# Binocular model for figure-ground segmentation in translucent and occluding images

**David Vernon,** MEMBER SPIE
CAPTEC Ltd.
Malahide
Co. Dublin, Ireland
E-mail: david@captec.ie

**Abstract.** A Fourier-based solution to the problem of figure-ground segmentation in short baseline binocular image pairs is presented. Each image is modeled as an additive composite of two component images that exhibit a spatial shift due to the binocular parallax. The segmentation is accomplished by decoupling each Fourier component in one of the resultant additive images into its two constituent phasors, allocating each to its appropriate object-specific spectrum, and then reconstructing the foreground and background using the inverse Fourier transform. It is shown that the foreground and background shifts can be computed from the differences of the magnitudes and phases of the Fourier transform of the binocular image pair. While the model is based on translucent objects, it also works with occluding objects. © *2002 Society of Photo-Optical Instrumentation Engineers.* [DOI: 10.1117/1.1504722]

## 1 Introduction

Many approaches have been proposed for segmentation of images using real or apparent motion arising from either binocular stereo or camera motion (e.g., Refs. 1–9). Many of these techniques accomplish the segmentation by implicitly or explicitly matching points or regions to establish stereo correspondence[10]* and then perform segmentation based on some similarity predicate (such as range or the parameters of affine motion). Such matching can be computationally expensive and it normally requires that the images exhibit a certain amount of texture to drive the matching process. In addition, few are able to deal with situations where there are multiple motions in each local region such as arise in transparent or translucent objects where the images are effectively additive.[11] Independent component analysis has been successfully applied to this problem,[12] but this approach requires that the images are statistically independent. It has been noted that frequency-domain techniques have advantages in dealing with complex imagery of this nature.[13]† One such frequency-domain technique for the separation of two images that have been additively combined (e.g., reflections in a window superimposed on a background scene) was presented in Ref. 14 and it has also been used to solve the figure/ground segmentation problem where the foreground occludes the background. This approach is based on a Fourier analysis of the composite images and the appearance of the images (spectrally or spatially) is not a limitation. However, the technique requires

four samples of the combined images to achieve the figure/ground segmentation. In this paper, we present a similar technique that requires only two samples. Computationally, this means the approach is efficient, requiring two fast Fourier transforms (FFTs), two inverse FFTs, and a linear time-complexity analysis and decomposition of constituent phasors in the spatial frequency domain.

## 2 Statement of the Problem

The requirement of four images in Ref. 14 arises directly from the formulation of the problem itself. Consider a composite image $f^{p_n}(x,y)$ acquired at position $p_n$ in a spatially ordered image sequence

$$f^{p_n}(x,y) = f_1^{p_n}(x,y) + f_2^{p_n}(x,y), \tag{1}$$

where $f_1^{p_n}(x,y)$ and $f_2^{p_n}(x,y)$ are the unknown foreground and background additive component images at position $p_n$. Assuming a uniform translatory motion or shift between each position $p_n$,

$$f_i^{p_n}(x,y) = f_i^{p_0}(x - n\,\delta x_i, y - \delta y_i), \tag{2}$$

where $(\delta x_i, \delta y_i)$ is the incremental spatial shift of the $i$th component image, the goal is to to recover each individual image $f_i^{p_0}(x,y)$.

The shift property of the Fourier transform states that the Fourier transform of a shifted function $f(x - n\,\delta x, y - n\,\delta y)$ is given by[15]

---

*Variation techniques for multiframe stereo reconstruction of smooth shapes are included in Ref. 10, which are complementary to traditional techniques and do not work in regions with strong texture.

†Reference 13 discusses optical snow dealing with the relative motion of highly complex objects in a scene and argues for the superiority of frequency-domain techniques.

$$\mathcal{F}[f(x-n\,\delta x, y-n\,\delta y)]=\mathcal{F}[f(x,y)]$$
$$\times \exp[-i(\omega_x n\,\delta x + \omega_y n\,\delta y)]. \tag{3}$$

Combining Eqs. (1) and (2) and taking the Fourier tranform, we have

$$\mathsf{F}^{p_n}(\omega_x,\omega_y)=\mathsf{F}_1^{p_0}(\omega_x,\omega_y)$$
$$\times\{\exp[-i(\omega_x n\,\delta x_1+\omega_y n\,\delta y_1)]\}^n$$
$$+\mathsf{F}_2^{p_0}(\omega_x,\omega_y)$$
$$\times\{\exp[-i(\omega_x n\,\delta x_2+\omega_y n\,\delta y_2)]\}^n$$
$$=\mathsf{F}_1^{p_0}(\Delta\Phi_1)^n+\mathsf{F}_2^{p_0}(\Delta\Phi_2)^n, \tag{4}$$

where $\Delta\Phi_i=\exp[-i(\omega_x\delta x_i+\omega_y\delta y_i)]$, a complex variable representing the frequency- and shift-dependent phase change. Consequently, we have four unknowns—$\mathsf{F}_1^{p_0}$, $\mathsf{F}_2^{p_0}$, $\Delta\Phi_1$, and $\Delta\Phi_2$—the two images and their respective incremental displacements. It has been shown that these four unknowns can be identified by solving a set of four simultaneous equations of the form of Eq. (4), each of which models the known additive combination $\mathsf{F}^{p_n}$ of the foreground and the background at displacements $p_0$, $p_1$, $p_2$, and $p_3$.

The purpose of this paper is to show how, in the case of frontoparallel binocular stereo, we can solve for $\mathsf{F}_1^{p_0}$, $\mathsf{F}_1^{p_0}$, $\Delta\Phi_1$, and $\Delta\Phi_2$ with just two images instead of four.

## 3 Solution Using Two Images

We assume that the binocular stereo configuration has a vergence angle of 0 deg, i.e., both cameras are pointing in the same direction so that their optical axes are parallel. This is often referred to as frontoparallel binocular stereo and is one of the most common stereo configurations. With this configuration, the epipolar lines are parallel to the line joining the optical centres of the two images. This means that the displacement of each pixel in the image is horizontal and, thus, $\delta y = 0$. Furthermore, since we have only two images, $n=0,1$. Thus, our model now becomes

$$\mathsf{F}^{p_0}(\omega_x,\omega_y)=\mathsf{F}_1^{p_0}(\omega_x,\omega_y)+\mathsf{F}_2^{p_0}(\omega_x,\omega_y),$$

$$\mathsf{F}^{p_1}(\omega_x,\omega_y)=\mathsf{F}_1^{p_0}(\omega_x,\omega_y)\exp[-i(\omega_x\delta x_1)]$$
$$+\mathsf{F}_2^{p_0}(\omega_x,\omega_y)\exp[-i(\omega_x\delta x_2)].$$

The goal is now to find $\delta x_1$, $\delta x_2$, $\mathsf{F}_1^{p_0}(\omega_x,\omega_y)$, and $\mathsf{F}_2^{p_0}(\omega_x,\omega_y)$.

To solve for these, first we separate the Fourier components into their magnitude and phase components [for the sake of brevity, we will drop the $(\omega_x,\omega_y)$ arguments]:

$$\mathsf{F}^{p_0}=|\mathsf{F}_1^{p_0}|e^{i\phi_1}+|\mathsf{F}_2^{p_0}|e^{i\phi_2}, \tag{5}$$

$$\mathsf{F}^{p_1}=|\mathsf{F}_1^{p_0}|\exp[i(\phi_1-\omega_x\delta x_1)]$$
$$+|\mathsf{F}_2^{p_0}|\exp[i(\phi_2-\omega_x\delta x_2)]. \tag{6}$$

The magnitude of each composite phasor at position $p_0$ is

$$|\mathsf{F}^{p_0}|=||\mathsf{F}_1^{p_0}|e^{i\phi_1}+|\mathsf{F}_2^{p_0}|e^{i\phi_2}|.$$

Rotating $\mathsf{F}^{p_0}$ to align it with one of the the component phasors does not alter its magnitude:

$$|\mathsf{F}^{p_0}|=|\mathsf{F}^{p_0}\exp(-i\phi_1)|$$
$$=||\mathsf{F}_1^{p_0}|e^{i\phi_1}\exp(-i\phi_1)+|\mathsf{F}_2^{p_0}|e^{i\phi_2}\exp(-i\phi_1)|$$
$$=||\mathsf{F}_1^{p_0}|+|\mathsf{F}_2^{p_0}|\exp[i(\phi_2-\phi_1)]|.$$

Similarly, rotating $\mathsf{F}^{p_1}$ to align it with one of the the component phasors:

$$|\mathsf{F}^{p_1}|=|\mathsf{F}^{p_1}\exp[-i(\phi_1-\omega_x\delta x_1)]|$$
$$=||\mathsf{F}_1^{p_0}|\exp[i(\phi_1-\omega_x\delta x_1)]\exp[-i(\phi_1-\omega_x\delta x_1)]$$
$$+|\mathsf{F}_2^{p_0}|\exp[i(\phi_2-\omega_x\delta x_2)]\exp[-i(\phi_1-\omega_x\delta x_1)]|$$
$$=||\mathsf{F}_1^{p_0}|+|\mathsf{F}_2^{p_0}|\exp\{i[\phi_2-\phi_1+\omega_x(\delta x_1-\delta x_2)]\}|.$$

Taking the difference between the magnitudes, we have:

$$|\mathsf{F}^{p_1}|-|\mathsf{F}^{p_0}|=||\mathsf{F}_1^{p_0}|+|\mathsf{F}_2^{p_0}|$$
$$\times\exp\{i[\phi_2-\phi_1+\omega_x(\delta x_1-\delta x_2)]\}|$$
$$-||\mathsf{F}_1^{p_0}|+|\mathsf{F}_2^{p_0}|\exp[i(\phi_2-\phi_1)]|.$$

For some given displacements $\delta x_1$ and $\delta x_2$, there exists a set of spatial frequencies $\omega_x'$ such that $|\mathsf{F}^{p_1}|-|\mathsf{F}^{p_0}|=0$. That is,

$$\exp\{i[\phi_2-\phi_1+\omega_x'(\delta x_1-\delta x_2)]\}=\exp[i(\phi_2-\phi_1)].$$

Hence,

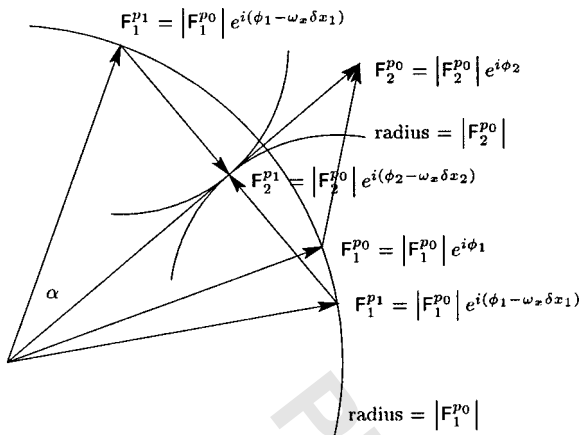$$\exp\{i[\omega_x'(\delta x_1-\delta x_2)]\}=1,$$

and thus,

$$\omega_x'(\delta x_1-\delta x_2)=2n\pi, n=0,1,2,\ldots.$$

Thus,

$$\omega_x'=\frac{2n\pi}{\delta x_1-\delta x_2}.$$

Thus, if we identify the set of spatial frequencies $\omega_x'$ for which $|\mathsf{F}^{p_1}|-|\mathsf{F}^{p_0}|=0$, we can then compute the difference between the two quantities we are seeking to find, i.e.,

$$\delta x_1-\delta x_2=\frac{2n\pi}{\omega_x'}. \tag{7}$$

Vernon: Binocular model for figure-ground segmentation . . .



**Fig. 1** Phase of the resultant $F^{p_0} = F_1^{p_0} + F_2^{p_0}$ will wrap at an angle in the range $2\pi \pm \alpha$, $\alpha = \sin^{-1}(|F_2^{p_0}|/|F_1^{p_0}|)$.
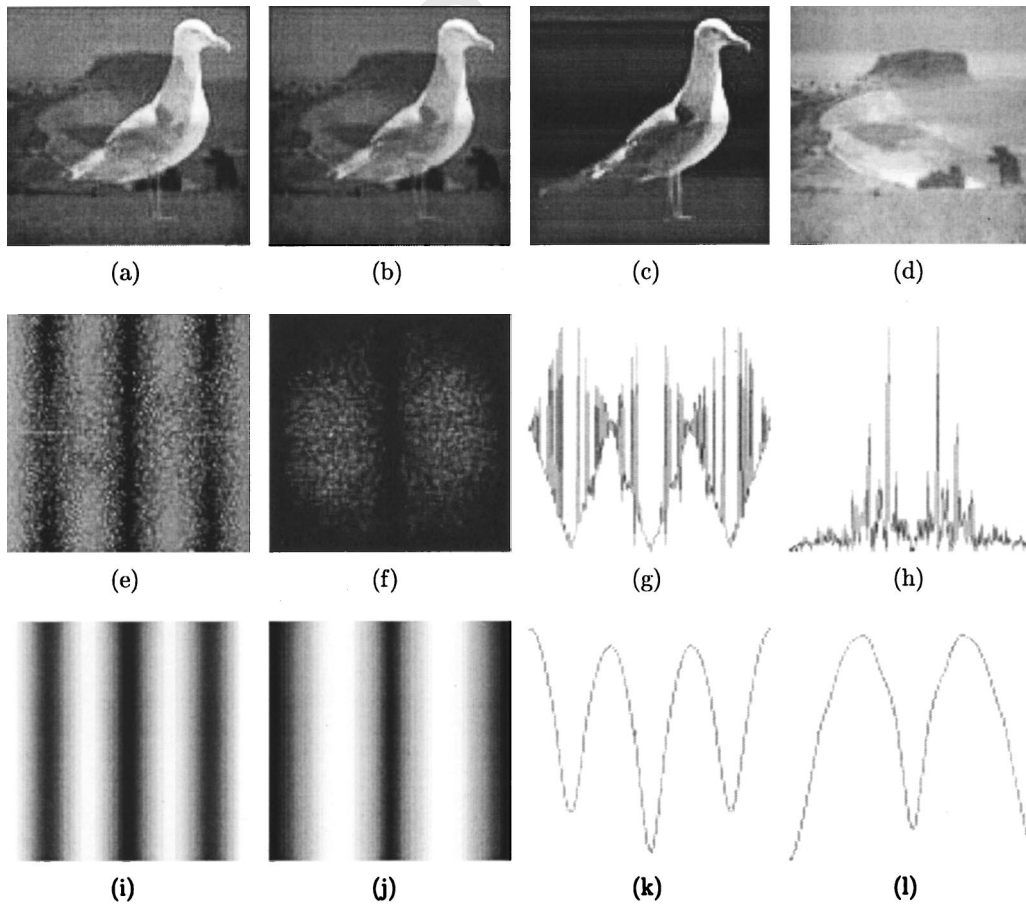
Since this equation depends only on the $\omega_x$ spatial frequencies, these zero components will occur at all $\omega_y$ frequencies. We use this fact to improve the robustness of the technique by summing all the magnitude differences along the $\omega_y$ axis to yield a 1-D signature $F^m$ of the difference of Fourier component magnitudes as a function of $\omega_x$:

$$F^m(\omega_x) = \sum_{\omega_y} |F^{p_1}(\omega_x, \omega_y)| - |F^{p_0}(\omega_x, \omega_y)|.$$

This signature periodically approaches zero at frequencies

$$\omega_x' = \frac{2n\pi}{\delta x_1 - \delta x_2}$$

(see Figs. 2–5 in Sec. 4).



**Fig. 2** Synthetic additive test: (a) left image, (b) right image, (c) segmented foreground, and (d) segmented background. The foreground displacement is 3 pixels and the background displacement is 1 pixel. Computation of displacements from the difference of phases and magnitudes of the resultants: (e) unprocessed phase difference, (f) unprocessed magnitude difference, (g) cross section through the phase difference image, (h) cross section through the magnitude difference image, (i) the phase difference averaged in vertical spatial frequency ($\omega_y$) direction, (j) the averaged magnitude difference, (k) cross section through the phase difference image, and (l) cross section through the magnitude difference image.
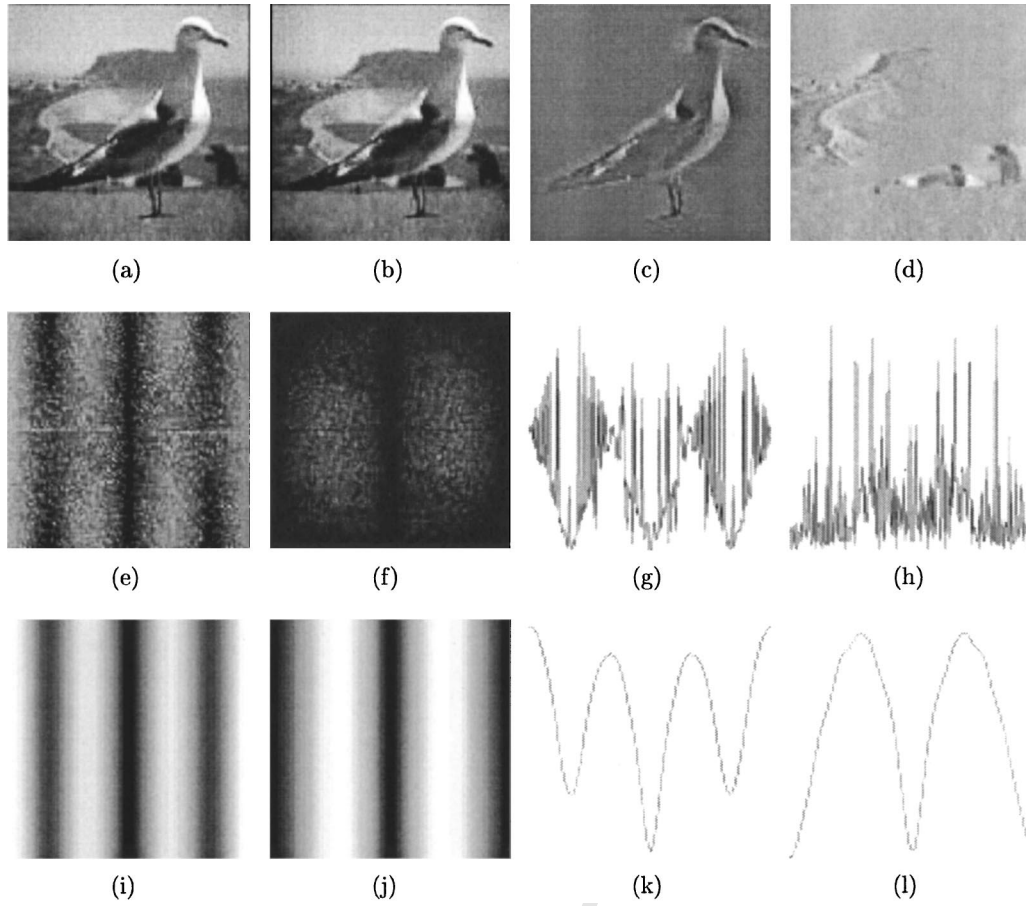
**Fig. 3** Synthetic occlusion test; see Fig. 2 for an explantion of images (a) to (l).

It now remains to compute either $\delta x_1$ or $\delta x_2$. We can do this as follows. Recall Eqs. (5) and (6). Let the phase angle of $\mathsf{F}^{p0}$ and $\mathsf{F}^{p1}$ be written $\angle \mathsf{F}^{p0}$ and $\angle \mathsf{F}^{p1}$, respectively. Without loss of generality, let $|\mathsf{F}^{p0}| \geqslant |\mathsf{F}^{p1}|$. By a geometric construction (refer to Fig. 1) we can see that, for any phase angle $(\phi_2 - \omega_x \delta x_2)$ of Fourier phasor component $\mathsf{F}_2^{p1} = |\mathsf{F}_2^{p0}| \exp[i(\phi_2 - \omega_x \delta x_2)]$, there will be a phase angle $(\phi_1 - \omega_x \delta x_1)$ of Fourier component $\mathsf{F}_1^{p1} = |\mathsf{F}_1^{p0}| \exp[i(\phi_1 - \omega_x \delta x_1)]$ for which $\angle \mathsf{F}^{p0} = \angle \mathsf{F}^{p1}$ if and only if

$$2n\pi - \alpha \leqslant \omega_x \delta x_1 \leqslant 2n\pi + \alpha, \quad n = 0,1, \ldots,$$

where

$$\alpha = \sin^{-1}\left( \frac{|\mathsf{F}_2^{p0}|}{|\mathsf{F}_1^{p0}|} \right).$$

This is true because if $\omega_x \delta x_1$ (the phase change of $\mathsf{F}_1^{p0}$) falls outside this range, the phasor $\mathsf{F}_2^{p1}$ cannot intersect the line directed along $\angle \mathsf{F}^{p0}$. That is, $\angle \mathsf{F}^{p0} = \angle \mathsf{F}^{p1}$ for some value of $\omega_x \delta x_1$, if and only if $2n\pi - \alpha \leqslant \omega_x \delta x_1 \leqslant 2n\pi + \alpha$.

Let $\omega_x^p$ denote the frequencies at which $\angle \mathsf{F}^{p0} = \angle \mathsf{F}^{p1}$. Then,

$$\frac{2n\pi - \alpha}{\delta x_1} \leqslant \omega_x^p \leqslant \frac{2n\pi + \alpha}{\delta x_1}. \tag{8}$$

Note that this condition does not depend on $\omega_y$. If $|\mathsf{F}_1^{p0}|$ and $|\mathsf{F}_2^{p0}|$ are not correlated, as is the case in real images, then the average frequency $\bar{\omega}_x^p = \Sigma_{\omega_y} \omega_x^p$ for which $\angle \mathsf{F}^{p0} = \angle \mathsf{F}^{p1}$ will be a good estimate of the central frequency in the range given in Eq. (8). That is,

$$\bar{\omega}_x^p \approx \frac{2n\pi}{\delta x_1}. \tag{9}$$

Algorithmically, we compute a phase difference image $\angle \mathsf{F}^{p1}(\omega_x,\omega_y) - \angle \mathsf{F}^{p0}(\omega_x,\omega_y)$ and sum along the $\omega_y$ axis to yield a 1-D signature $\mathsf{F}^p(\omega_x)$:

$$\mathsf{F}^p(\omega_x) = \sum_{\omega_y} \angle \mathsf{F}^{p1}(\omega_x,\omega_y) - \angle \mathsf{F}^{p0}(\omega_x,\omega_y).$$

This signature approaches zero at frequencies $\bar{\omega}_x^p \approx (2n\pi \delta x_1)$, $n = 0,1,\ldots$. To identify $\delta x_1$, we find $\bar{\omega}_x^p$ and compute $\delta x_1$ from Eq. (9) and $\delta x_2$ fom Eq. (7).

Knowing $\delta x_1$ and $\delta x_2$, we can now compute $\Delta \Phi_1(\omega_x,\omega_y)$ and $\Delta \Phi_2(\omega_x,\omega_y)$ from
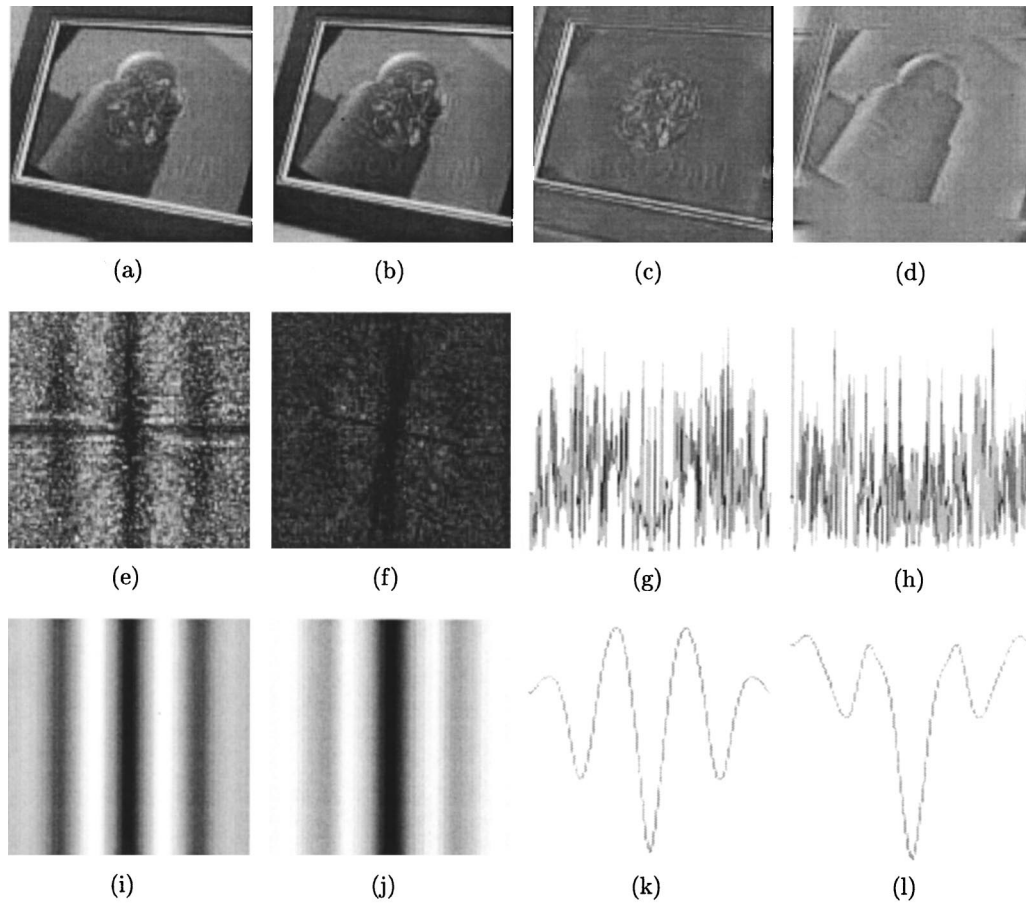
**Fig. 4** Real additive test; see Fig. 2 for an explantion of images (a) to (l).

$$\Delta \Phi_j(\omega_x, \omega_y) = \exp(-i\omega_x \delta x_j),$$

and, in principle, we can do so for all frequencies $(\omega_x, \omega_y)$.

It then remains to solve for $F^{p0}(\omega_x, \omega_y)$ and $F^{p0}(\omega_x, \omega_y)$, which we can get directly from Eq. (4) by solving two simultaneous equations of the form

$$F^{p0}(\omega_x, \omega_y) = F_1^{p0}(\omega_x, \omega_y) + F_2^{p0}(\omega_x, \omega_y), \tag{10}$$

$$F^{p1}(\omega_x, \omega_y) = F_1^{p0}(\omega_x, \omega_y)\Delta\Phi_1(\omega_x, \omega_y)$$
$$+ F_2^{p0}(\omega_x, \omega_y)\Delta\Phi_2(\omega_x, \omega_y), \tag{11}$$

for all spatial frequencies $(\omega_x, \omega_y)$. The images of the figure and ground can then be reconstructed in the same manner described in Ref. 14 by taking the inverse Fourier transform of $F_1^{p0}(\omega_x, \omega_y)$ and $F_2^{p0}(\omega_x, \omega_y)$.

## 4 Application of the Technique

Figures 3–5 show the result of applying the technique to four data sets:

1. a synthetic scene comprising two additive images
2. a synthetic scene comprising two occluding objects
3. a real scene comprising two occluding binocular images

4. a real scene comprising two additive binocular images.

The necessity to average in the direction of the vertical spatial frequency $\omega_y$ to make the minima more accessible can be clearly seen in all of these figures.

## 5 Discussion

We presented a significant simplification of the theory developed in Ref. 14, reducing the number of images or samples required to accomplish the segmentation from four to two. However, one should recognize that this simplification applies only in the case of fronto-parallel binocular imaging (either using a single translating camera or two static cameras). It does not apply in the case of two independently moving objects such, as was discussed in Ref. 14. The approach described in this paper depends on the fact that the phase shifts of both constituent phasors is confined to one axis (in this case the $x$ axis), thereby removing one component of displacement (the $y$ component in this case). This then leaves just two independent unknown quantities: the two $x$ components of the displacement of each object. These conditions are satisfied if and only if the displacements of the two images/objects are parallel, which is indeed the case for frontoparallel binocular vision. That said, it is plausible to extend the approach somewhat, relaxing the assumption that the variation is aligned with the
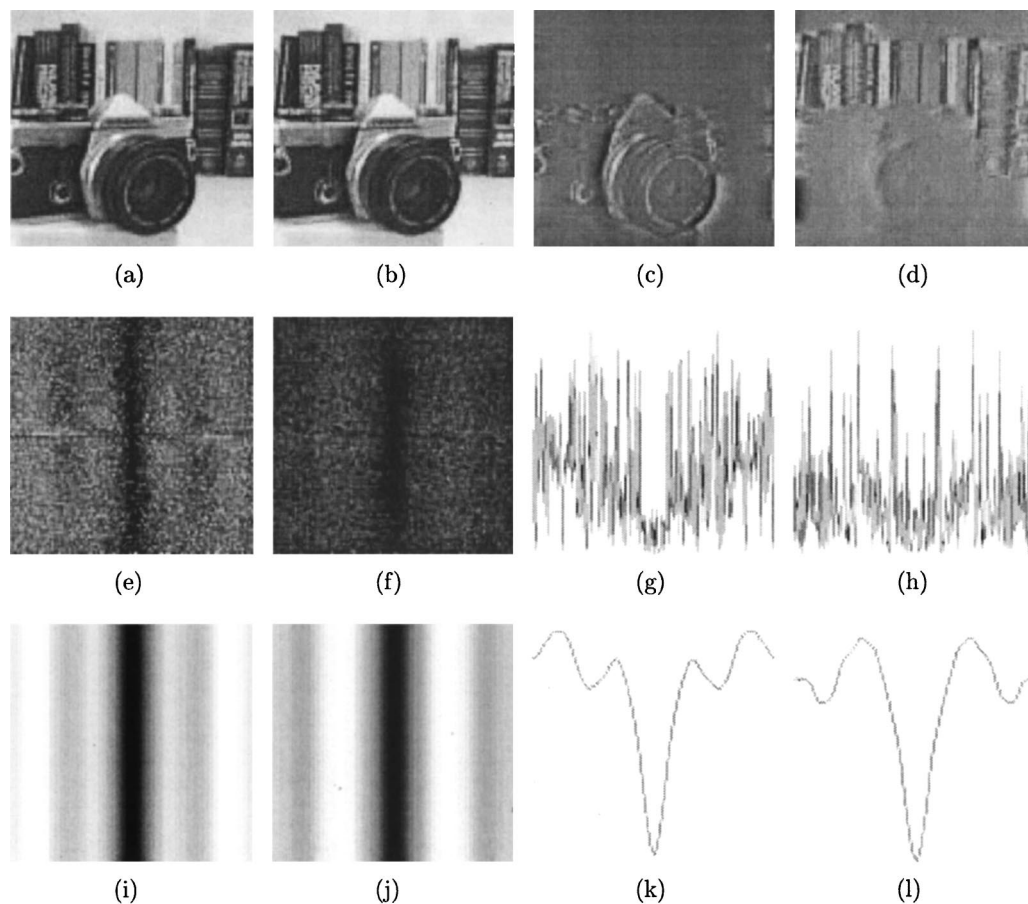
**Fig. 5** Real occlusion test; see Fig. 2 for an explantion of images (a) to (l).

$x$ axis, by allowing an arbitrary but common axis of variation: in this case, we must first find the axis of variation (typically by identifying the orthogonal axis along which there is no variation in phase difference). Unfortunately, the general case where the objects move independently will not yield to this line of attack because the phase differences will vary along every direction.

We said in the preceding that $\mathsf{F}^{P_0}(\omega_x,\omega_y)$, $\mathsf{F}^{P_0}(\omega_x,\omega_y)$, $\Delta\Phi_1(\omega_x,\omega_y)$, and $\Delta\Phi_2(\omega_x,\omega_y)$ can be computed for all spatial frequencies. However, this is not strictly correct and there are some spatial frequencies for which we cannot solve. These are the frequencies at $\omega_x=2n\pi, n=0,1,\dots$. In these cases $\Delta\Phi_{1,2}(\omega_x,\omega_y)$ $=\exp(-i\omega_x\delta x_{1,2})=1$ and Eqs. (10) and (11) become degenerate. Consequently, all such frequencies (including the dc component) are omitted from the inverse Fourier transform, and this accounts for the alteration in the appearance of the reconstructed images compared with the original in Figs. 2 and 3. Furthermore, in the case of occluded image segmentation, frequencies close to integer multiples of $2\pi$ (i.e., $\omega_x\approx2n\pi, n=0,1,\dots$) are ill-conditioned and are attenuated before taking the inverse Fourier transform (see Ref. 14 for details). Note that the exent of this ill-conditioning increases with object displacement and, consequently, the quality of the segmentation disimproves with longer binocular baselines. Thus, the approach is best

suited to short-baseline stereo where the occluded signal is less significant.

Current work is directed at incorporating an explicit model for occlusion using, for example, the recent advances described in Ref. 16. In addition, we are aiming to extend the approach to segment more than two layers; this will require an increase in the number of observations or input images.

## Acknowledgments

## References

1. P. Bouthemy and P. Lalande, ''Recovery of moving object masks in an image sequence using local spatiotemporal contextual information,'' *Opt. Eng.* **32**(6), 1205–1212 (1993).
2. J. Y. A. Wang and E. H. Adelson, ''Representing moving images with layers,'' *IEEE Trans. Image Process.* **3**(5), 625–638 (1994).
3. J. Santos-Victor and G. Sandini, ''Uncalibrated obstacle detection using normal flow,'' *Mach. Vision Appl.* **9**, 130–137 (1996).
4. D. Tzovaras, N. Grammalidis, and M. G. Strintzis, ''3-D camera motion estimation and foreground/background separation for stereoscopic image sequences,'' *Opt. Eng.* **36**(2), ■■■ (1997).
5. J. Santos-Victor and G. Sandini, ''Visual behaviors for docking,'' *Comput. Vis. Image Underst.* **67**(3), 223–238 (1997).
6. G. D. Borshukov, G. Bozdagi, Y. Altunbasak, and A. M. Tekalp, ''Motion segmentation by multistage affine classification,'' *IEEE Trans. Image Process.* **6**(11), 1591–1594 (1997).
7. M. M. Chang, A. M. Tekalp, and M. I. Sezan, ''Simultaneous motion

estimation and segmentation," *IEEE Trans. Image Process.* **6**(9), 1326–1333 (1997).

8. S.-W. Lee, J. G. Choi, and S.-D. Kim, "Scene segmentation using a combined criterion of motion and intensity," *Opt. Eng.* **36**(8), 2346–2352 (1997).

9. M. Irani and P. Anandan, "A unified approach to moving object detection in 2D and 3D scenes," *IEEE Trans. Pattern Anal. Mach. Intell.* **20**(6), 577 (1998).

10. A. J. Yezzi and S. Soatto, "Stereoscopic segmentation," in *Proc. Int. Conf. Computer Vision*, pp. 59–66 (July 2001).

11. J. R. Bergen, P. J. Burt, R. Hingorani, and S. Peleg, "A three-frame algorithm for estimating two-component image motion," *IEEE Trans. Pattern Anal. Mach. Intell.* **14**(9), 886–895 (1992).

12. J.-F. Cardoso, "Blind signal separation: statistical principles," *Proc. IEEE* **86**(10), 2009–2025 (1998).

13. M. S. Langer and R. Mann, "Dimensional analysis of image motion," *Proc. Int. Conf. on Computer Vision*, pp. 155–162 (July 2001).

14. D. Vernon, "Decoupling Fourier components of dynamic image sequences: a theory of signal separation, image segmentation, and estimation of optical flow," in *Proc. 5th Eur. Conf. on Computer Vision—ECCV'98*, Vol. 2, pp. 69–85 Springer-Verlag (1998).

15. A. D. Poularikas, Ed. *The Transforms and Applications Handbook*, IEEE Press, (1996).

16. S. S. Beauchemin and J. L. Barron, "The frequency structure of one-dimensional occluding image signals," *IEEE Trans. Pattern Anal. Mach. Intell.* **22**(2), 200–206 (2000).

**David Vernon** graduated in engineering from the University of Dublin, Trinity College, Ireland, in 1979 and began his professional career as a software engineer with Westinghouse Electric. He was appointed as a lecturer in the Department of Computer Science, Trinity College, in 1983, completing his PhD in robot vision in 1985. From 1991 to 1993, he worked as a scientific officer in the European Commission and he was appointed to the Chair of Computer Science at the National University of Ireland, Maynooth, in 1995. He now works for Science Foundation Ireland and he is also the coordinator of ECVision, the European Network on Cognitive Computer Vision at CAPTEC Ltd, Ireland. He has authored over seventy-five papers and three books on computer vision and image processing. He is a past Fellow of Trinity College and he is a Chartered Engineer.