

Wavelet-Based Correlation for Stereopsis

Maureen Clerc

CERMICS, INRIA
BP 92, 06902 Sophia-Antipolis
France

Abstract. Position disparity between two stereoscopic images, combined with camera calibration information, allow depth recovery. The measurement of position disparity is known to be ambiguous when the scene reflectance displays repetitive patterns. This problem is reduced if one analyzes scale disparity, as in shape from texture, which relies on the deformations of repetitive patterns to recover scene geometry from a single view.

These observations lead us to introduce a new correlation measure based not only on position disparity, but on position and scale disparity. Local scale disparity is expressed as a change in the scale of wavelet coefficients. Our work is related to the spatial frequency disparity analysis of Jones and Malik (ECCV92). We introduce a new wavelet-based correlation measure, and we show its application to stereopsis. We demonstrate its ability to reproduce perceptual results which no other method of our knowledge had accounted for.

Introduction

Position disparity is a well-known cue used for stereoscopic depth reconstruction, but it is not the only one. Jones and Malik have demonstrated the importance of frequency disparity information in stereopsis [4,5], and perceptual results have shown that stereo image pairs created with bandpass filtered noise, and designed to contain frequency disparity information, but no position disparity information, lead to slant perception [12]. In shape from texture, which studies the recovery of scene depth from the monocular image of a homogeneously textured object, the shape cue comes from the scale or the frequency disparity between different positions in the image [3,6].

Many stereoscopic algorithms which produce dense disparity maps require a measure of similarity between regions from the image pair, and define disparity as the local translation maximizing this similarity. In order not to be sensitive to noise, the similarity measure is usually based on averages of image intensities over a region. A geometrical problem arises when the local scene element being viewed is not fronto-parallel: because of projective distortion, the shapes and sizes of corresponding image patches are not identical in the two images. Since the distortion parameters depend on the local orientation of the scene whose shape is being computed, an iterative scheme must be used, in order to improve the similarity measure by using the scene shape calculated at the previous step [2].

Repetitive texture (regular grating, woven fabric), is difficult to handle with classical similarity measures because they produce many local maxima. We introduce a new similarity measure, incorporating a scale disparity constraint, which we observe in practice to lift the repetitive texture ambiguity.

Wavelets, whose parameters are position and scale, are ideally suited to define a correlation measure based on position and scale disparity. Wavelet methods have been proposed for stereopsis in the context of coarse-to-fine disparity measurement [9,10]. The method proposed here is related to [8], in which an affine transformation between image patches is measured at the output of a set of filters.

We demonstrate through numerical examples in 1D and 2D that our wavelet-based correlation is smooth and well-behaved and generally displays only one local maximum in the presence of repetitive texture. This makes it a good ingredient for stereopsis algorithms which rely on a similarity measure [13].

1 Wavelet Analysis of Distortion

In this paper, we neglect occlusion effects, which can be taken into account at a higher level, for instance in a cooperative stereopsis algorithm [13]. We suppose the stereoscopic pair of images I_l and I_r to satisfy

$$I_l(x) = I_r(d(x)) .$$

where

$$d(x) = d(x_1, x_2) = (d_1(x_1, x_2), d_2(x_1, x_2))$$

is a continuous map, which we call the disparity¹ between I_l and I_r . Let $\psi(x)$ be a wavelet, i.e. a two-dimensional oscillating function whose spatial support is localized around $(0, 0)$ and whose spatial frequency support concentrates around a frequency $\xi \neq (0, 0)$. A Gabor wavelet, which is a Gaussian window modulated to oscillate at a frequency ξ , satisfies these requirements. Let u be a position in \mathbb{R}^2 and let S denote a positive definite 2×2 matrix. An affine transformation

$$\psi_{u,S}(x) = (\det S)^{-1} \psi(S^{-1}(x - u)) .$$

modifies the space-scale localization of the wavelet to position u and frequency $S^T \xi$, where S^T is the matrix transpose of S .

The wavelet coefficients of I are defined by

$$W(u, S) = \langle I, \psi_{u,S} \rangle = \int I(x) \psi_{u,S}^*(x) dx . \quad (1)$$

The squared amplitude of wavelet coefficients $|W(u, S)|^2$ measures the energy contained in a surface patch of the image I centered at u , around spatial frequency $S^T \xi$.

¹ Disparity is generally defined as $d(x) - x$.

Let us compare the wavelet coefficients of I_l (denoted $W_l(u, S)$) to those of I_r (denoted $W_r(u, S)$), supposing for the moment that the disparity d between the two images is an affine transformation

$$d(x) = d(u) + J \times (x - u) ,$$

where J is a 2×2 matrix. A simple change of variable² in (1) yields

$$W_l(u, S) = W_r(d(u), J S) .$$

This relationship between wavelet coefficients makes apparent the position and scale disparities between the two images.

In the case of a general, no longer affine, disparity, if d is differentiable, it can be approximated, locally around a position u , by its first-order Taylor approximation

$$d(x) \approx d(u) + J(u) \times (x - u) ,$$

where $J(u)$ is the 2×2 Jacobian matrix of d at position u

$$J(u) = \begin{pmatrix} \frac{\partial d_1}{\partial u_1}(u) & \frac{\partial d_1}{\partial u_2}(u) \\ \frac{\partial d_2}{\partial u_1}(u) & \frac{\partial d_2}{\partial u_2}(u) \end{pmatrix} .$$

If I_l were a smooth image, it would be possible to extend the Taylor approximation to the image, and obtain

$$I_l(d(x)) \approx I_l(d(u) + J(u) \times (x - u))$$

Unfortunately, even in the absence of sharp discontinuities, an image cannot be assumed to be smooth on account of measurement noise. One can however model the image as the realization of a stochastic process, whose covariance is smooth away from the diagonal. Then one can show [1] that the variances of wavelet coefficients satisfy

$$E[|W_l(u, S)|^2] \approx E[|W_r(d(u), J(u) S)|^2] . \quad (2)$$

2 Correlation Measure

We choose a correlation measure between images I_l and I_r of the form

$$\rho = \frac{2 (\text{Feat}_l, \text{Feat}_r)}{(\text{Feat}_l, \text{Feat}_l) + (\text{Feat}_r, \text{Feat}_r)} \quad (3)$$

² The change of variable is the motivation for the L^1 normalization of the wavelet, instead of the more classical L^2 normalization.

where Feat_l and Feat_r are features relative to the two images, and (\cdot, \cdot) is an inner product in feature parameter space. Clearly, $\rho \leq 1$ and $\rho = 1$ if and only if $\text{Feat}_l = \text{Feat}_r$. Section 4.2 comments on this choice of correlation ratio.

Relationship (2) allows to derive a correlation measure which combines position disparity d and scale disparity J . Given a collection of scaling matrices S_i (typically less than 5), let

$$\text{Feat}_l(u, S_i) = E [|W_l(u, S_i)|^2]$$

and for a given position disparity d and scale disparity J , let

$$\text{Feat}_r^{(d,J)}(u, S_i) = E [|W_r(d(u), J(u) S_i)|^2] .$$

We define the features at a position u as the collection of wavelet coefficient variances at preselected scales:

$$\text{Feat}_l(u) = \{ \text{Feat}_l(u, S_i) \}_{\{S_i\}} \tag{4}$$

$$\text{Feat}_r^{(d,J)}(u) = \{ \text{Feat}_r^{(d,J)}(u, S_i) \}_{\{S_i\}} . \tag{5}$$

The inner product (\cdot, \cdot) is then simply

$$(\text{Feat}_l, \text{Feat}_r) = \sum_{S_i} \text{Feat}_l(u, S_i) \cdot \text{Feat}_r(u, S_i) .$$

We finally obtain a correlation measure which depends on position disparity d and scale disparity J :

$$\rho(u, d, J) = \frac{2 \sum_{S_i} \text{Feat}_l(u, S_i) \cdot \text{Feat}_r^{(d,J)}(u, S_i)}{\sum_{S_i} (\text{Feat}_l(u, S_i))^2 + (\text{Feat}_r^{(d,J)}(u, S_i))^2} . \tag{6}$$

Note that one is in practice limited to a unique realization of the images. In order to estimate the variance of wavelet coefficients, we rely on estimation results from [1]. We estimate $E[|W(u, S)|^2]$ by averaging $|W(v, S)|^2$ for v in a neighborhood $B(u)$ of u . This estimation procedure, which is proved to be consistent, justifies the use of wavelet filters rather than Gaussian filters in the correlation measure.

For a given scaling matrix S_i , $\text{Feat}_l(u, S_i)$ is estimated with

$$\widehat{\text{Feat}}_l(u, S_i) = \int_{B(u)} |W_l(v, S_i)|^2 dv$$

and the corresponding $\text{Feat}_r^{(d,J)}(u, S_i)$ is estimated with

$$\widehat{\text{Feat}}_r^{(d,J)}(u, S_i) = \int_{B(u)} |W_r(d(v), J(v) S_i)|^2 dv .$$

3 Local Shape Measurement

We have relaxed the classical stereopsis problem by introducing a new parameter for image matching: scale disparity J . At first view, the problem may appear more difficult to solve, because finding the best local match between I_l and I_r now requires to maximize ρ (defined in (6)) over d and J , instead of d only.

This difficulty disappears when recalling the objective, which is to measure scene depth. At a given position u in image I_l , position disparity $d(u)$ and scale disparity $J(u)$ are both a function of the local position and orientation of the scene element being viewed. Consider a simplified 2D geometry displayed in Figure 3, where local orientation is defined by a unique angle θ , and position is determined by the signed distance p from the surface tangent to a reference point O . Appendix A gives the expressions of position disparity $d_{p,\theta}$ and scale disparity $J_{p,\theta}$. In 3D, local geometry would be expressed by three scalar parameters (p, θ_1, θ_2) .

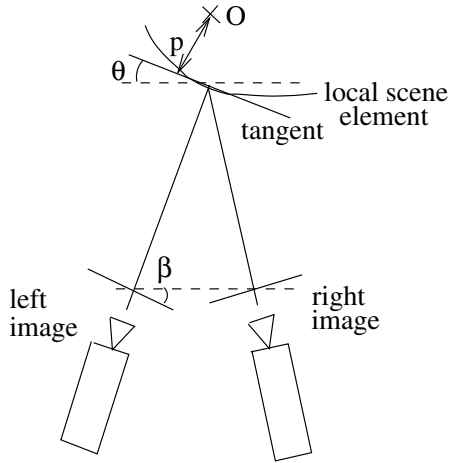


Fig. 1. Simplified 2D geometry.

The correlation measure ρ can then be reformulated as a function of p and θ instead of d and J :

$$\rho(u, p, \theta) = \frac{2 \sum_{S_i} \text{Feat}_l(u, S_i) \cdot \text{Feat}_r^{(d_{p,\theta}, J_{p,\theta})}(u, S_i)}{\sum_{S_i} (\text{Feat}_l(u, S_i))^2 + (\text{Feat}_r^{(d_{p,\theta}, J_{p,\theta})}(u, S_i))^2}. \quad (7)$$

4 Numerical Results

In this section, we compare the correlation measure $\rho(u, p, \theta)$ defined in (7) to the classical area correlation $\rho_0(u, p, \theta)$ defined by

$$\rho_0(u, p, \theta) = \frac{\int_{B(u)} \tilde{I}_l(v) \tilde{I}_r(d_{p,\theta}(v)) dv}{\left(\int_{B(u)} \tilde{I}_l(v)^2 dv \int_{B(u)} \tilde{I}_r(d_{p,\theta}(v))^2 dv \right)^{1/2}} ,$$

where

$$\tilde{I}_l(v) = I_l(v) - \frac{1}{\text{area}(B(u))} \int_{B(u)} I_l(v) dv$$

and

$$\tilde{I}_r(d_{p,\theta}(v)) = I_r(d_{p,\theta}(v)) - \frac{1}{\text{area}(B(u))} \int_{B(u)} I_r(d_{p,\theta}(v)) dv .$$

4.1 One-Dimensional Results

Consider the following synthetic example: a one-dimensional signal covering a straight line with position and orientation parameters p_0 and θ_0 is viewed by two cameras, leading to two distorted one-dimensional signals I_l and I_r (Figure 2). The local scene geometry is specified by $p_0 = 0$ and $\theta_0 = .15$ radians. The width of the portion of image being viewed is .69 units in the real scene. The widths of the perspective projections are .15 units in the left image, and .14 units in the right image.

The following three examples compare ρ and ρ_0 for repetitive texture, for the extreme case with only scale disparity, and for non-repetitive texture.

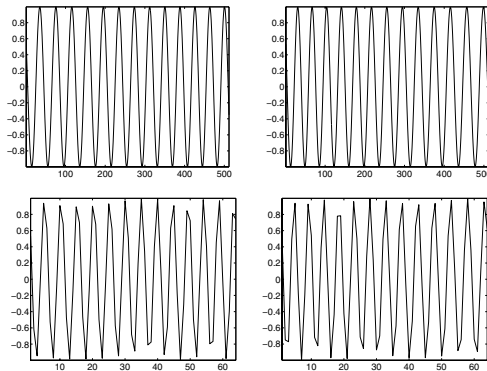


Fig. 2. Two stereo “images” I_l (left) and I_r (right), high resolution (top), and low resolution (bottom).

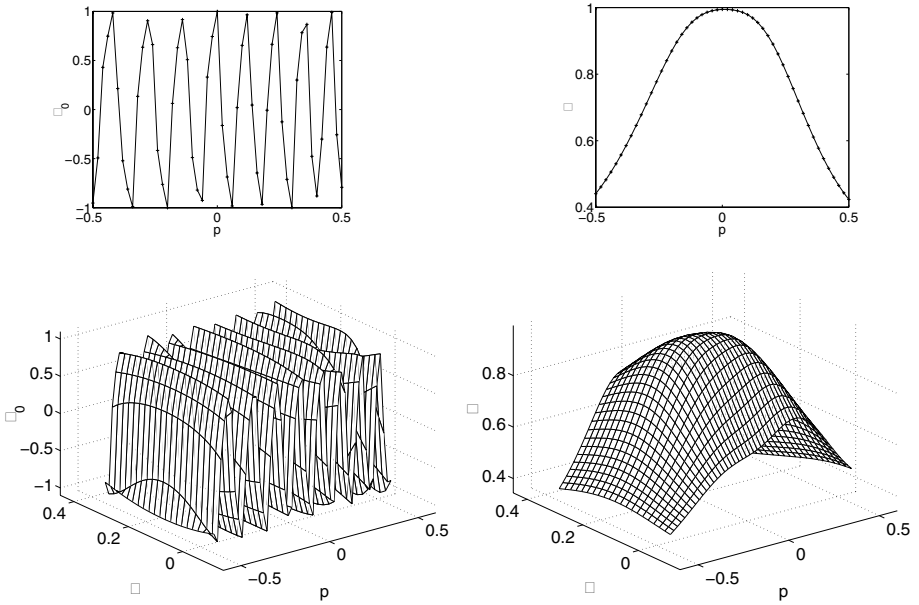


Fig. 3. Correlation for the high-resolution stereo pair of Figure 2 (Example 1, modality 1). Left: classical correlation measure $\rho_0(u_0, p, \theta_0)$ for fixed $\theta = \theta_0$ (top) and $\rho_0(u_0, p, \theta)$ (bottom). Right: new correlation measure $\rho(u_0, p, \theta_0)$ for fixed $\theta = \theta_0$ (top) and $\rho(u_0, p, \theta)$ (bottom). The new correlation ρ is smoother than ρ_0 and displays a unique local maximum at the correct value $(p, \theta) = (p_0, \theta_0)$.

Example 1: repetitive texture. A 1D repetitive texture stereo pair is displayed in Figure 2 both with a high resolution (512 pixels) and with a low resolution (64 pixels).

We display ρ and ρ_0 measured at a fixed position u_0 in the middle of the left image for $-0.05 \leq \theta \leq 0.35$ and $-0.5 \leq p \leq 0.5$ in three different modalities:

1. High-resolution images I_l and I_r with 512 pixels (Figure 3);
2. Low-resolution images I_l and I_r with 64 pixels (Figure 4);
3. High-resolution images I_l and I_r corrupted by two distinct realizations of an additive white Gaussian noise with variance equal to 1/15 of the standard deviation of I_l (Figure 6).

The wavelet used is the one-dimensional Gabor wavelet

$$\psi(x) = \exp(-x^2) \exp(-i\xi x) .$$

We select 5 scales $s_i = 0.05 \times (1.1)^i$ for $i = 0, \dots, 4$, which are relatively coarse compared to the width of image I_l which is .15. The width of the averaging window $B(u_0)$ is .015.

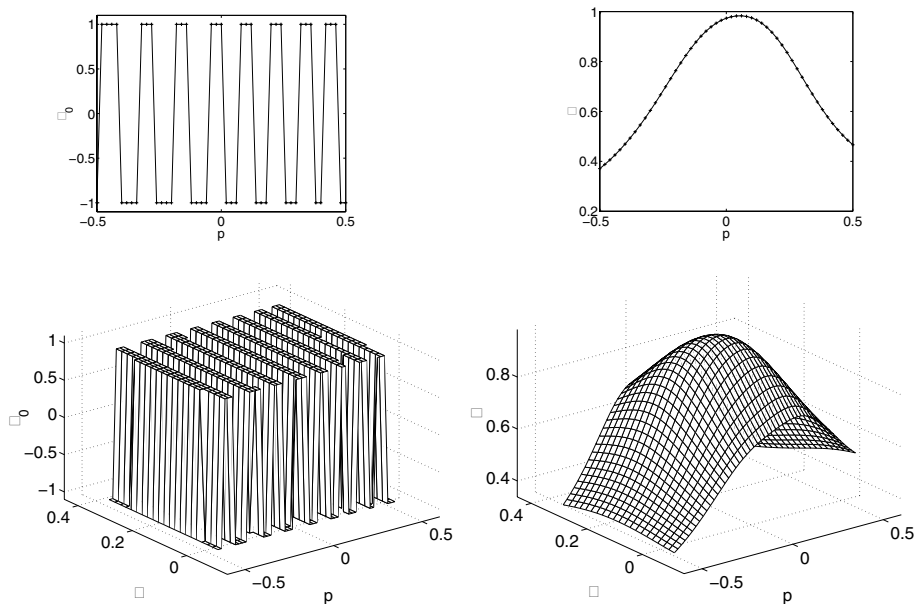


Fig. 4. Correlation for the low-resolution stereo pair of Figure 2 (Example 1, modality 2). Left: classical correlation measure $\rho_0(u_0, p, \theta_0)$ for fixed $\theta = \theta_0$ (top) and $\rho_0(u_0, p, \theta)$ (bottom). Right: new correlation measure $\rho(u_0, p, \theta_0)$ for fixed $\theta = \theta_0$ (top) and $\rho(u_0, p, \theta)$ (bottom). The classical correlation ρ_0 is blocky because of the low resolution, whereas ρ is smooth, with a unique local maximum (which is not as precise as in the first modality because of the low resolution).

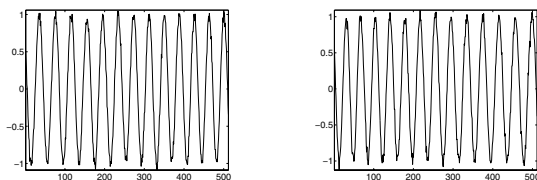


Fig. 5. Stereo pair I_l (left) and I_r (right) corrupted by additive Gaussian noise.

Example 2: scale disparity only. We demonstrate that our correlation measure accounts for a stereoscopic depth perception experiment [12], in which there is a spatial frequency disparity between images I_l and I_r , but no position disparity. Images I_l and I_r in Figure 7 are created by projecting two independent realizations of a colored noise according to a stereoscopic geometry with parameters $(p_0, \theta_0) = (0, .15)$. The numerical results in Figure 8 show that the classical correlation measure ρ_0 , based on position disparity alone, has many lo-

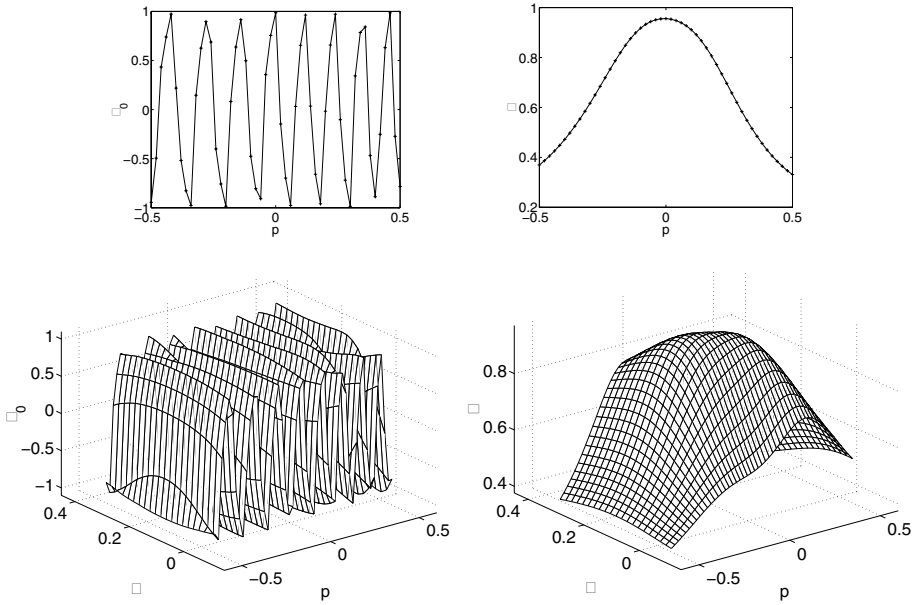


Fig. 6. Correlation for the noisy stereo pair of Figure 5 (Example 1, modality 3). Left: classical correlation measure $\rho_0(u_0, p, \theta_0)$ for fixed $\theta = \theta_0$ (top) and $\rho_0(u_0, p, \theta)$ (bottom). Right: new correlation measure $\rho(u_0, p, \theta_0)$ for fixed $\theta = \theta_0$ (top) and $\rho(u_0, p, \theta)$ (bottom). Observe the smoothness of ρ compared to ρ_0 .

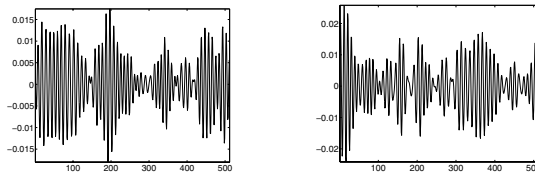


Fig. 7. Two stereo “images” I_l (left) and I_r (right) containing scale disparity information, but no position disparity information.

cal maxima, whereas the correlation measure ρ which is based on scale disparity information displays a unique local maximum at the correct position p_0 .

The wavelets and scales used in this example are the same as in Example 1.

Example 3: non-repetitive texture. Finally, we show the advantage of the new correlation measure ρ over the classical one ρ_0 in the case of a stereo pair with no repetitive texture, but high-frequency oscillations, displayed in Figure 9. Because of the band-pass filtering performed by the Gabor wavelets, the correlation measure ρ is smoother than the classical one, as displayed in Figure

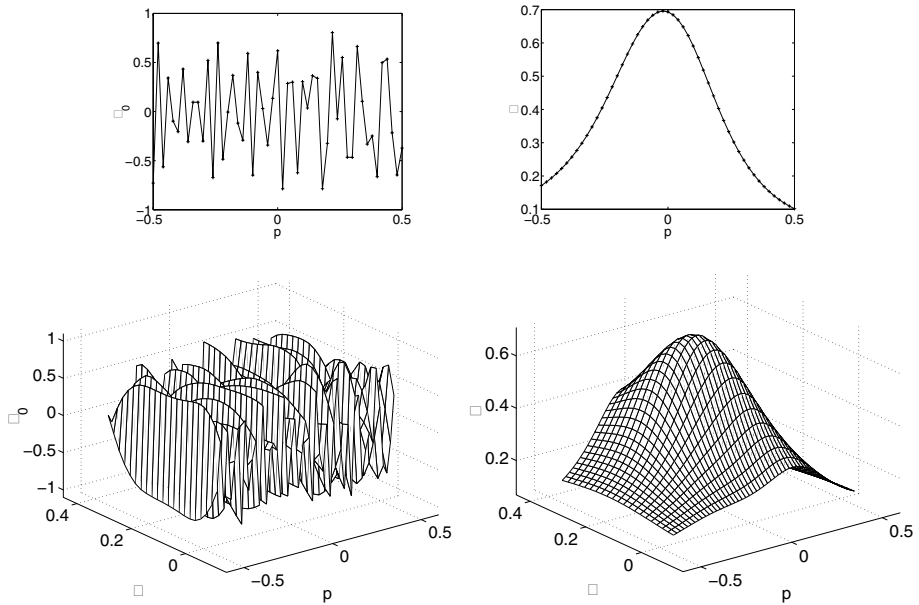


Fig. 8. Left: with only scale disparity information (stereo pair of Figure 7), the classical correlation measure ρ_0 fails. Right: on the other hand, the new correlation measure ρ displays a unique local maximum at $p = p_0$ and $\theta = \theta_0$.

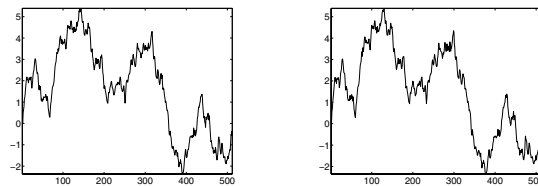


Fig. 9. A stereo pair I_l (left), I_r (right) with non-repetitive texture.

10. Note that this improvement could also have been obtained by applying a low-pass filter to the stereo image pair, before computing the classical correlation. However, this example is noteworthy because it shows the consistency of the wavelet-based correlation measure, which can be applied successfully to different types of images.

In this example we used a set of finer scales than in the two previous examples: $s_i = 0.015 \times (1.1)^i$, for $i = 0, \dots, 4$.

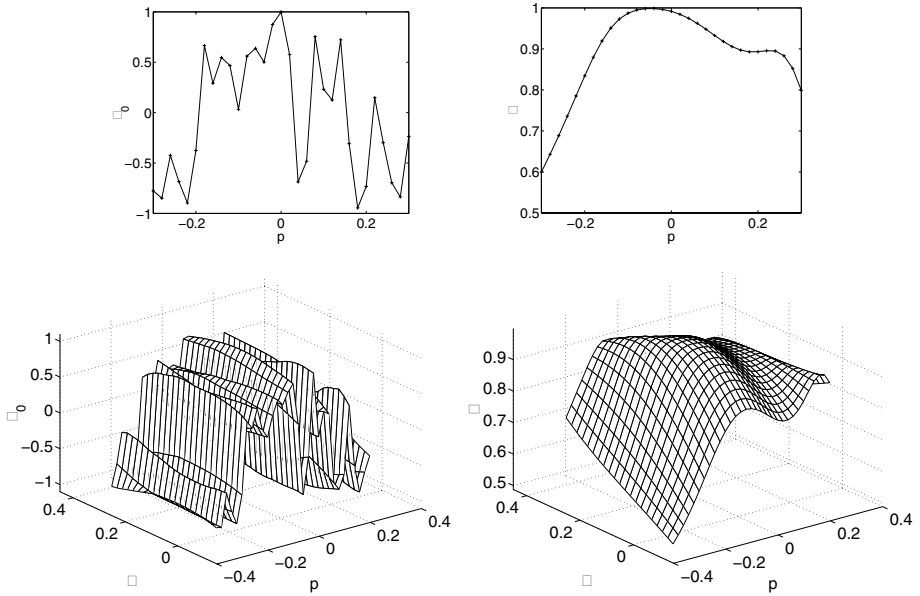


Fig. 10. Correlation for the non-repetitive stereo pair of Figure 9. Left: classical correlation measure $\rho_0(u_0, p, \theta_0)$ for fixed $\theta = \theta_0$ (top) and $\rho_0(u_0, p, \theta)$ (bottom). Right: new correlation measure $\rho(u_0, p, \theta_0)$ for fixed $\theta = \theta_0$ (top) and $\rho(u_0, p, \theta)$ (bottom). Right: $\rho(u_0, p, \theta_0)$. Even in the absence of repetitive texture, ρ is smoother than ρ_0 and does not display spurious local maxima.

4.2 Choice of Correlation Ratio

We found the correlation measure (3) to be more discriminant than the widespread correlation measure

$$\frac{(\text{Feat}_l, \text{Feat}_r)}{(\text{Feat}_l, \text{Feat}_l)^{1/2} (\text{Feat}_r, \text{Feat}_r)^{1/2}} \tag{8}$$

The above ratio (8) is equal to one as soon as Feat_l and Feat_r are *collinear* in feature space, whereas (3) is not equal to one unless Feat_l and Feat_r are *equal* in feature space. The possible advantage of (8) over (3) could be its immunity to shading variations between the images, but we observed that shading variations between images bring the maximum value of (3) down from one, without reducing its high contrast.

4.3 Two-Dimensional Results

Two-dimensional stereoscopic pairs created synthetically by projecting a planar image at position p_0 and with orientation θ_0 onto two cameras in the simplified

geometry of Figure 3. We compare the classical correlation measure ρ_0 to the wavelet-based correlation measure ρ , for repetitive texture (the metallic panel in Figure 11), and non-repetitive texture (the dog hair in Figure 13).

The correlation measure ρ is computed using two-dimensional Gabor wavelets

$$\psi_{u,s}(x) = \det S^{-1} \exp\left(-\|S^{-1}(x-u)\|^2\right) \exp(-i\xi \cdot (S^{-1}(x-u))) .$$

The selected scaling matrices are of the form

$$S_i = (1.05)^i \times \begin{pmatrix} 0.05 & 0 \\ 0 & 0.05 \end{pmatrix} \quad \text{for } i = 0, \dots, 4.$$

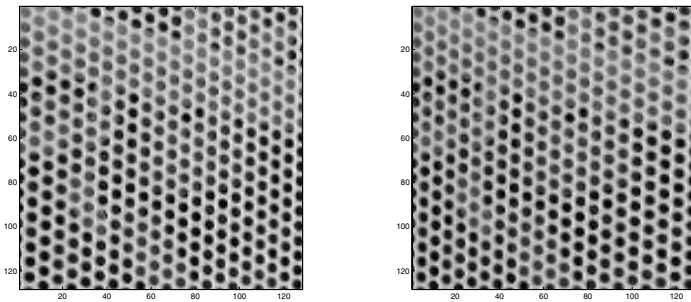


Fig. 11. A synthetic stereo pair, created from a photograph of a metallic panel.

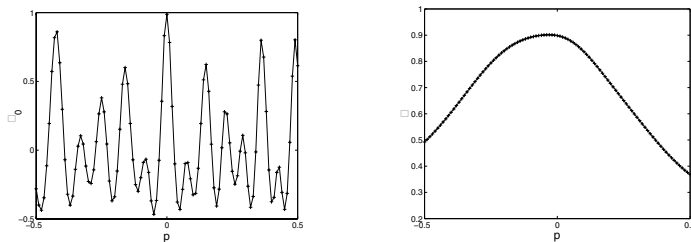


Fig. 12. Left: $\rho_0(u_0, p, \theta_0)$ computed with the stereo pair of Figure 11. Right: $\rho(u_0, p, \theta_0)$. The correct value for p is $p_0 = 0$.

The numerical results displayed in Figures 12 and 14 for fixed $\theta = \theta_0$ show that the new correlation measure is smoother than the classical area-based correlation, and displays a unique local maximum at the correct position $p = p_0$, for repetitive texture and non-repetitive texture alike.

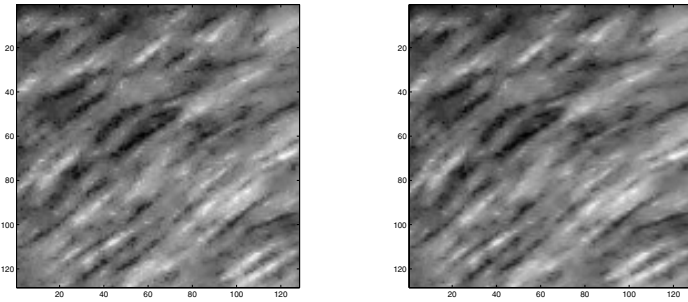


Fig. 13. Synthetic stereo pair, created from a photograph of dog hair.

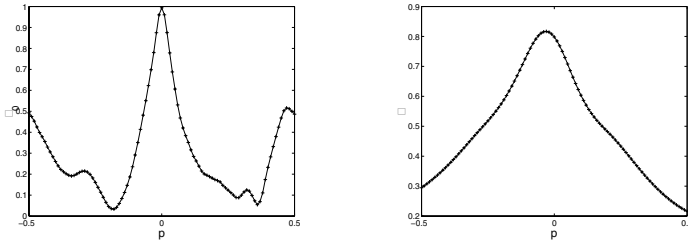


Fig. 14. Left: $\rho_0(u_0, p, \theta_0)$ computed with the stereo pair of Figure 13. Right: $\rho(u_0, p, \theta_0)$. The correct value for p is $p_0 = 0$.

Conclusion

We have introduced a new correlation measure for stereopsis based on wavelet coefficients of images, which has several interesting properties: it lifts the ambiguity on disparity measurement due to the presence of repetitive texture, it shows good performance at low resolution, in the presence of noise, as well as for non-repetitive texture with high-frequency components. Moreover, it is able to reproduce results on stereoscopic depth perception from scale disparity, in the absence of any position disparity. These promising results indicate that, if incorporated in a stereoscopy algorithm which deals with occlusions, this correlation measure could significantly improve its performance. Further work includes testing the algorithm on real data, and investigating the influence of surface curvature.

A Disparity as a Function of Local Surface Position and Orientation

We give the expressions of $d(u)$ and $J(u)$ as a function of p and θ in the case of a locally planar surface element. Let the origin of the Cartesian coordinates

be at the midpoint between the two optical centers, and let the line joining the two optical centers define the x -axis. We assume both cameras to have the same focal length, and we denote by c the half-length between the two optical centers. We suppose that the viewing angles of the cameras with respect to the x -axis are β and $-\beta$. In the simplified geometry of Figure 3, the position and orientation of a local scene element are described by

- p , the distance between the surface tangent and a reference point O with coordinates (x_O, y_O) ,
- θ , the angle between the surface tangent and the x -axis.

If the scene element is locally flat (i.e. neglecting its curvature), the distortion map d is a homography

$$d(u) = \frac{A + Bu}{C + Du} \quad (9)$$

where

$$\begin{aligned} A &= -\cos(\theta - \beta) X' + \cos(\theta + \beta) X \\ B &= -\sin(\theta - \beta) X' - \cos(\theta + \beta) Y' \\ C &= -\sin(\theta + \beta) X - \cos(\theta - \beta) Y' \\ D &= -\sin(\theta - \beta) Y' + \sin(\theta + \beta) Y \end{aligned}$$

with

$$\begin{aligned} X &= \cos \beta (c + x_O - p \sin \theta) - \sin \beta (y_O - p \cos \theta) \\ Y &= \sin \beta (c + x_O - p \sin \theta) + \cos \beta (y_O - p \cos \theta) \end{aligned}$$

and X', Y' are obtained by replacing c by $-c$ and β by $-\beta$ in the expressions of X, Y above. The Jacobian $J(u)$ is calculated by differentiating (9) with respect to u :

$$J(u) = \frac{BC - AD}{(C + Du)^2}. \quad (10)$$

References

1. Clerc, M. and Mallat, S. (2000) Estimating Deformations of Stationary Processes. Research Report no. 192, CERMICS, ENPC.
2. Devernay, F. and Faugeras, O. (1994) Computing Differential Properties of 3D Shapes from Stereoscopic Images without 3D Models. Research Report No. 2304, INRIA, July 1994.
3. Gårding, J. (1992). Shape from Texture for Smooth Surfaces under Perspective Projection. *Journal of Mathematical Imaging and Vision* **2**, pp. 327-350.
4. Jones, D.G. and Malik, J. (1992). A Computational Framework for Determining Stereo Correspondence from a Set of Linear Spatial Filters. Proc. ECCV'92, pp. 395-410.
5. Jones, D.G. and Malik, J. (1992). Determining Three-Dimensional Shape from Orientation and Spatial Frequency Disparities. Proc. ECCV'92, pp. 661-669.
6. Malik, J. and Rosenholtz, R. (1997). Computing Local Surface Orientation and Shape From Texture for Curved Surfaces. *Int. J. of Computer Vision* **23-2**, pp. 149-168.

7. Mallat, S. (1997). *A wavelet tour of signal processing*. Academic Press.
8. Manmatha, R. (1994). Measuring the Affine Transform using Gaussian Filters. *Proc. 3rd European Conf. on Computer Vision*, pp. 159-164, Stockholm, Sweden.
9. Pan, H.P. (1996). General Stereo Image Matching using Symmetric Complex Wavelets. *Proceedings of SPIE Wavelet Applications in Signal and Image Processing IV*.
10. Perrin, J., Torr sani, B., and Fuchs, P. (1999). A Localized Correlation Function for Stereoscopic Image Matching. *Traitement du Signal* **16-1**.
11. Schaffalitzky, F., Zisserman, A. (2001). Viewpoint invariant Texture Matching and Wide Baseline Stereo. *Proceedings of ICCV*.
12. Tyler, C.W. and Sutter, E.E. (1979). Depth from spatial frequency difference: an old kind of stereopsis? *Vision Research* 19:859-865.
13. Zitnick, C.L. and Kanade, T. (2000). A Cooperative Algorithm for Stereo Matching and Occlusion Detection. *IEEE Trans. Pat. Anal. and Mach. Intell.* **22-7**, pp. 675-684.