

Estimating Multiple Depths in Semi-transparent Stereo Images

M. Borga, H. Knutsson
Computer Vision Laboratory
Department of Electrical Engineering
Linköping University
SE-581 83 Linköping, Sweden

Abstract

A stereo algorithm that can estimate multiple depths in semi-transparent images is presented. The algorithm is based on a combination of phase analysis and canonical correlation analysis. The algorithm adapts filters in each local neighbourhood of the image in a way which maximizes the correlation between the filtered images. The adapted filters are then analysed to find the disparity. This is done by a simple phase analysis of the scalar product of the filters. For images with different but constant depths, a simple reconstruction procedure is suggested.

Keywords: Stereo, Phase analysis, Canonical correlation analysis, Reconstruction.

1 Introduction

An important feature of binocular vision systems is *disparity*, which is a measure of the shift between two corresponding neighbourhoods in a pair of stereo images. The disparity is related to the angle the eyes (cameras) must be rotated relative to each other in order to focus on the same point in the 3-dimensional outside world. The corresponding process is known as *vergence*.

The stereo problem is closely related to motion estimation where there are two (or more) consecutive images from an image sequence rather than a stereo pair. The difference is, of course, that in stereo there is only a one-dimensional translation whereas motion estimation requires the estimation of translation in two dimensions.

The problem of estimating disparity between pairs of stereo images is not a new one [2]. Early approaches often used matching of some feature in the two images [8]. The simplest way to calculate the disparity is to correlate a region in one image with all horizontally shifted regions on the same vertical position and then to find the shift that gave maximum correlation. This is, however, a computationally very expensive method.

Later approaches have been more focused on using the phase information given by for example Gabor filters or quadrature filters [9, 11, 7, 10]. An advantage of

phase-based methods is that phase is a continuous variable that allows for sub-pixel accuracy. In phase-based methods, the disparity can be estimated as a ratio between the phase difference between corresponding vertical line/edge filter outputs from the two images and the instantaneous frequency.

A problem that standard phase-based methods can not handle, however, is to estimate multiple disparities in one position. This is the case at depth discontinuities and in semi-transparent images, i.e. images that are sums of images with different depths. Such images are typical in many medical applications such as x-ray images. An every-day example of this kind of image is obtained by looking through a window with reflection. (The effect on the intensity of a light- or X-ray when passing two objects is in fact multiplicative, but a logarithmic transfer function is usually applied when generating X-ray images which makes the images additive.)

This problem (both at depth discontinuities and in semi-transparent images) can be solved with a technique that combines phase analysis and *canonical correlation analysis* (CCA) [3, 4]. This technique is based on a method that uses CCA for combining filters to design feature detectors in images [5].

In the following section, we give a brief overview of the theory of canonical correlation analysis. In section 3, the CCA-based stereo algorithm is described in detail. In section 4 is explained how multiple depth estimates are obtained. Some experimental results are presented in section 5 and, finally, in section 6 we summarize and make some concluding remarks.

2 Canonical correlation analysis

Consider two random variables, \mathbf{x} and \mathbf{y} , from a multi-normal distribution:

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim N \left(\begin{bmatrix} \mathbf{x}_0 \\ \mathbf{y}_0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_{xx} & \mathbf{C}_{xy} \\ \mathbf{C}_{yx} & \mathbf{C}_{yy} \end{bmatrix} \right), \quad (1)$$

where $\mathbf{C} = \begin{bmatrix} \mathbf{C}_{xx} & \mathbf{C}_{xy} \\ \mathbf{C}_{yx} & \mathbf{C}_{yy} \end{bmatrix}$ is the covariance matrix. \mathbf{C}_{xx} and \mathbf{C}_{yy} are nonsingular matrices and $\mathbf{C}_{xy} = \mathbf{C}_{yx}^T$.

Consider the linear combinations, $x = \mathbf{w}_x^T(\mathbf{x} - \mathbf{x}_0)$ and $y = \mathbf{w}_y^T(\mathbf{y} - \mathbf{y}_0)$, of the two variables respectively. The correlation between x and y is given by the following example:

$$\rho = \frac{\mathbf{w}_x^T \mathbf{C}_{xy} \mathbf{w}_y}{\sqrt{\mathbf{w}_x^T \mathbf{C}_{xx} \mathbf{w}_x \mathbf{w}_y^T \mathbf{C}_{yy} \mathbf{w}_y}}, \quad (2)$$

see for example [1]. The correlation ρ is a function of \mathbf{w}_x and \mathbf{w}_y . The extremum points of equation 2 are given by the solutions to an eigenvalue problem [3]:

$$\begin{bmatrix} \mathbf{C}_{xx} & [0] \\ [0] & \mathbf{C}_{yy} \end{bmatrix}^{-1} \begin{bmatrix} [0] & \mathbf{C}_{xy} \\ \mathbf{C}_{yx} & [0] \end{bmatrix} \begin{pmatrix} \hat{\mathbf{w}}_x \\ \hat{\mathbf{w}}_y \end{pmatrix} = \rho \begin{pmatrix} \lambda_x \hat{\mathbf{w}}_x \\ \lambda_y \hat{\mathbf{w}}_y \end{pmatrix} \quad (3)$$

where: $\rho, \lambda_x, \lambda_y > 0$ and $\lambda_x \lambda_y = 1$. Equation (3) can be rewritten as:

$$\begin{cases} \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \hat{\mathbf{w}}_y = \rho \lambda_x \hat{\mathbf{w}}_x \\ \mathbf{C}_{yy}^{-1} \mathbf{C}_{yx} \hat{\mathbf{w}}_x = \rho \lambda_y \hat{\mathbf{w}}_y \end{cases} \quad (4)$$

Solving (4) gives N solutions $\{\rho_n, \hat{\mathbf{w}}_{xn}, \hat{\mathbf{w}}_{yn}\}$, $n = \{1..N\}$. N is the minimum of the input dimensionality and the output dimensionality. The linear combinations, $x_n = \hat{\mathbf{w}}_{xn}^T \mathbf{x}$ and $y_n = \hat{\mathbf{w}}_{yn}^T \mathbf{y}$, are termed *canonical variates* and the correlations, ρ_n , between these variates are termed the *canonical correlations* [6]. An important aspect in this context is that the canonical correlations are *invariant to affine transformations* of \mathbf{x} and \mathbf{y} . Also note that the canonical variates corresponding to the different roots of (4) are uncorrelated, implying that:

$$\begin{cases} \mathbf{w}_{xn}^T \mathbf{C}_{xx} \mathbf{w}_{xm} = 0 \\ \mathbf{w}_{yn}^T \mathbf{C}_{yy} \mathbf{w}_{ym} = 0 \\ \mathbf{w}_{xn}^T \mathbf{C}_{xy} \mathbf{w}_{ym} = 0 \end{cases} \quad \text{if } n \neq m \quad (5)$$

It should be noted that (3) is a special case of the *generalized eigenproblem* [3]:

$$\mathbf{A} \mathbf{w} = \lambda \mathbf{B} \mathbf{w}.$$

3 The stereo algorithm

The basic idea behind the stereo algorithm described here is to let the system adapt filters to fit the disparity in question instead of using fixed filters. The algorithm consists of two parts: CCA and phase analysis. Both are performed for each disparity estimate. Canonical correlation analysis is used to create adaptive linear combinations of quadrature filters. These linear combinations are new quadrature filters that are adapted in frequency response and spatial position in order to maximize the correlation between the filter outputs from the two images.

These new filters are then analysed in the phase analysis part of the algorithm. The coefficients given by the canonical correlation vectors are used as weighting coefficients in a pre-computed table that allows for an efficient phase-based search for disparity.

It is, of course, possible to use other basis functions than quadrature filters, or even use the pixel base itself, in the canonical correlation analysis. The advantage of having complex basis filters such as quadrature filters is that it allows for the phase-based search which is efficient and can give sub-pixel accuracy.

In the following two subsections, the two parts of the stereo algorithm are described in more detail.

3.1 Canonical correlation analysis part

The input \mathbf{x} and \mathbf{y} to the CCA come from the left and right images respectively. Each input is a vector with outputs from a set of M quadrature filters:

$$\mathbf{x} = \begin{pmatrix} q_{x1} \\ \vdots \\ q_{xM} \end{pmatrix} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} q_{y1} \\ \vdots \\ q_{yM} \end{pmatrix}, \quad (6)$$

where q_i is the (complex) filter output for the i th quadrature filter in the filter set. In the implementation described here, the filter set consists of two identical one-dimensional (horizontal) quadrature filters with two pixels relative displacement. (Other and larger sets of filters can be used including, for example, filters with different bandwidths, different centre frequencies, different positions, etc.)

The data is sampled from a neighbourhood \mathcal{N} around the point for the disparity estimate. The choice of neighbourhood size is a compromise between noise sensitivity and locality. The covariance matrix \mathbf{C} is calculated using the vectors \mathbf{x} and \mathbf{y} in \mathcal{N} . The fact that quadrature filters have zero DC component simplifies this calculation to an outer product sum:

$$\mathbf{C} = \sum_{\mathcal{N}} \begin{pmatrix} \mathbf{x}_i \\ \mathbf{y}_i \end{pmatrix} \begin{pmatrix} \mathbf{x}_i \\ \mathbf{y}_i \end{pmatrix}^T \quad (7)$$

The first canonical correlation ρ_1 and the corresponding vectors \mathbf{w}_x and \mathbf{w}_y are then calculated by solving (4). In the case where only two filters are used, this calculation becomes very simple. If very large sets of filters are used, the covariance matrix gets very big and an analytical calculation of the canonical correlation becomes computationally very expensive. In such a case, an iterative $\mathcal{O}(n)$ algorithm, that avoids outer products and matrix inverses, can be used [3, 5].

The canonical correlation vectors \mathbf{w}_x and \mathbf{w}_y define two new filters, $\mathbf{f}_x = \sum_{i=1}^N w_{xi} \mathbf{f}_i$ and $\mathbf{f}_y = \sum_{i=1}^N w_{yi} \mathbf{f}_i$ where \mathbf{f}_i are the basis filters, N is the number of filters in the filter set and w_{xi} and w_{yi} are the components

in the first pair of canonical correlation vectors. This means that the new filters \mathbf{f}_x and \mathbf{f}_y have maximally correlated output in \mathcal{N} , given the set of basis filters \mathbf{f}_i .

3.2 Phase analysis part

The key idea of this part is to search for the disparity that corresponds to a real-valued correlation between the output of the two new filters. This idea is based on the fact that canonical correlations are real valued [3]. In other words, find the disparity δ such that

$$\text{Im} [\text{Corr} (q_y(\xi + \delta), q_x(\xi))] = \text{Im} [c(\delta)] = 0, \quad (8)$$

where q_x and q_y are the left and right filter outputs respectively and ξ is the spatial (horizontal) coordinate.

A calculation of the correlation over \mathcal{N} for all δ would be very expensive. A much more efficient solution is to assume that the signal \mathbf{s} can be described by a covariance matrix \mathbf{C}_{ss} . Under this assumption, the correlation between the left filter convolved with the signal \mathbf{s} and the right filter convolved with the same signal shifted a certain amount δ can be measured. But convolving a filter with a shifted signal is the same as convolving a shifted filter with the non-shifted signal. Hence, the correlation $c(\delta)$ can be calculated as the correlation between the left filter convolved with \mathbf{s} and a shifted version of the right filter convolved with the same signal \mathbf{s} .

Under the assumption that the signal \mathbf{s} has the covariance matrix \mathbf{C}_{ss} , the correlation in equation 8 can be written as

$$\begin{aligned} c(\delta) &= \frac{E[q_x^* q_y(\delta)]}{\sqrt{E[|q_x|^2] E[|q_y|^2]}} \\ &= \frac{E[(\mathbf{s}^* \mathbf{f}_x)^* (\mathbf{s}^* \mathbf{f}_y(\delta))]}{\sqrt{E[(\mathbf{s}^* \mathbf{f}_x)^* (\mathbf{s}^* \mathbf{f}_x)] E[(\mathbf{s}^* \mathbf{f}_y)^* (\mathbf{s}^* \mathbf{f}_y)]}} \\ &= \frac{E[\mathbf{f}_x^* \mathbf{s} \mathbf{s}^* \mathbf{f}_y(\delta)]}{\sqrt{E[\mathbf{f}_x^* \mathbf{s} \mathbf{s}^* \mathbf{f}_x] E[\mathbf{f}_y^*(\delta) \mathbf{s} \mathbf{s}^* \mathbf{f}_y(\delta)]}} \\ &= \frac{\mathbf{f}_x^* \mathbf{C}_{ss} \mathbf{f}_y(\delta)}{\sqrt{\mathbf{f}_x^* \mathbf{C}_{ss} \mathbf{f}_x \mathbf{f}_y^* \mathbf{C}_{ss} \mathbf{f}_y}}, \end{aligned} \quad (9)$$

where $\mathbf{f}_y(\delta)$ is a shifted version of \mathbf{f}_y . Note that the quadrature filter outputs have zero mean, which is necessary for the first equality.

A lot of the computations needed to calculate $c(\delta)$ can be saved since

$$\begin{aligned} \mathbf{f}_x^* \mathbf{C}_{ss} \mathbf{f}_y(\delta) &= \left(\sum_{i=1}^M w_{xi} \mathbf{f}_i \right)^* \mathbf{C}_{ss} \left(\sum_{j=1}^M w_{yj} \mathbf{f}_j(\delta) \right) \\ &= \sum_{i=1}^M \sum_{j=1}^M w_{xi} w_{yj} \mathbf{f}_i^* \mathbf{C}_{ss} \mathbf{f}_j(\delta) = \sum_{ij} v_{ij} g_{ij}(\delta), \end{aligned} \quad (10)$$

where

$$g_{ij}(\delta) = \mathbf{f}_i^* \mathbf{C}_{ss} \mathbf{f}_j(\delta). \quad (11)$$

The function $g_{ij}(\delta)$ does not depend on the result from the CCA and can therefore be calculated in advance for different disparities δ and stored in a table. The denominator in equation 9 can be treated in the same way but does not depend on δ :

$$\mathbf{f}_x^* \mathbf{C}_{ss} \mathbf{f}_x = \sum_{ij} v_{ij}^x g_{ij}(0) \quad \text{and} \quad \mathbf{f}_y^* \mathbf{C}_{ss} \mathbf{f}_y = \sum_{ij} v_{ij}^y g_{ij}(0), \quad (12)$$

where $v_{ij}^x = w_{xi}^* w_{xj}$ and $v_{ij}^y = w_{yi}^* w_{yj}$. Note that the filter vectors \mathbf{f} must be padded with zeros at both ends to enable the scalar product between a filter and a shifted filter δ . (The zeros do not, of course, affect the result of equation 10.) In the case of two basis filters, the table contains four rows and eight constants.

Hence, for a given disparity a (complex) correlation $c(\delta)$ can be computed as a normalized weighted sum:

$$c(\delta) = \frac{\sum_{ij} v_{ij} g_{ij}(\delta)}{\sqrt{\sum_{ij} v_{ij}^x g_{ij}(0) \sum_{ij} v_{ij}^y g_{ij}(0)}}. \quad (13)$$

The aim is to find the δ for which the correlation $c(\delta)$ is real valued. This is done by finding the zero crossings of the phase of the correlation. A very coarse quantization of δ can be used in the table since the phase is, in general, rather linear near the zero crossing (as opposed to the imaginary part which in general is not linear). Hence, first a coarse estimate of the zero crossing is obtained. Then the derivative of the phase at the zero crossing is measured, using two neighbouring samples. Finally, the error in the coarse estimate is compensated for by using the actual phase value and the phase derivative at the estimated position:

$$\delta = \delta_c - \frac{\varphi(\delta_c)}{\partial \varphi / \partial \delta}, \quad (14)$$

where δ_c is the coarse estimate of the zero crossing and $\varphi(\delta_c)$ is the complex phase of $c(\delta_c)$ (see figure 1).

If the signal model is uncorrelated white noise, \mathbf{C}_{ss} is the identity matrix and the calculations of the values in the table reduce to a simple scalar product: $g_{ij}(\delta) = \mathbf{f}_i^* \mathbf{f}_j(\delta)$. There is no computational reason to choose white noise as signal model if there is a better model, since the table is calculated only once. However, experiments show that the results are in practice almost identical when using the identity matrix and when using a covariance matrix estimated from the image [3].

4 Multiple disparities

If more than one zero crossing are detected, the magnitudes of the correlations can be used to select a solution. Since the CCA searches for maximum correlation,

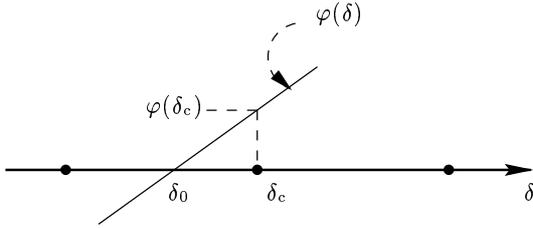


Figure 1: The estimation of the coordinate δ_0 of the phase zero crossing using the coarse estimate δ_c of the zero crossing, the phase value $\varphi(\delta_c)$ and the derivative at the coarse estimate. The black dots illustrate the sampling points of the phase given by the table $g_{ij}(\delta)$.

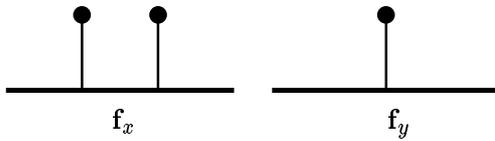


Figure 2: A simple example of a pair of filters that have two correlation peaks.

the zero crossing with maximum correlation $c(\delta)$ is most likely to be the best estimate. If two zero crossings have approximately equal magnitude (and the canonical correlation ρ is high), both disparity estimates can be considered to be correct within the neighbourhood, which indicates either a depth discontinuity or that there really exist two disparities.

In the case of discontinuities there are in practice only one or two disparities at each point. In the case of semi-transparent images, however, there could of course be more than two disparities. Still each disparity is associated with a certainty measure given by the magnitude of the correlation at the zero crossing. So even if the algorithm gives several disparity estimates in each point, most of the estimates will in general be associated with a very low certainty measure which means that they should have very little influence to whatever decision is based on the disparity estimate. If, however, an estimate of the number of disparities actually is required, a threshold can be applied on the relative certainties. There is of course no general way to decide this threshold. To simplify the discussion, we will in the following assume *two* disparities in each point.

Note that both disparity estimates are represented by the same canonical correlation solution. This means that the CCA must generate filters that have correlation peaks for *two* different disparities. To see how this can be done, consider the simple filter pair illustrated in figure 2. The autocorrelation function (or convolution) between these two filters is identical to the left

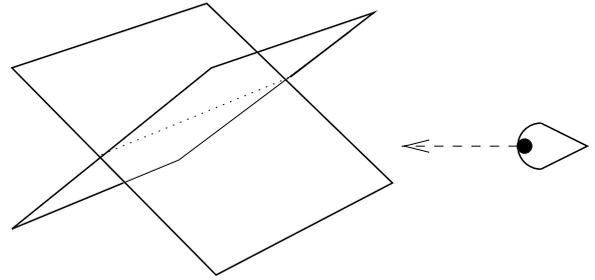


Figure 3: The test image scene for semi-transparent images.

filter, which consists of two impulses. The example is much simplified, but illustrates the possibility of having a pair of filters with two correlation peaks. If the CCA was used directly on the pixel data instead of on the quadrature filter outputs, such a filter pair *could* develop. In the present method, the image data are represented by using other, complex, basis functions (the quadrature filters of the basis filter set), but it is still possible to construct filters with two correlation peaks.

5 Experimental results

In the experiments presented here a basis filter set have been used consisting of two one-dimensional horizontally oriented quadrature filters, both with a centre frequency of $\pi/4$ and a bandwidth of two octaves. The filters have 15 coefficients in the spatial domain and are shifted two pixels relative to each other. The frequency function is a quadratic cosine on a log scale:

$$F(u) = \cos^2(k \ln(u/u_0)) \quad (15)$$

where $k = \pi / (2 \ln(2))$ and $u_0 = \pi/4$.

5.1 Crossing planes

The test images in this experiment were generated as a sum of two images with white uncorrelated noise. The images were tilted in opposite directions around the horizontal axis. The disparity range was $+/- 5$ pixels. Figure 3 illustrates the test scene. The stereo pair is shown in figure 4. Here, the averaging or fusion performed by the human visual system for small disparities can be seen in the middle of the image. A neighbourhood \mathcal{N} of 31×3 pixels was used for the CCA. The result is shown in figure 5. The results show that the disparities of both the planes are approximately estimated. In the middle, where the disparity difference is small, the result is an average between the two disparities.



Figure 4: The stereo image pair for the crossing planes.



Figure 6: The stereo image pair for the summed real images.

5.2 Real images

The second experiment setup consists of two images that are shifted horizontally relative each other and added together. The shift is ± 2 pixels so the total disparity is 4 pixels. The stereo pair is showed in figure 6. Here, a large neighbourhood of 100×100 pixels was used since the shift was constant over the image. The disparity estimates are shown in the histogram in figure 7. The peak at zero is caused by edge effects. the other two peaks are at -2.13 and $+2.56$ pixels. This means that the disparity is slightly over-estimated. To see how this error effects the separation of the images, a simple reconstruction operation has been performed. The recon-

struction is performed by differential operation followed by an integration. First the images are shifted half the estimated disparity relative to each other and added together. This results in a image where one of the original images is differentiated and the other image is more or less cancelled out. Such a differential image is shown in figure 8. The differential image is then integrated using a cumulative summation over each line. The procedure is then repeated with a shift in the opposite direction to get the other image. The resulting images are shown in figures 9 and 10. This simple reconstruction procedure is, of course, not optimal since the differentiation operator is not $(-1, 1)$ but have the $+1$ and -1 separated by the shift. A more optimal reconstruction is possible but

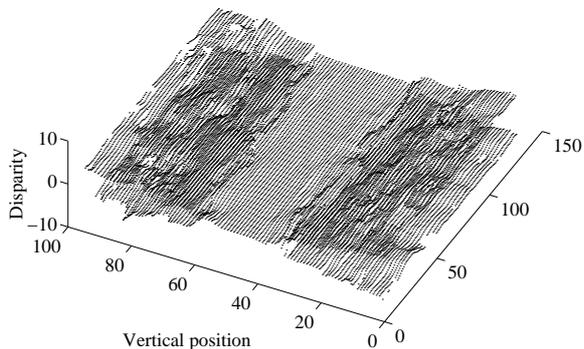


Figure 5: The result for the semi-transparent images. The disparity estimates are coloured to simplify the visualization.

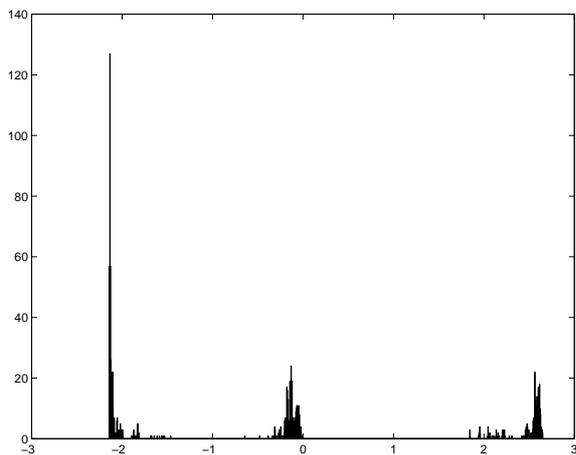


Figure 7: The disparity estimates for the stereo pair in figure 6.

much more complex.

6 Summary and discussion

We have presented a stereo algorithm with sub-pixel accuracy that can handle multiple depths in semi-transparent images. The algorithm combines canonical correlation analysis and phase analysis. So far we have only used a basis filter set of two identical filters shifted two pixels. A larger filter set can be used which may contain filters with different spatial positions as well as filters with other frequency functions. Such a set would allow for a wider range of disparities, more simultaneous estimates and higher resolution.



Figure 8: One of the two differential images in the reconstruction from the images in figure 6.



Figure 9: The first reconstructed image from the images in figure 6.

The choice of neighbourhood \mathcal{N} for the CCA is of course important for the result. If there is a priori knowledge of the shape of the regions that have relatively constant depths, the neighbourhood should, of course, be chosen accordingly. This means that if the



Figure 10: The second reconstructed image from the images in figure 6.

disparity is known to be relatively constant along the horizontal axis, for example, the shape of the neighbourhood should be elongated horizontally, as in the experiment on artificial data in the previous section. It is also possible to let the algorithm select a suitable neighbourhood shape automatically. One way to accomplish this is to measure the canonical correlation for a few different neighbourhood shapes. These shapes could be, for example, one horizontally elongated, one vertically elongated and one square. The algorithm should then use the result from the neighbourhood that gave the highest canonical correlation to estimate the disparity.

Another natural extension of the algorithm is to include also vertical shifts. Such an extended algorithm could for example be used for motion estimation.

References

- [1] T. W. Anderson. *An Introduction to Multivariate Statistical Analysis*. John Wiley & Sons, second edition, 1984.
- [2] S. T. Barnard and M. A. Fichsler. Computational Stereo. *ACM Comput. Surv.*, 14:553–572, 1982.
- [3] M. Borga. *Learning Multidimensional Signal Processing*. PhD thesis, Linköping University, Sweden, SE-581 83 Linköping, Sweden, 1998. Dissertation No 531, ISBN 91-7219-202-X.
- [4] M. Borga and H. Knutsson. An adaptive stereo algorithm based on canonical correlation analysis. In B. Verma, Z. Liu, A. Sattar, T. Zurawski, and J. You, editors, *Proceedings of the Second IEEE International Conference on Intelligent Processing Systems*, pages 177–182, Gold Coast, Australia, August 1998. IEEE. Also as report: LiTH-ISY-R-2013.
- [5] M. Borga, H. Knutsson, and T. Landelius. Learning Canonical Correlations. In *Proceedings of the 10th Scandinavian Conference on Image Analysis*, Lappeenranta, Finland, June 1997. SCIA.
- [6] H. Hotelling. Relations between two sets of variates. *Biometrika*, 28:321–377, 1936.
- [7] A. D. Jepson and D. J. Fleet. Scale-space singularities. In O. Faugeras, editor, *Computer Vision-ECCV90*, pages 50–55. Springer-Verlag, 1990.
- [8] D. Marr. *Vision*. W. H. Freeman and Company, New York, 1982.
- [9] T. D. Sanger. Stereo disparity computation using gabor filters. *Biological Cybernetics*, 59:405–418, 1988.
- [10] C-J. Westelius. *Focus of Attention and Gaze Control for Robot Vision*. PhD thesis, Linköping University, Sweden, SE-581 83 Linköping, Sweden, 1995. Dissertation No 379, ISBN 91-7871-530-X.
- [11] R. Wilson and H. Knutsson. A multiresolution stereopsis algorithm based on the Gabor representation. In *3rd International Conference on Image Processing and Its Applications*, pages 19–22, Warwick, Great Britain, July 1989. IEE.