the pedestrian is due to the distance between the camera and them. The closer to the camera the objects are, the bigger their optical flow is. In a future work we could improve this point by exploring the possibility of normalizing the SSD by the optical flow. This would result in a better ratio between far and close obstacles.

Stereo generated occupancy grid is represented on subfigure (c). It shows the pedestrian in the middle of the grid. The grid is not very dense after 3 m but give information until 7 m. We can also see a false detection on the top right of the grid.

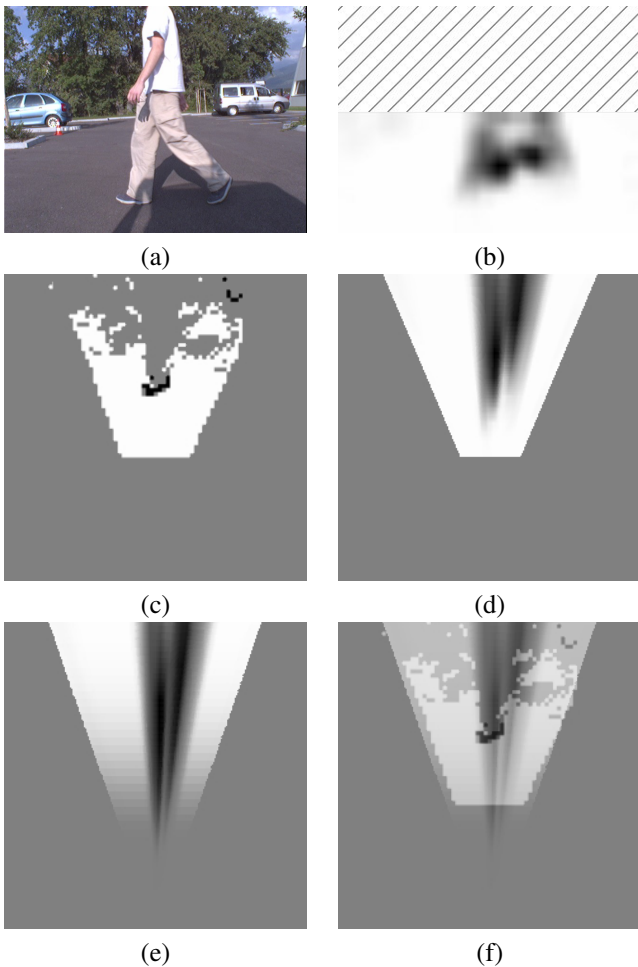


Fig. 9. (a) is the left image of the stereo camera, (b) is the detected obstacle from optical flow, (c) is the occupancy grid generated from the 3D point cloud, (d) is the projection of (b) on the ground plane, (e) is the improved model presented in III-B and (f) is the final occupancy grid which is the fusion of the two occupancy grids (c) and (e)

Subfigure (d) is the naive projection of subfigure (b), and subfigure (e) is the projection we defined in III-B.

Finally, subfigure (f) shows the global result of the fusion. We used the same confidence for both algorithms. We can see that the area where the stereo does not provide a dense information are supplemented by the optical flow algorithm. The area where the pedestrian is, is also reinforced. The false detection on the top right of the grid is minimised. After the fusion step, this false detection has a probability which means the occupancy is unknown. The cells in front of the obstacle are a little degraded.

VI. CONCLUSION AND FUTURE WORK

In this paper, we proposed a real-time method to detect obstacles using theoretical models of the ground plane using the 3D point cloud given by a stereo camera, and an optical flow field given by one of the stereo pair's camera.

The performance of the global process is better than the stereo detection or the optical flow detection alone. We could improve the quality of the occupancy grid by adding more

sensors (other cameras, laser range finders, ...) and/or more camera modalities (colour segmentation, ...).

The next step will be to perform time integration to remove some ambiguities (especially ambiguities related to monocular camera algorithms).

At least we could integrate our algorithms in a whole SLAM process.

VII. ACKNOWLEDGEMENTS

This work was done in the context of a cooperation between the CSIRO ICT Centre of Brisbane (Australia) and the INRIA Rhône-Alpes in Grenoble (France).

REFERENCES

- [1] J.L. Barron, D.J. Fleet, and S.S. Beauchemin. Performance of optical flow techniques. In *IJCV*, volume 12, pages 43–77, 1994.
- [2] T. Brox, A. Bruhn, N. Papenberg, and J. Weickert. High accuracy optical flow estimation based on a theory for warping. *Proc. 8th European Conference on Computer Vision*, 3024:25–36, may 2004.
- [3] A. Bruhn, J. Weickert, C. Feddern, T. Kohlberger, and C. Schnrr. Variational optical flow computation in real-time. *IEEE Transactions on Image Processing*, 14/5:608–615, 2005.
- [4] A. Bruhn, J. Weickert, and C. Schnrr. Lucas/kanade meets horn/schunck: Combining local and global optic flow methods. *International Journal of Computer Vision*, 61/3:211–231, 2005.
- [5] R. Collins, A. Lipton, H. Fujiyoshi, and T. Kanade. Algorithms for cooperative multisensor surveillance. *Proceedings of the IEEE*, 89(10):1456 – 1477, October 2001.
- [6] E.D. Dicksmanns. The development of machine vision for road vehicles in the last decade. *IEEE Intelligent Vehicles Symposium*, 2002.
- [7] A. Elfes. Using occupancy grids for mobile robot perception and navigation. *Computer*, 22(6):46–57, 1989.
- [8] B.K.P. Horn and B.G. Schunck. Determining optical flow. In *artificial intelligence*, volume 17, pages 185–203, 1981.
- [9] Q. Ke and T. Kanade. Transforming camera geometry to a virtual downward-looking camera: robust ego-motion estimation and ground layer detection. *Proc. of IEEE International Conference on Computer Vision and Pattern Recognition*, 2003.
- [10] K. Konolige. Improved occupancy grids for map building. *Autonomous Robots*, 4:351–367, 1997.
- [11] H.C. Longuet-Higgins. The visual ambiguity of a moving plane. In *Royal Society London*, London, Great Britain, 1984.
- [12] J. A. Nelder and R. Mead. A simplex method for function minimization. *Comput. J.* 7, pages 308–313, 1965.
- [13] C. Stauffer and W.E.L. Grimson. Adaptive background mixture models for real-time tracking. *Proc. of the International Conference on Computer Vision and Pattern Recognition*, January 1998.
- [14] A. Talukder and L. Matthies. Real-time detection of moving objects from moving vehicle using dense stereo and optical flow. *Proc. of the International Conference on Intelligent Robots and Systems*, October 2004.
- [15] K. Young-Geun and K. Hakil. Layered ground floor detection fo vision-based mobile robot navigation. In *International Conference on Robotics and Automation*, pages 13–18, New Orleans, april 2004.