

# Stereo Vision, Models of

## Synonyms

Models of stereopsis, models of 3D vision.

## Definition

Stereo vision refers to the perception of depth based on slight disparities between the images seen by the two eyes, due to their different location in space. This article reviews models of stereo vision based on real nervous systems. Stereo vision algorithms have also been developed independently within the machine vision literature, sometimes biologically inspired.

## Detailed description

Stereo vision has two main steps: (1) Extracting disparity from the retinal images, and (2) Perceiving depth structure from disparity. Although psychophysics experiments have probed both aspects, much more is known about the neuronal mechanisms supporting the extraction of disparity, whereas little is known about how and where in the brain disparity is converted into a perception of depth.

## Extracting disparity from the retinal images

In primates, the first step is believed to begin in primary visual cortex (V1). V1 is the first visual area to contain neurons which receive inputs from both eyes, many of which are tuned to binocular disparity (Poggio, Motter et al. 1985). The most influential model of disparity tuning in V1 is the stereo energy model, put forward by Izumi Ohzawa, Greg DeAngelis and Ralph Freeman (1990), and based in turn on similar models in the motion literature (Adelson and Bergen 1985).

## Stereo energy model: basic unit

The basic idea behind the stereo energy model is that of a binocular neuron with a receptive field in each eye (Figure 1). The input from each eye depends on the inner product of each eye's image with the corresponding receptive field. For example, the input from the left eye is

$$L = \int_{\text{retina}} dx dy I_L(x,y) \rho_L(x,y)$$

The function  $I_L(x,y)$  represents the retinal image in the left retina, with luminance expressed relative to the mean, e.g. +1 = white, -1 = black, 0=gray. The function  $\rho_L(x,y)$  represents the classical receptive field in the left retina, with 0=no response (i.e. areas of the retina that are outside the receptive field), positive values corresponding to ON regions (areas where light stimulation tends to increase the cell's firing rate, and dark suppresses it) and negative values to OFF regions (where dark features tend to increase the cell's firing and light suppresses it).

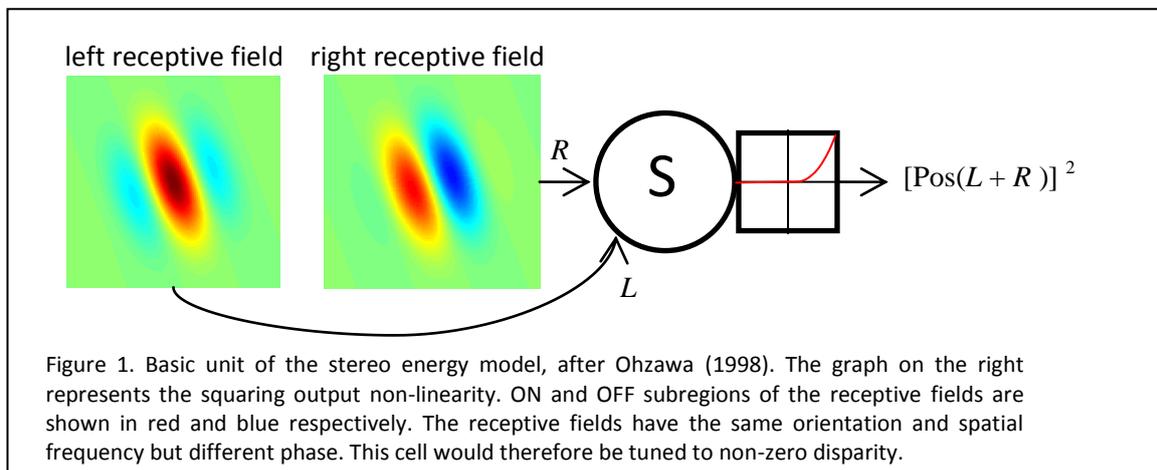
The energy model assumes that inputs from each eye are summed linearly. If the result is negative, the neuron is inhibited and does not fire. If positive, the neuron fires at a rate proportional to the square of the inputs:

$$S = [L + R]^2$$

Equation 1

This model neuron is tuned to disparity both in monocularly visible features, e.g. the bar stimuli used in Ohzawa et al. (1990), and also in “cyclopean” stimuli such as dynamic random-dot patterns (Julesz 1971). A detailed analysis shows that the cyclopean disparity tuning is due to the output non-linearity (Read 2005). A unit which simply summed left and right inputs linearly ( $L+R$ ) would be tuned to the disparity of a bar stimulus, but only as a side-effect of its tuning to the position of monocular features. It would respond equally well, on average, to broadband noise images with any disparity. random-dot but not to cyclopean disparity. Thus, a non-linearity is essential to produce cyclopean disparity tuning, i.e. to ensure that the cell responds more on average to broadband noise images with its preferred disparity. The threshold-and-squaring non-linearity in Equation 1 is used in many models of V1 neurons because it describes their responses well. However, from a mathematical point of view, almost any non-linearity would suffice to produce cyclopean disparity tuning.

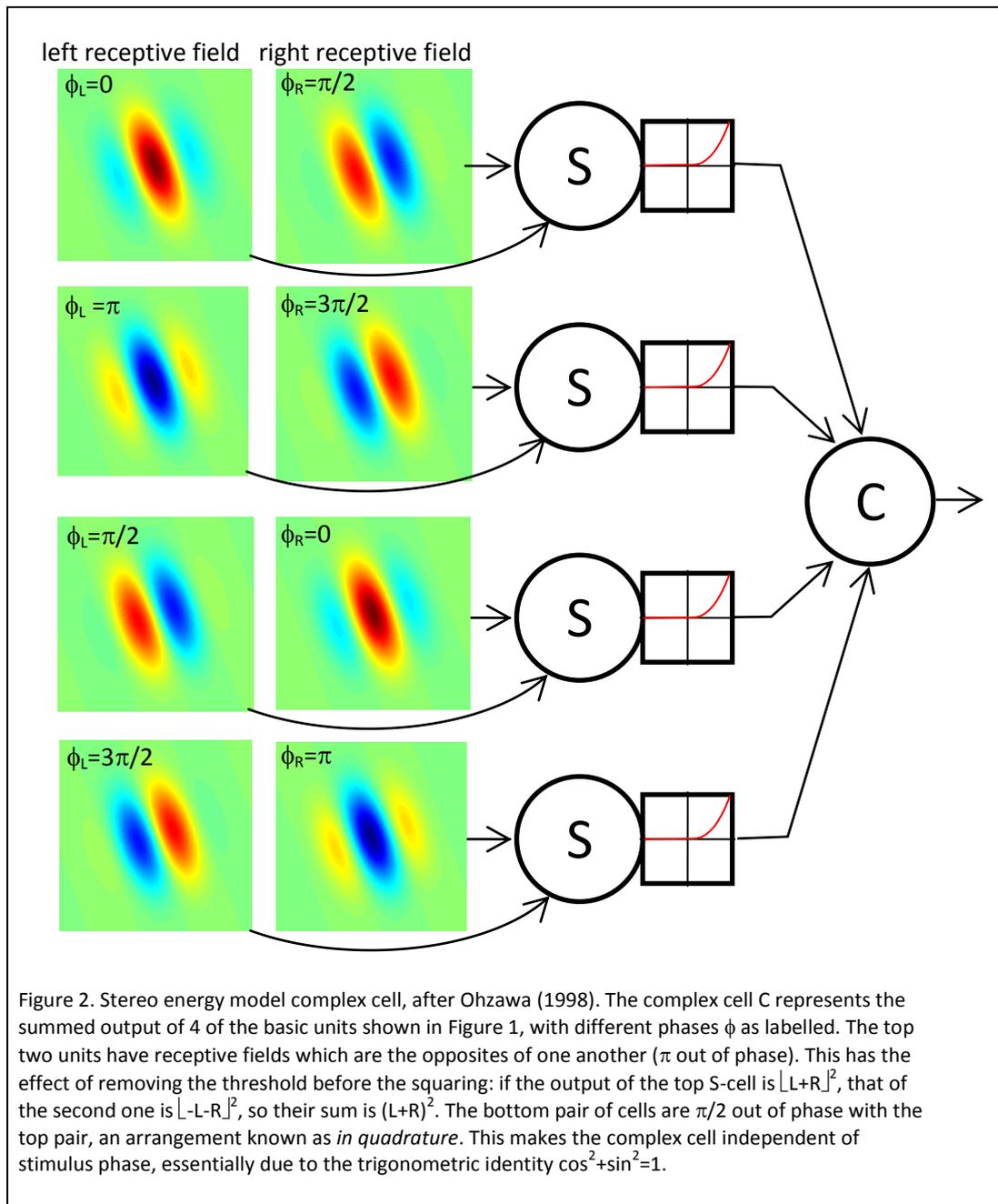
The receptive fields are often represented by Gabor functions. Empirically, the spatial frequency and orientation tuning is similar between eyes (Bridge and Cumming 2001; Read and Cumming 2003), and is usually modeled as being identical.



## Position and phase disparity

The spatial location and phase (i.e. the location of ON and OFF subregions) of the two receptive fields may vary. These determine the preferred disparity of the unit. Position disparity shifts the disparity curve without altering its shape; phase disparity alters the shape. A phase disparity of 0 produces a curve which is symmetrical about a central peak, known as tuned-excitatory. A phase disparity of  $\pi$  inverts this to give a curve with a central trough (tuned-inhibitory). A phase disparity

of  $\pm\pi/2$  produces odd-symmetric curves with equally large peaks and troughs. Several models have been built using pure phase disparity (Sanger 1988; Fleet, Jepson et al. 1991; Qian 1994; Qian and Mikaelian 2000). However, V1 contains cells with both position and phase disparity (Anzai, Ohzawa et al. 1997; Prince, Cumming et al. 2002). Several models therefore incorporate both position and phase (Chen and Qian 2004; Read and Cumming 2007).



## Complex cells

The unit in Equation 1 is a simple cell, in that it is tuned to the phase of a grating stimulus (Movshon, Thompson et al. 1978). To build a complex cell, which is insensitive to grating phase (Hubel and Wiesel 1962; Movshon, Thompson et al. 1978), we can sum many such units with

different receptive-field phases (Qian and Zhu 1997). A common computational short-cut is to sum just four such units, whose receptive fields differ in phase by multiples of  $\pi/2$  (**Figure 2**). This produces a complex cell whose response is perfectly independent of grating phase. Note that this difference in phase must not be confused with the phase disparity discussed above. Phase disparity refers to a difference between the left and right-eye receptive fields of a given energy model unit ( $\phi_R - \phi_L$  for each unit in **Figure 2**). Here, we are talking about a difference in phase between receptive fields of different units, in the same eye (the 4 different  $\phi_L$  in **Figure 2**).

## Modifications of the energy model

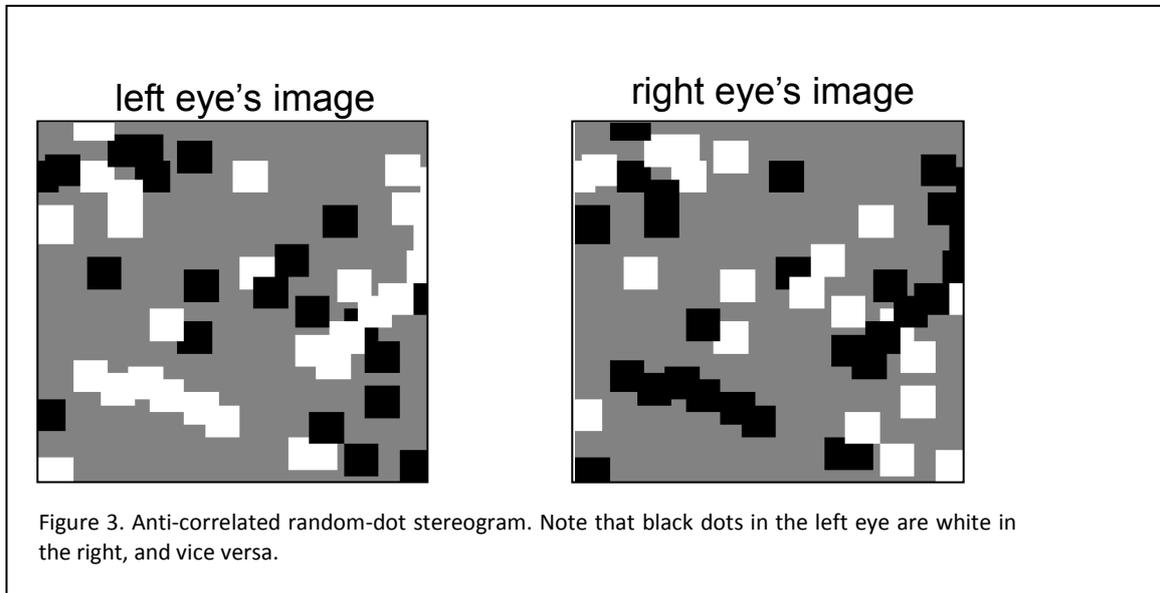
Several modifications have been proposed to the original energy model outlined above. One problem with the original energy model is that it responds equally strongly to disparity in anti-correlated random-dot patterns. These are non-physical stimuli in which one eye's image is the photographic negative of the other (Figure 3). This manipulation all but destroys depth perception (Cogan, Lomakin et al. 1993; Cumming, Shapiro et al. 1998; Read and Eagle 2000). Since this corresponds to changing the sign of  $I_R$  relative to  $I_L$ , it simply inverts the disparity tuning of an energy model unit (Equation 1). However, the response of real V1 neurons is attenuated as well as inverted (Cumming and Parker 1997). To capture the attenuation, Read et al. (2002) proposed thresholding the inner products L and R before summing them in Equation 1. A similar approach was also proposed by Tanaka et al. (2006) to explain second-order stereopsis. Lippert and Wagner (2001) proposed changing the squaring non-linearity. Adding an expansive non-linearity to tuned-excitatory cells, or a compressive non-linearity to tuned-inhibitory cells, can produce the desired effect. Haefner & Cumming (2008) showed that a similar approach captures the response of V1 cells to a wide range of non-physical stimuli. Samonds et al. (2013) incorporated recurrent connections within V1 to account for the temporal dynamics of disparity tuning, as well as the attenuated response to anti-correlated images.

A second problem is that all these versions of the energy model predict a strong relationship between a cell's preferred orientation with grating stimuli and its tuning to 2D disparity. Binocular disparity is in principle 2D since the two eyes' images can be offset both horizontally and vertically, although in normal viewing the horizontal component is much larger (Read, Phillipson et al. 2009). The energy model predicts that cells should be much more sensitive to disparities perpendicular to their preferred orientation than to those parallel to it. This prediction is not borne out: Regardless of preferred orientation, V1 neurons are much more sensitive to vertical than to horizontal disparity (Cumming 2002). One way to modify the energy model to account for this property is to include multiple subunits scattered more widely horizontally than vertically (Read and Cumming 2004).

## The energy model and time

The original energy model included no dependence on time. A more realistic model would include receptive fields whose response depends on time as well as spatial location. The first such model was proposed by Qian & Andersen (1997). Real V1 neurons generally have temporally band-pass tuning, i.e. they respond better to flickering than to constant stimuli. It is difficult to reconcile this with the tuning of real V1 neurons to noise stimuli with interocular delay. In real neurons, interocular delay tends to simply reduce the amplitude of disparity tuning (Read and Cumming

2005). Yet if one simply adds temporally band-pass receptive fields to the energy model, the model predicts that disparity tuning will invert for certain interocular delays.



### Solving the correspondence problem

V1 cells are local, seeing only a small region of each retina. They respond to “false matches”, image features which look similar but were caused by different objects in space, as well as to corresponding views of the same object. To disambiguate these, it is necessary to solve the “correspondence problem”, correctly matching up image features in both eyes’ images (Ohzawa 1998). This process is not completed in V1. Anti-correlated stereograms are one piece of evidence for this. As noted, V1 cells continue to respond to these non-physical stimuli, albeit with attenuation, whereas cortical areas further down the processing pathway, for example IT, do not (Janssen, Vogels et al. 2003).

Models of stereo correspondence predate detailed knowledge of disparity encoding in V1 (Marr and Poggio 1979; Pollard, Mayhew et al. 1985). More recently, modelers have sought to understand how the correspondence problem could be solved by appropriately combining populations of V1-like neurons (Qian 1994; Qian 1997).

One approach is to consider a population of cells which differ only in their tuning to disparity, and take the stimulus disparity to be the preferred disparity of the maximally-responding V1 cell (Qian 1994). However, often this will not be the correct answer, for two main reasons. First, the energy model is not sensitive to the precise arrangement of luminance within the receptive field; this information is lost when the inner product is computed (Equation 1), and its response depends on image contrast. Therefore, an energy-model unit will respond well to a false match which happens to stimulate left and right receptive fields strongly, even if the image within each receptive field is quite different. To overcome these problems, it helps first to normalise out the cell’s response to contrast (Sanger 1988; Read and Cumming 2006; Read 2010), e.g.

$$C = \frac{2LR}{L^2 + R^2}$$

Such divisive normalisation has been found in many brain areas, and is so widespread that it has been proposed as a canonical computation of cortex (Carandini and Heeger 2012). In the present case, the division transforms unnormalised energy, similar to a cross-correlation function (Qian and Zhu 1997), to a normalised correlation, bounded between -1 and +1. Now, the preferred disparity of the maximally-responding C-unit is guaranteed to be the correct stimulus disparity, if the stimulus disparity is exactly constant over the cells' receptive fields (Read and Cumming 2007). However, for more realistic situations, where disparity varies across the image, this approach too can fail.

Greater robustness is obtained by expanding the population under consideration to include cells tuned to a range of orientations and spatial frequencies / spatial scales, as well as to a range of disparities. Fleet and colleagues (Fleet 1994; Fleet, Wagner et al. 1996) proposed a linear pooling of neuronal responses across orientations and scales. Allenmark & Read (2011) showed that by pooling across orientations and spatial frequencies and then normalising as in Equation 2, one effectively computes the cross-correlation of left and right retinal image-patches. This linked the physiologically-based stereo energy model to a class of more abstract models based on windowed cross-correlation of the retinal images, which had successfully captured several aspects of human vision (Tyler 1974; Tyler 1975; Tyler 1978; Banks, Gepshtein et al. 2004; Nienborg, Bridge et al. 2004; Nienborg, Bridge et al. 2005; Filippini and Banks 2009).

Tsai & Victor (2003) used a template-matching approach pioneered in this area by Lehky & Sejnowski (1990). Their population includes neurons tuned to a range of spatial frequency and phase disparities. They compute the mean response of this population to "template" noise stimuli with uniform disparity. The disparity in any given stimulus is taken to be that of the template which best matches the population response. Their model also allows for the perception of *transparency* (multiple planes at different depths) when multiple templates match similarly well, capturing several aspects of human perception (McKee and Verghese 2002). Other aspects of transparency appear to require excitatory and inhibitory interactions between disparity-tuned neurons (Tsirlin, Allison et al. 2012)

## Monocular occlusions

A feature of binocular vision first studied by Leonardo da Vinci is that some image features may be visible to only one eye. These monocular occlusions present a special challenge for establishing correspondence. A few biologically-inspired models suggest how these regions may be handled in the brain (Watanabe and Fukushima 1999; Hayashi, Maeda et al. 2004; Assee and Qian 2007).

## Perceiving depth structure from disparity

Once the stereo correspondence problem has been solved, we have a "disparity map" specifying image disparity at every point in the visual field. This must then be converted into a perception of depth. Relative depth (i.e. "object A is in front of object B") can be deduced immediately from the relative disparity between them. Humans are much more sensitive to the relative disparity between objects or surfaces than to the absolute disparity of an isolated object (Westheimer 1979; Parker 2007).

Further evidence that neurons in primary visual cortex (V1) do not directly support depth perception is provided by the fact that their response is determined by the absolute disparity of the stimulus within their receptive field, not by the disparity relative to other objects in the scene (Cumming and Parker 2000). The extraction of relative disparity appears to begin in cortical area V2, where neurons are found that are sensitive to relative disparity and are specifically tuned for disparity edges (von der Heydt, Zhou et al. 2000; Thomas, Cumming et al. 2002; Bredfeldt and Cumming 2006). These neurons' responses can be modeled by combining the output of different energy model units (Bredfeldt, Read et al. 2009).

For metric depth (i.e. "object A is 10cm in front of object B"), disparity must be calibrated by a knowledge of eye position. In theory, this can be obtained from extra-retinal signals (proprioception, efference copy), or from purely retinal signals if 2D disparity is available. As noted above, disparity in natural viewing is overwhelmingly horizontal, but vertical disparities also occur, in a pattern which is dependent on eye position and largely independent of the scene viewed. In theory, therefore, eye position can be recovered from the 2D disparity map (Longuet-Higgins 1982). There is psychophysical evidence that vertical disparity is indeed detected by the human visual system and used to guide perception (Rogers and Bradshaw 1993; Garding, Porrill et al. 1995; Rogers and Bradshaw 1995; Backus, Banks et al. 1999). However, little is known about the underlying neuronal mechanisms or the cortical areas involved. Detailed neuronally-based models have therefore not been constructed.

- Adelson, E. H. and J. R. Bergen (1985). "Spatiotemporal energy models for the perception of motion." *J Opt Soc Am [A]* **2**(2): 284-299.
- Allenmark, F. and J. C. A. Read (2011). "Spatial stereoresolution for depth corrugations may be set in primary visual cortex" *PLoS Computational Biology* **7**(8): e1002142.
- Anzai, A., I. Ohzawa, et al. (1997). "Neural mechanisms underlying binocular fusion and stereopsis: position vs. phase." *Proc Natl Acad Sci U S A* **94**: 5438-5443.
- Assee, A. and N. Qian (2007). "Solving da Vinci stereopsis with depth-edge-selective V2 cells." *Vision Res* **47**(20): 2585-2602.
- Backus, B. T., M. S. Banks, et al. (1999). "Horizontal and vertical disparity, eye position, and stereoscopic slant perception." *Vision Res* **39**(6): 1143-1170.
- Banks, M. S., S. Gepshtein, et al. (2004). "Why is spatial stereoresolution so low?" *J Neurosci* **24**(9): 2077-2089.
- Bredfeldt, C. E. and B. G. Cumming (2006). "A simple account of cyclopean edge responses in macaque V2." *J Neurosci* **26**(29): 7581-7596.
- Bredfeldt, C. E., J. C. Read, et al. (2009). "A quantitative explanation of responses to disparity defined edges in macaque V2." *J Neurophysiol* **101**(2): 701-713.
- Bridge, H. and B. G. Cumming (2001). "Responses of macaque V1 neurons to binocular orientation differences." *J Neurosci* **21**(18): 7293-7302.
- Carandini, M. and D. J. Heeger (2012). "Normalization as a canonical neural computation." *Nat Rev Neurosci* **13**(1): 51-62.
- Chen, Y. and N. Qian (2004). "A coarse-to-fine disparity energy model with both phase-shift and position-shift receptive field mechanisms." *Neural Comput* **16**(8): 1545-1577.

- Cogan, A. I., A. J. Lomakin, et al. (1993). "Depth in anticorrelated stereograms: effects of spatial density and interocular delay." Vision Res **33**(14): 1959-1975.
- Cumming, B. and A. Parker (2000). "Local disparity not perceived depth is signaled by binocular neurons in cortical area V1 of the Macaque." J Neurosci **20**(12): 4758-4767.
- Cumming, B. G. (2002). "An unexpected specialization for horizontal disparity in primate primary visual cortex." Nature **418**(6898): 633-636.
- Cumming, B. G. and A. J. Parker (1997). "Responses of primary visual cortical neurons to binocular disparity without depth perception." Nature **389**: 280-283.
- Cumming, B. G., S. E. Shapiro, et al. (1998). "Disparity detection in anticorrelated stereograms." Perception **27**(11): 1367-1377.
- Filippini, H. R. and M. S. Banks (2009). "Limits of stereopsis explained by local cross-correlation." J Vis **9**(1): 8 1-18.
- Fleet, D., H. Wagner, et al. (1996). "Neural encoding of binocular disparity: energy models, position shifts and phase shifts." Vision Res **36**: 1839-1857.
- Fleet, D. J. (1994). Disparity from local weighted phase-correlation. Systems, Man, and Cybernetics, 1994. 'Humans, Information and Technology', 1994 IEEE International Conference on, San Antonio, TX, USA.
- Fleet, D. J., A. D. Jepson, et al. (1991). "Phase-based disparity measurement." Computer Vision, Graphics and Image Processing: Image Understanding **53**(2): 198-210.
- Garding, J., J. Porrill, et al. (1995). "Stereopsis, vertical disparity and relief transformations." Vision Res **35**(5): 703-722.
- Haefner, R. M. and B. G. Cumming (2008). "Adaptation to natural binocular disparities in primate V1 explained by a generalized energy model." Neuron **57**(1): 147-158.
- Hayashi, R., T. Maeda, et al. (2004). "An integrative model of binocular vision: a stereo model utilizing interocularly unpaired points produces both depth and binocular rivalry." Vision Res **44**(20): 2367-2380.
- Hubel, D. H. and T. N. Wiesel (1962). "Receptive fields binocular interaction and functional architecture in the cat's visual cortex." J. Physiol. **160**: 106.
- Janssen, P., R. Vogels, et al. (2003). "At least at the level of inferior temporal cortex, the stereo correspondence problem is solved." Neuron **37**(4): 693-701.
- Julesz, B. (1971). Foundations of cyclopean perception. Chicago, University of Chicago Press.
- Lehky, S. R. and T. J. Sejnowski (1990). "Neural model of stereoacuity and depth interpolation based on a distributed representation of stereo disparity [published erratum appears in J Neurosci 1991 Mar;11(3):following Table of Contents]." J Neurosci **10**(7): 2281-2299.
- Lippert, J. and H. Wagner (2001). "A threshold explains modulation of neural responses to opposite-contrast stereograms." Neuroreport **12**(15): 3205-3208.
- Longuet-Higgins, H. C. (1982). "The role of the vertical dimension in stereoscopic vision." Perception **11**(4): 377-386.
- Marr, D. and T. Poggio (1979). "A computational theory of human stereo vision." Proc R Soc Lond B Biol Sci **204**(1156): 301-328.
- McKee, S. P. and P. Verghese (2002). "Stereo transparency and the disparity gradient limit." Vision Res **42**(16): 1963-1977.
- Movshon, J., I. Thompson, et al. (1978). "Spatial summation in the receptive fields of simple cells in the cat's striate cortex." J Physiol **283**: 53-77.
- Movshon, J. A., I. D. Thompson, et al. (1978). "Receptive field organization of complex cells in the cat's striate cortex." J Physiol **283**: 79-99.
- Nienborg, H., H. Bridge, et al. (2004). "Receptive Field Size in V1 Neurons Limits Acuity for Perceiving Disparity Modulation." J Neurosci **24**(9): 2065-2076.
- Nienborg, H., H. Bridge, et al. (2005). "Neuronal computation of disparity in V1 limits temporal resolution for detecting disparity modulation." J Neurosci **25**(44): 10207-10219.

- Ohzawa, I. (1998). "Mechanisms of stereoscopic vision: the disparity energy model." Curr Opin Neurobiol **8**(4): 509-515.
- Ohzawa, I., G. C. DeAngelis, et al. (1990). "Stereoscopic depth discrimination in the visual cortex: neurons ideally suited as disparity detectors." Science **249**: 1037-1041.
- Parker, A. J. (2007). "Binocular depth perception and the cerebral cortex." Nat Rev Neurosci **8**(5): 379-391.
- Poggio, G. F., B. C. Motter, et al. (1985). "Responses of neurons in visual cortex (V1 and V2) of the alert macaque to dynamic random-dot stereograms." Vision Research **25**(3): 397-406.
- Pollard, S. B., J. E. Mayhew, et al. (1985). "PMF: a stereo correspondence algorithm using a disparity gradient limit." Perception **14**(4): 449-470.
- Prince, S. J., B. G. Cumming, et al. (2002). "Range and mechanism of encoding of horizontal disparity in macaque V1." J Neurophysiol **87**(1): 209-221.
- Qian, N. (1994). "Computing stereo disparity and motion with known binocular cell properties." Neural Computation **6**: 390-404.
- Qian, N. (1997). "Binocular disparity and the perception of depth." Neuron **18**(3): 359-368.
- Qian, N. and R. A. Andersen (1997). "A physiological model for motion-stereo integration and a unified explanation of Pulfrich-like phenomena." Vision Res **37**(12): 1683-1698.
- Qian, N. and S. Mikaelian (2000). "Relationship between phase and energy methods for disparity computation." Neural Comput **12**(2): 279-292.
- Qian, N. and Y. Zhu (1997). "Physiological computation of binocular disparity." Vision Res **37**(13): 1811-1827.
- Read, J. C. A. (2005). "Early computational processing in binocular vision and depth perception." Progress in Biophysics and Molecular Biology **87**: 77-108.
- Read, J. C. A. (2010). "Vertical binocular disparity is encoded implicitly within a model neuronal population tuned to horizontal disparity and orientation." PLoS Computational Biology **6**(4): e1000754.
- Read, J. C. A. and B. G. Cumming (2003). "Testing quantitative models of binocular disparity selectivity in primary visual cortex." J Neurophysiol **90**(5): 2795-2817.
- Read, J. C. A. and B. G. Cumming (2004). "Understanding the cortical specialization for horizontal disparity." Neural Comput **16**: 1983-2020.
- Read, J. C. A. and B. G. Cumming (2005). "The effect of interocular delay on disparity selective V1 neurons: relationship to stereoacuity and the Pulfrich effect." Journal of Neurophysiology **94**: 1541-1553.
- Read, J. C. A. and B. G. Cumming (2006). "Does visual perception require vertical disparity detectors?" Journal of Vision **6**(12): 1323-1355.
- Read, J. C. A. and B. G. Cumming (2007). "Sensors for impossible stimuli may solve the stereo correspondence problem." Nat Neurosci **10**(10): 1322-1328.
- Read, J. C. A. and R. A. Eagle (2000). "Reversed stereo depth and motion direction with anti-correlated stimuli." Vision Res **40**(24): 3345-3358.
- Read, J. C. A., A. J. Parker, et al. (2002). "A simple model accounts for the reduced response of disparity-tuned V1 neurons to anti-correlated images." Vis Neurosci **19**: 735-753.
- Read, J. C. A., G. P. Phillipson, et al. (2009). "Latitude and longitude vertical disparities." J Vis **9**(13): 1111-1137.
- Rogers, B. J. and M. F. Bradshaw (1993). "Vertical disparities, differential perspective and binocular stereopsis." Nature **361**(6409): 253-255.
- Rogers, B. J. and M. F. Bradshaw (1995). "Disparity scaling and the perception of frontoparallel surfaces." Perception **24**(2): 155-179.
- Samonds, J. M., B. R. Potetz, et al. (2013). "Recurrent connectivity can account for the dynamics of disparity processing in v1." J Neurosci **33**(7): 2934-2946.
- Sanger, T. (1988). "Stereo disparity computation using Gabor filters." Biological Cybernetics **59**(405-418).

- Tanaka, H. and I. Ohzawa (2006). "Neural basis for stereopsis from second-order contrast cues." J Neurosci **26**(16): 4370-4382.
- Thomas, O. M., B. G. Cumming, et al. (2002). "A specialization for relative disparity in V2." Nat Neurosci **5**(5): 472-478.
- Tsai, J. J. and J. D. Victor (2003). "Reading a population code: a multi-scale neural model for representing binocular disparity." Vision Res **43**(4): 445-466.
- Tsirlin, I., R. S. Allison, et al. (2012). "Perceptual asymmetry reveals neural substrates underlying stereoscopic transparency." Vision Res **54**: 1-11.
- Tyler, C. W. (1974). "Depth perception in disparity gratings." Nature **251**(5471): 140-142.
- Tyler, C. W. (1975). "Spatial organization of binocular disparity sensitivity." Vision Res **15**(5): 583-590.
- Tyler, C. W. (1978). "Binocular cross-correlation in time and space." Vision Res **18**(1): 101-105.
- von der Heydt, R., H. Zhou, et al. (2000). "Representation of stereoscopic edges in monkey visual cortex." Vision Res **40**(15): 1955-1967.
- Watanabe, O. and K. Fukushima (1999). "Stereo algorithm that extracts a depth cue from interocularly unpaired points." Neural Netw **12**(4-5): 569-578.
- Westheimer, G. (1979). "Cooperative neural processes involved in stereoscopic acuity." Exp Brain Res **36**(3): 585-597.