

Image Segmentation based on Motion/Luminance Integration and Oscillatory Correlation

Erdogan Çesmeli¹ and Deliang L. Wang²

¹Biomedical Engineering Center

²Department of Computer and Information Science and Center for Cognitive Science
The Ohio State University, Columbus, OH 43210, USA
{cesmeli, dwang}@cis.ohio-state.edu

Abstract

An image segmentation method is proposed based on the integration of motion and luminance information. The method is composed of two parallel pathways that process motion and luminance, respectively. Inspired by the visual system, the motion pathway has two stages. The first stage estimates local motion at locations with reliable information. The second stage groups locations based on their motion estimates. In the parallel pathway, the input scene is segmented based on luminance. In the subsequent integration stage, motion estimates are refined to obtain the final segmentation result in the motion pathway. For segmentation, LEGION (Locally Excitatory Globally Inhibitory Oscillator Networks) is employed whereby the phases of oscillators are used for region labeling. Results on synthetic and real image sequences are provided.

1. Introduction

A central problem to computational investigation of motion perception is the selective integration of local motion estimates. Most approaches assume that motion integration can be addressed using only motion information. However, recent studies, e.g. [10, 4], show that the assumption may not be valid. Their investigation is based on a well known visual stimulus, namely plaids, which are constructed by the superimposition of two moving gratings at different orientations. When presented, observers report either a *coherent* motion corresponding to that of the plaid or a pair of *component* motions belonging to the gratings. Without changing the underlying motion, when the luminance at the intersection of the gratings is varied to support that gratings belong to two different surfaces, component motions are more frequently perceived. Consistent with their observation, other studies [9,12] also indicate that non-motion cues, e.g. stereo, might play a role in motion integration. Similarly, engineering studies, e.g. [2], demonstrate that the inclusion of luminance information improves motion-based segmentation.

Motivated by these studies, we propose a segmentation method based on the integration of motion and luminance

information using neural oscillator networks. Our method consists of two parallel pathways for motion and luminance, respectively. The motion pathway has two stages. First, local motions are estimated by an adaptive temporal block matcher, a variation of the Reichardt detector [7]. In the second stage, locations are grouped based on motion similarity in a multilayered LEGION network. LEGION is based on the idea of oscillatory correlation [5, 14, 11], whereby phases of neural oscillators encode region labeling. In order to complement the initial motion segmentation, a stationary scene analysis is performed in parallel in the luminance pathway also using a LEGION network [14]. Next, the integration stage refines motion estimates. The final segmentation is obtained in the motion network based on the refined estimates.

The following section describes the two building blocks of our method. Next is the detailed presentation of our method. Finally, its performance is demonstrated on both synthetic and real scenes, and conclusions are drawn.

2. Background

2.1 Temporal block matcher (TBM)

TBM detects different velocities by comparing luminance within a block, N_B , with those located at different distances and directions on the previous snapshot (frame) of the scene as shown in Figure 1. At location (x, y) and time t , the correlation corresponding to displacement, $\mathbf{r} = (r_x, r_y)$, is:

$$\tilde{v}_{\mathbf{r}}(x, y, t) = \sum_{(j, k) \in N_B(x, y)} \frac{I(j, k, t)}{|I(j, k, t) - I(j - r_x, k - r_y, t - 1)|} \quad (1)$$

where $I(x, y, t)$ is the luminance at location (x, y) in frame t and N_B is centered at location (x, y) . As in that of (1), denominators in all expressions include a small quantity to avoid division by zero. Note that a large correlation implies a high probability for the displacement, $\mathbf{r} = (r_x, r_y)$ at (x, y) .

2.2 LEGION

LEGION is based on the idea of oscillatory correlation [5, 14, 11], where the phases of the neural oscillators encode the grouping of locations with similar stimuli. The building

block of a LEGION network is a single relaxation oscillator, i , and is defined as a feedback loop between an excitatory unit x_i and an inhibitory unit y_i [14,11]:

$$dx_i/dt = 3x_i - x_i^3 + 2 - y_i + S_i + \rho \quad (2a)$$

$$dy_i/dt = \varepsilon \left(\alpha \left(1 + \tanh\left(\frac{x_i}{\beta}\right) \right) - y_i \right) \quad (2b)$$

Here S_i denotes coupling, ρ is the variance of Gaussian noise, and α and β are system parameters. The parameter ε is chosen to be a small positive number so that (2) defines a relaxation oscillator where the x-nullcline ($dx/dt = 0$) is a cubic curve and the y-nullcline ($dy/dt = 0$) is a sigmoid, as shown in Figure 2A. Only when these curves intersect along the middle branch of the x-nullcline, the system is oscillatory. Oscillator i travels along the left branch (LB) reaching the left knee (LK) and jumps to the right branch (RB), where it becomes *active*. After traveling along RB, it reaches the right knee (RK), where it jumps back to LB completing a limit cycle. A temporal trace of several limit cycles of an oscillator is depicted in Figure 2B.

S_i includes a variable called lateral potential and coupling from neighboring oscillators and a global inhibitor:

$$S_i = \sum_{k \in N(i)} W_{ik} H(x_k) + W_p H(p_i - 0.5) - W_z H(z - 0.5) \quad (3)$$

Here, W_{ik} is the connection weight from oscillator k to oscillator i , $H(\cdot)$ is the Heaviside step function, W_p and W_z are the weights for the potential, p_i , and the global inhibition, z , respectively. N represents a local coupling neighborhood, e.g. four nearest neighbors. Each oscillator is assigned the potential variable, p_i , to distinguish homogeneous regions from noisy ones:

$$dp_i/dt = (1 - p_i) H \left[\sum_{k \in N_p(i)} H(x_k) - \theta_p \right] - \varepsilon p_i$$

where $N_p > N$ is the potential neighborhood. The potential of an oscillator is initially high but continuously decays. When an oscillator is active and has a number of active neighbors more than θ_p in its N_p , its potential rises to 1. Oscillators that maintain high potential are referred to as leaders while others are called followers. Groups of strongly coupled oscillators are candidates to form segments. However, due to the potential term, only when a group has a leader can it form a segment. Oscillators that cannot become active form background.

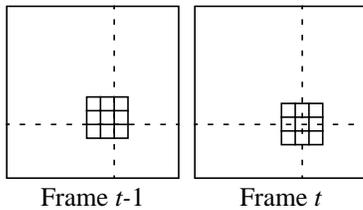


Figure 1. TBM cross-correlates luminance within $N_B = 3 \times 3$ in two subsequent frames to detect the motion corresponding to $r = (1, 1)$ pixel/frame.

Whether an oscillator can become active also depends on the global inhibition, z :

$$dz/dt = H \left[\sum_k H(x_k) - 1 \right] - z$$

When no oscillator is active, z decays to 0, otherwise it rises to 1. Sole leadership allows an oscillator to jump only when $z < 0.5$, as indicated by the last term in (3). Otherwise, an oscillator can become active only when it has strong couplings with currently active oscillators. Parameters including ε , ρ , α , β , θ_p , and the weights W_p and W_z are constant and W_{ik} is defined for a particular feature space, e.g. luminance, motion. A common form of a LEGION network is a two dimensional (2D) array of oscillators and a global inhibitor (see Fig. 4B below).

3. The Method: Adaptive TBM and LEGION

Our method is shown in Figure 3 where the motion pathway, outlined in a box, performs motion analysis. The parallel pathway segments the scene based on luminance. Performing occlusion analysis, estimates are refined to obtain the final segmentation in the motion network.

3.1 Motion Pathway

Consistent with psychology and neurophysiology (see [8] for a review), our method has two stages for motion estimation and grouping, respectively.

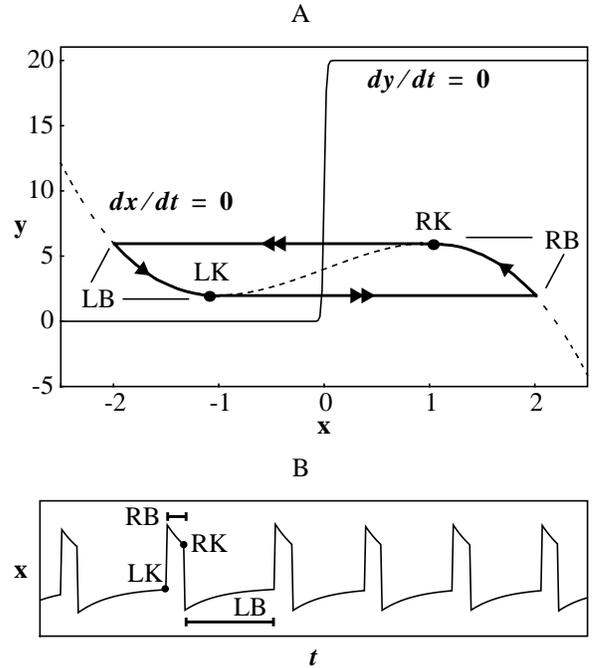


Figure 2. A) Phase plane diagram of a single oscillator. The x-nullcline is the dotted line and the y-nullcline is the solid line. The oscillator produces a limit cycle, drawn as thick solid line. B) The temporal activity of the oscillator in A, where $\varepsilon = 0.005$, $\alpha = 10.0$, $\beta = 0.02$, $I_i = 2.0$, and $\rho = 0.02$.

We apply an adaptive TBM to calculate the correlation corresponding to displacement $\mathbf{r} = (r_x, r_y)$:

$$v_r(x, y, t) = \tilde{v}_r(x, y, t) + \tilde{v}_r(x, y, t+1) \quad (4)$$

where $\tilde{v}_r(x, y, f)$ is the correlation at (x, y) between the frames f and $f-1$ as given in (1). We have a total number of $L = (2R+1)^2$ different velocities corresponding to a set of displacements varying from $-R$ to R in the x - and y -directions. In estimation, we also employ local spatial correlation surfaces (SCS). An SCS is obtained by applying (1) within the same frame. Replacing luminance blocks with SCSs, a cross-correlation, $c_r(x, y, t)$, is obtained using (4). Consequently, the temporal correlation at (x, y, t) for displacement $\mathbf{r} = (r_x, r_y)$ is defined as:

$$\hat{V}_r(x, y, t) = v_r(x, y, t) + c_r(x, y, t) \quad (5)$$

The shape of a cross-correlation surface depends on the underlying luminance structure at a location. In order to obtain a unique solution in the presence of aperture problem, as in the case of a straight border, we multiply the cross-correlation surface, \hat{V} , obtained using (5) with a 2D Gaussian centered at zero velocity:

$$V_r(x, y, t) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{|\mathbf{r}|^2}{2\sigma^2}\right) \hat{V}_r(x, y, t) \quad (6)$$

where σ is large enough to allow the maximum velocity considered. We assume the displacement, \mathbf{e} , yielding the peak correlation, V_e , to be the local motion estimate at (x, y, t) . Note that a peak obtained for a location along a straight border may not be as well-localized as that of a corner. In order to quantify these differences, we define an estimate certainty along a direction axis which is perpendicular to the speed axis passing through origin and \mathbf{e} in the cross-correlation space. Dropping (x, y, t) from the expressions, the certainty, ω_e , of \mathbf{e} is given by:

$$\omega_e = \frac{(V_e - V_{e-k})(V_e - V_{e+k})}{2V_e - V_{e-k} - V_{e+k}} \quad (7)$$

Here, $\mathbf{e}-\mathbf{k}$ and $\mathbf{e}+\mathbf{k}$ are the nearest neighboring displacements to \mathbf{e} along the direction axis and correspond to correlations V_{e-k} and V_{e+k} , respectively. According to (7), the larger ω_e is, the sharper V_e is and the more certain the estimate, \mathbf{e} , is.

An estimate and its certainty are obtained at locations that satisfy a reliability criterion quantified by a *mobility value*:

$$M(x, y, t) = \frac{\sum_{(j,k) \in N_M(x,y)} |I(j, k, t) - I(j, k, t-1)|}{\sum_{(j,k) \in N_M(x,y)} I(j, k, t)} \quad (8)$$

where N_M is the mobility neighborhood. A large value at a location indicates a high probability of motion at that location. Based on mobility values, N_B is selected at each location. Initially small N_B is expanded until it includes sufficiently large total and average mobility or its size reaches an upper limit. Unless the limit is reached, (6) and (7) are employed to obtain an estimate and its certainty.

In order to allow for multiple overlapping motions, i.e. motion transparency [3], our method represents each velocity with a LEGION network as depicted in Figure 4A. The coupling weight between two oscillators, i and k , on a velocity layer \mathbf{r} is given by:

$$W_{r,ik} = \frac{V_r(i) + V_r(k)}{|V_r(i) - V_r(k)|}$$

When oscillators have similar correlations for a particular displacement, they are strongly coupled in the corresponding velocity layer. At locations without estimates, couplings are set to zero. We replace S_i in (3) by $S_{ri} = H(\tilde{S}_{ri} - \theta_M)$ where θ_M is a threshold and

$$\begin{aligned} \tilde{S}_{ri} &= \sum_{k \in N(i)} W_{r,ik} H(x_k) - W_z H(z - 0.5) \\ &+ W_p H(p_{ri} - 0.5) H(M_i - 0.25) H\left[\sum_{q=1}^L H(V_{ri} - V_{r_q i}) - L\right] \end{aligned}$$

The second $H(\cdot)$ multiplying W_p ensures that only locations with large mobility values can become leaders. The last $H(\cdot)$ only allows for a single leader within each velocity column. Due to this competition within each column, the oscillator with the largest correlation becomes a winner. Provided that a winner has sufficiently high mobility and potential, it becomes a leader and starts

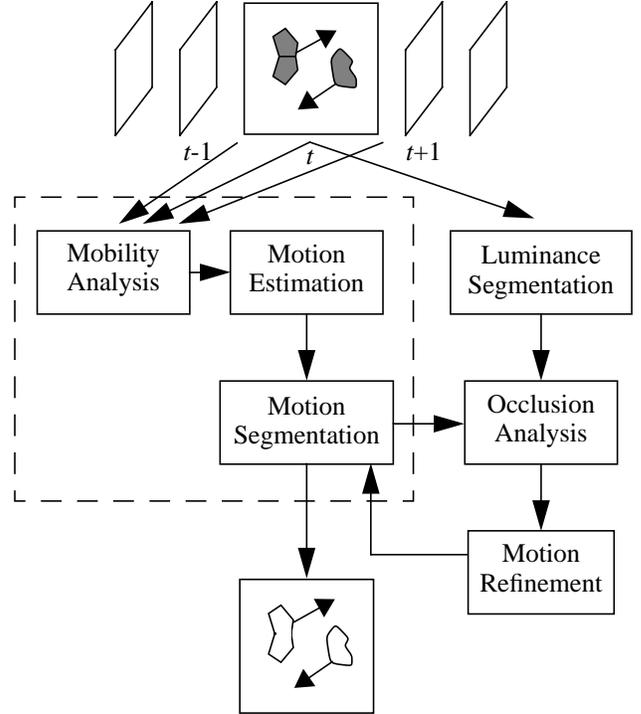


Figure 3. The flow diagram of our method. Processing progresses downward. Occlusion analysis facilitates the refinement of motion estimates based on which the final segmentation is obtained in the motion network.

forming its segment. A leader recruits oscillators on its layer through local couplings when they have similar correlations. As in the single layer LEGION network, recruited oscillators become active simultaneously (synchronization). Because of the global inhibition across all layers, leaders form their segments at different times (desynchronization). Also, note that an oscillator can become active when it is a leader or a follower recruited by a leader on its layer. Since there could be several followers in the velocity column of a leader, more than one oscillator at a single location can become active representing different motions. This ability has a key role in the representation of motion transparency.

3.2 Luminance Pathway

The parallel luminance pathway processes the middle frame of the sequence analyzed in the motion pathway. In this pathway, where also a LEGION network is employed as depicted in Figure 4B, we assume that each region is depicted approximately homogeneous. The coupling weight between two oscillators, i and k , is defined as:

$$W_{ik} = \frac{I(i) + I(k)}{|I(i) - I(k)|}$$

When locations have relatively similar luminance, a strong coupling weight results. We employ the network model given in (3) by replacing S_i with $H(S_i - \theta_B)$ where θ_B is a threshold. Strongly coupled oscillators with at least one leader become synchronized. Oscillators corresponding to regions with different luminances become desynchronized. When oscillators correspond to textured regions, they tend not to have strong couplings and thus, leaders. Lacking leaders, they do not form segments and are distinguished from homogeneous regions, a distinction that has a significant role in the integration stage.

3.3 Integration Stage

The first step of the integration stage is an occlusion analysis. It first considers motion segments. When all locations in a motion segment belong to a textured region, estimates in this segment are not changed. However, when the majority of locations belong to multiple homogeneous regions, an occlusion relationship among these regions is obtained by detecting T- and X-junctions. A T-junction detected among three homogeneous regions indicates the occluding opaque region and the occluded ones. An X-junction shows which two of the four homogeneous regions form a transparent occluding surface [6, 1]. Our method detects T- and X-junctions by applying a set of templates to the luminance segments.

When a segment includes both textured and homogeneous regions, the occlusion relationship among them is resolved by determining two types of motion distributions in the homogeneous regions. The first one includes estimates at all locations. The second one considers only those along the boundary of a textured region. When the peak of the first distribution is the same as that of the second in a homogeneous region, the textured region is assumed to

occlude the homogeneous one. Otherwise, the two regions are assumed to move together.

In the second step of the integration stage, estimates along an occluding boundary are eliminated in occluded regions. The remaining estimates in each luminance segment, B , interact iteratively and result in a segment velocity, r_B , for that segment:

$$r_B^\tau = \left(\sum_{(x,y) \in B} \omega^\tau(x,y) r(x,y) \right) / \Omega_B^\tau \quad (9a)$$

$$\omega^{\tau+1}(x,y) = \frac{\omega^\tau(x,y)}{\Omega_B^\tau} \left(1 + \frac{r(x,y) \cdot r_B^\tau}{\sqrt{\|r(x,y)\| \|r_B^\tau\|}} \right) \quad (9b)$$

Here, r_B^τ and Ω_B^τ are the segment velocity and the sum of certainties, $\omega^\tau(x,y)$, in B at iteration step τ , respectively. $a \cdot b$ is the dot product of vectors a and b , and $\|a\|$ is the magnitude of a . In (9a), r_B^τ is determined by weighing the estimates in B by their certainties. In (9b), $\omega^{\tau+1}(x,y)$

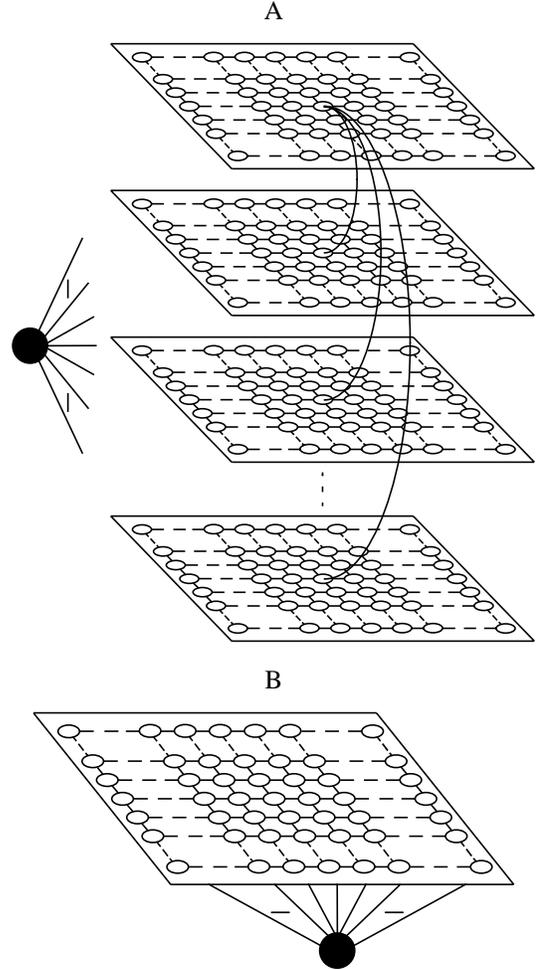


Figure 4. Neural network architectures. Small ellipses represent oscillators. The black circle is the global inhibitor. **A)** Multilayered motion network where each velocity layer is a LEGION network. **B)** 2D luminance network.

increases when $r(x, y)$ and r_B^τ have a similar direction. Finally, when $\omega^{\tau+1}(x, y)$ s in B do not change, r_B^τ is assumed to have converged to its final value, r_B . As a result, r_B is filled-in at all locations in B .

The motion interaction in (9) takes place in all luminance segments. Since textured regions, which already have reliable estimates, do not form luminance segments, the motion interaction does not take place in these regions.

Also, noting that large homogeneous regions touching image borders tend to be perceived stationary, prior to (9), estimates of zero velocity are assigned to locations without estimates in luminance segments along the image borders.

Following the integration stage, couplings in the motion network are updated based on the refined estimates and the final segmentation result is obtained.

4. Results

We demonstrate the performance of our method using both synthetic and real scenes. Figure 5A shows an input scene where two vertical rectangles are moving toward each other while a horizontal one is moving downward. In the motion pathway, first, the mobility image is obtained using (8) as shown in Figure 5B. Next, local motions and their certainties are estimated at locations with large mobility values, as depicted in Figures 5C-D. As shown in Figure 5E, the initial motion segmentation is based on these estimates. Note that erroneous estimates in occluded regions along occluding boundaries cause inaccurate segmentation. However, the segmentation result in the luminance pathway, as shown in Figure 5F, matches well with the input scene. In the subsequent integration stage, the occlusion analysis is performed using these segmentation results. Having detected T-junctions as depicted in Figure 5G, it is found that the large vertical rectangle occludes the background and the horizontal rectangle, which, in turn, occludes the background and the small vertical rectangle. Figures 5H-I show the remaining estimates and their certainties after removing the ones in the occluded regions along the occluding boundaries. Finally, the motion interaction takes place in luminance segments, filling-in their locations with their resulting segment velocities as illustrated in Figure 5J. The final segmentation result based on the refined estimates is shown in Figure 5K and compares well with the regions and their motions in the input scene. Note that the homogenous background is assigned zero velocity due to the introduction of estimates along the image borders.

The input scene in Figure 6A is composed of two square regions. The upper region moves in the right and downward direction while the lower one is having a left- and downward motion. In the center of the scene, the regions overlap transparently. Similar to the example in Figure 5, the scene is processed and initial segmentation results are obtained in the two pathways. By detecting X-junctions, transparency is inferred and locations in each region are filled-in with their

segment velocities as shown in Figure 6B. Note that locations in the overlapping area are assigned both velocities. Thus, the final result includes two overlapping square segments as depicted in Figure 6C.

An intriguing visual illusion occurs when gratings move behind an aperture. Gratings in Figures 7A and C move vertically upward. When the aperture is circular, the motion is perceived to be in the perpendicular direction to the grating orientation, consistent with our result shown in Figure 7B. When the aperture is rectangular, our method results in a distribution of velocities which are parallel to the longer axis of the aperture mimicking the barber pole

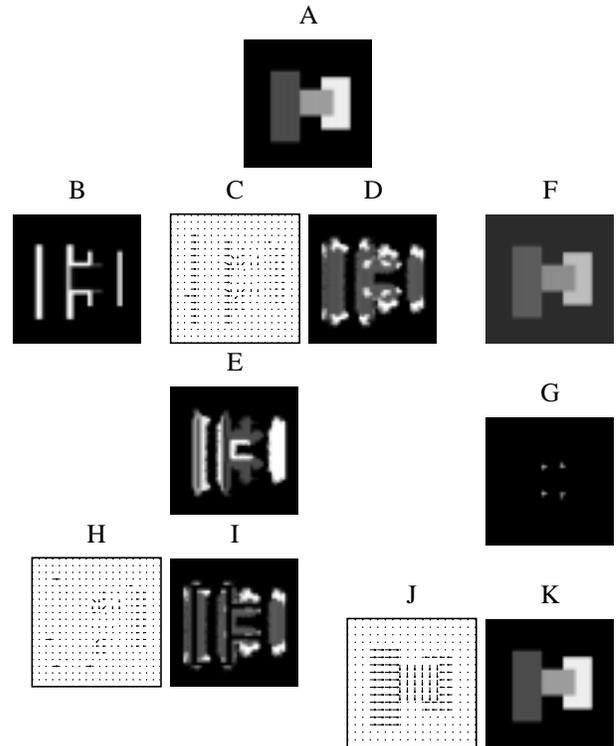


Figure 5. A segmentation example. **A)** The middle frame of the input sequence. **B)** Mobility analysis. **C-D)** Motion estimates and their certainties. **E)** Initial motion segmentation. **F)** Luminance Segmentation. **G)** Occlusion analysis where T-junctions are detected. **H-I)** Unreliable estimates and their certainties are removed. **J)** The result of motion interaction where estimates are refined and filled-in. **K)** Final motion segmentation.

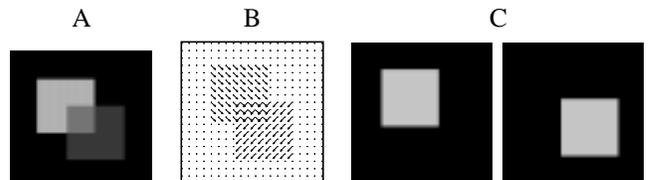


Figure 6. Motion transparency. **A)** The middle frame of the input sequence. **B)** Refined motion estimates. **C)** Overlapping segments represent motion transparency.

illusion [13] as shown in Figure 7D.

Our method also performs well with real scenes. In Figure 8A, a rider and his motorcycle move diagonally left and downward while the background appears to move in the opposite direction due to the camera motion. Our method eliminates erroneous estimates in the occluded regions, e.g. homogenous background occluded by the motorcycle. The final result segregates the rider and his motorcycle from the background as shown in Figures 8B. Similarly, in Figure 8C, the camera pans to the right. Due to their different distances to the camera, a woman and two dish antennas appear to have different motions. Our method is able to deal with neighboring homogenous regions in the scene, e.g the woman's hat and the antenna, as depicted in Figure 8D.

In addition to the network parameters given in Figure 2, we used $N = 3 \times 3$, $N_z = 7 \times 7$, $W = 0.74$, $\theta_z = 36.75$, and $W_z = 1.00$. In $p(8)$, $N_M = 5^p \times 5$, which is also the initial size of N_B . For all results, the threshold set, $(\theta_{M,1}, \theta_{M,2}, \theta_B)$, is $(20, 10, 10)$, except for Figure 8B where $\theta_B = 50$. Numerical subscripts, 1 and 2, in θ_M correspond to initial and final segmentations, respectively. Motion estimates in Figures 5-7 are spatially subsampled

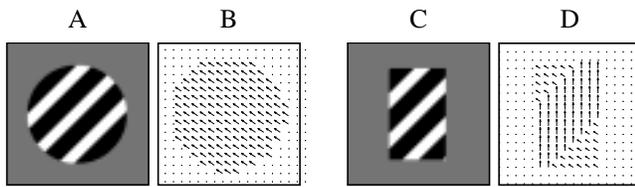


Figure 7. Barber-pole illusion. The input scenes and the resulting motion distributions when the aperture is **A-B)** circular or **C-D)** rectangular. Even though gratings have the same motion in A and C, the results in B and D depend on the aperture shape.

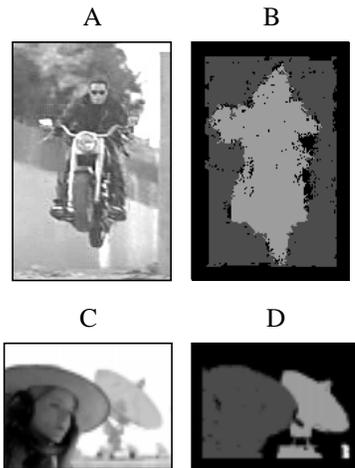


Figure 8. Real scenes. **A-B)** The motorcycle and the rider have a different motion than the background. **C-D)** Due to their different distances to the camera, the woman and the antennas appear to move differently.

for visual clarity.

5. Conclusion

We proposed a segmentation method based on the integration of motion and luminance information. Our method is able to represent motion transparency and deal with occlusion in both synthetic and real scenes. It can also mimic visual illusions. Since our neural network architecture is biologically plausible, our method has implications for new psychophysical experiments to study human visual system.

Acknowledgments

We thank J. T. Todd and D. T. Lindsey for many helpful discussions. This work was supported in part by an NSF grant (IRI-9423312) and an ONR Young Investigator Award (N00014-96-1-0676) to DLW.

References

- [1] J. Beck, K. Pradny, and R. Ivry, "The perception of transparency with achromatic colors," *Percept. Psychophys.*, 35:407-422, 1984.
- [2] M. J. Black and A. D. Jepson, "Estimating optical flow in segmented images using variable-order parametric models with local deformations," *IEEE Trans. Pattern. Anal. Machine Intell.*, 18(10):972-986, 1996.
- [3] E. J. Gibson, J. J. Gibson, O. W. Smith, and H. Flock, "Motion parallax as a determinant of perceived depth," *J. Exp. Psychol.*, 58(1):40-51, 1959.
- [4] D. T. Lindsey and J. T. Todd, "On the relative contributions of motion energy and transparency to the perception of moving plaids," *Vision Res.*, 36(2):207-222, 1996.
- [5] C. v. d. Malsburg, "The correlation theory of brain functions," *Internal Report #81-2*, Max-Planck-Institut for Biophysical Chemistry, Göttingen, FRG, 1981.
- [6] F. Metelli, "The perception of transparency," *Sci. Amer.*, 230(4):90-98, 1974.
- [7] W. Reichardt, "Autokorrelationsauswertung als Funktionsprinzip des Zentralnervensystems" *Z. Naturforsch.*, 12b:447-457, 1957.
- [8] M. E. Sereno, *Neural Computation of pattern motion: modeling stages of motion analysis in the primate visual cortex*. Cambridge MA: MIT press, 1993.
- [9] S. Shimojo, G. H. Silverman, and K. Nakayama, "Occlusion and the solution to the aperture problem for motion," *Vision Res.*, 29(5):619-626, 1989.
- [10] G. R. Stoner, T. D. Albright, and V. S. Ramachandran, "Transparency and coherence in human motion perception," *Nature*, 344:153-155, 1990.
- [11] D. Terman and D. L. Wang, "Global competition and local cooperation in a network of neural oscillators," *Physica D*, 81:148-176, 1995.
- [12] J. C. Trueswell and M. M. Hayhoe, "Surface segmentation mechanism and motion perception," *Vision Res.*, 33(3):313-328, 1993.
- [13] H. Wallach, "Über visuell wahrgenommene Bewegungsrichtung," *Psychologische Forschung*, 20:325-380, 1935.
- [14] D. L. Wang and D. Terman, "Locally excitatory globally inhibitory oscillator networks," *IEEE Trans. Neural Networks*, 6:283-286, 1995.