

- [9] E. Rivlin, S.J. Dickinson, and A. Rosenfeld. Recognition by functional parts. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 267–274, Seattle, WA, June 1994.
- [10] E. Rivlin, A. Rosenfeld, and D. Perlis. Recognition of object functionality in goal-directed robotics. In *Proceedings of the AAAI Workshop on Reasoning about Function*, 1993.
- [11] F. Solina and R. Bajcsy. Shape and function. In *Proceedings of the SPIE Conference on Intelligent Robots and Computer Vision*, Volume 726, pages 284–291, 1983.
- [12] L. Stark and K. Bowyer. Achieving generalized object recognition through reasoning about association of function to structure. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1097–1104, Maui, HI, June 1991.
- [13] L. Stark and K. Bowyer. Generic recognition through qualitative reasoning about 3-d shape and object function. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 251–256, Maui, HI, 1991.
- [14] L. Stark and K. Bowyer. Indexing function-based categories for generic recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 795–797, Champaign, IL, June 1992.
- [15] L. Stark, A. Hoover, D. Goldgof, and K. Bowyer. Function based recognition from incomplete knowledge of shape. In *Proceedings of the IEEE Workshop on Qualitative Vision*, pages 11–22, New York, NY, 1993.
- [16] S. Ullman and R. Basri. Recognition by linear combinations of models. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13:992–1006, 1991.
- [17] A. Verri and T. Poggio. Against quantitative optical flow. In *Proceedings of the International Conference on Computer Vision*, pages 171–180, London, England, June 1987.
- [18] P.H. Winston, T.O. Binford, B. Katz, and M. Lowry. Learning physical descriptions from functional descriptions, examples, and precedents. In *Proceedings of the National Conference on Artificial Intelligence*, pages 433–439, 1983.

Seven sequences of images were used to demonstrate the approach. Function understanding from motion was established in all seven cases. In the first three sequences, motion was used to discriminate between three cutting actions: stabbing, chopping and jabbing. In the last two pairs of sequences we used motion information to differentiate between two different functionalities of the same object: scooping and hitting with a shovel, and hammering and tightening with a wrench. These examples of double usage are typical instances of improvisation; motion provides clear information for a correct interpretation of the action that is taking place.

Natural extensions of this work include the analysis of more complex objects. Complexity can be expressed in terms of either the shapes of the parts or the way in which the parts are connected. An interesting area is the analysis of articulated objects. The different types of connections between the parts constrain the possible relative motions of the parts. A pair of pliers or a pair of scissors is a simple case, with only a single articulated connection (one degree of freedom in relative motion of the parts). We see our work as a step toward action perception of moving objects, which could lead to a better understanding of perceiving the actions of moving agents.

References

- [1] L. Bogoni and R. Bajcsy. Active investigation of functionality. In *Proceedings of the IAPR Workshop on Visual Behaviors*, Seattle, WA, June 1984.
- [2] M. Brady, P.E. Agre, D.J. Braunegg, and J. Connell, II. The mechanic's mate. In *Proceedings of the Sixth European Conference on Artificial Intelligence*, pages 79–94, 1984.
- [3] P. Freeman and A. Newell. A model for functional reasoning in design. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 621–640, August 1971.
- [4] K. Gould and M. Shah. The trajectory primal sketch: A multi-scale scheme for representing motion characteristics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 79–85, June 1989.
- [5] B.K.P. Horn and B.G. Schunck. Determining optical flow. *Artificial Intelligence*, 17:189–203, 1981.
- [6] J.R. Kender and D.G. Freudenstein. What is a degenerate view? In *Proceedings of the DARPA Image Understanding Workshop*, pages 589–598, 1987.
- [7] H. Murase and S.K. Nayar. Learning object models from appearance. In *Proceedings of the National Conference on Artificial Intelligence*, pages 836–843, Washington, DC, July 1993.
- [8] R. Polana and R. Nelson. Detecting activities. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2–7, New York, NY, June 1993.

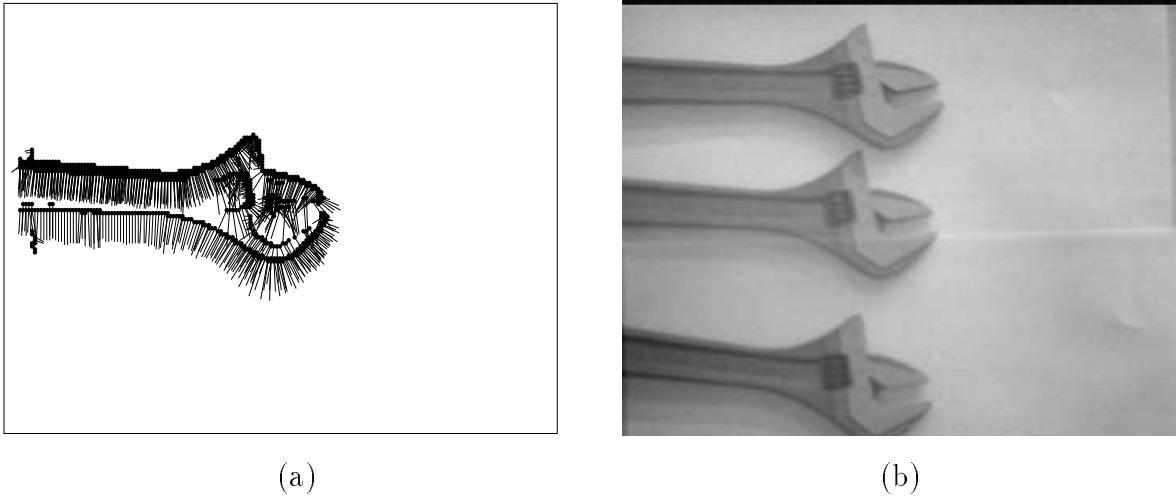


Figure 14: (a) Flow vectors for hammering with a wrench. (b) Hammering motion.

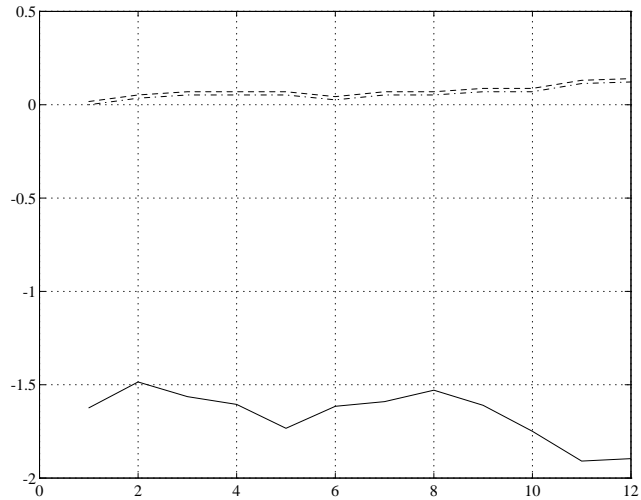


Figure 15: Angles α , β , and θ for hammering with a wrench. α is given by a dashed line, β is given by a solid line, and θ is given by a dash-dot line.

7 Conclusions

Perceiving function from motion provides an understanding of the way an object is being used by an agent. To accomplish this we combined information on the shape of the object, its motion, and its relation to the actee (the object it is acting on). Assuming a decomposition of the object into primitive parts, we analyzed a part's motion relative to its principal axes. Primitive motions (translation and rotation relative to the principal axes of the object) were dominating factors in the analysis. We used a frame of reference relative to the actee. Once such a frame is established, it can have major implications for the functionality of an action.

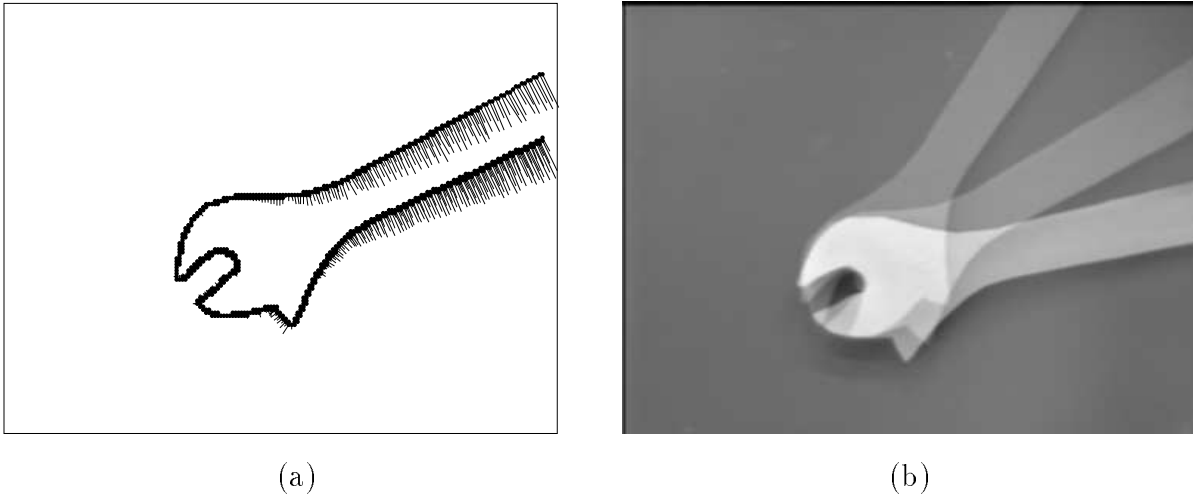


Figure 12: (a) Flow Vectors for tightening with a wrench. (b) Tightening motion.

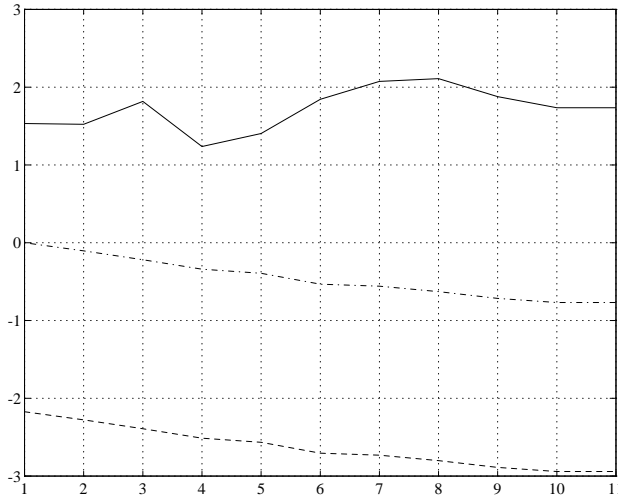


Figure 13: Angles α , β , and θ for tightening with a wrench. α is given by a dashed line, β is given by a solid line, and θ is given by a dash-dot line.

motion dominates over the rotational part of the motion and the direction of translation is approximately orthogonal to the direction of the medial axis of the tool. Figure 14 shows the flow vectors taken from the 6th sample and a composite image of the wrench taken from the 1st, 6th and 11th samples of the hammering with a wrench experiment. Figure 15 shows a plot of the triple (α, β, θ) with respect to time (frame numbers). We can see that the values of α are small and that $\theta \approx 0$ while β is close to $-\pi/2$.

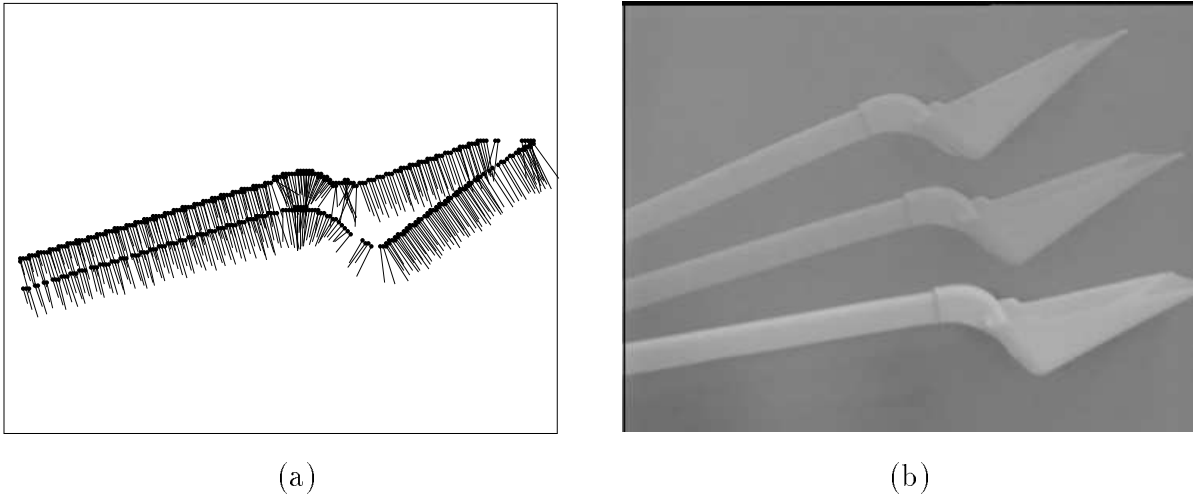


Figure 10: (a) Flow vectors for hitting with a shovel. (b) Hitting motion.

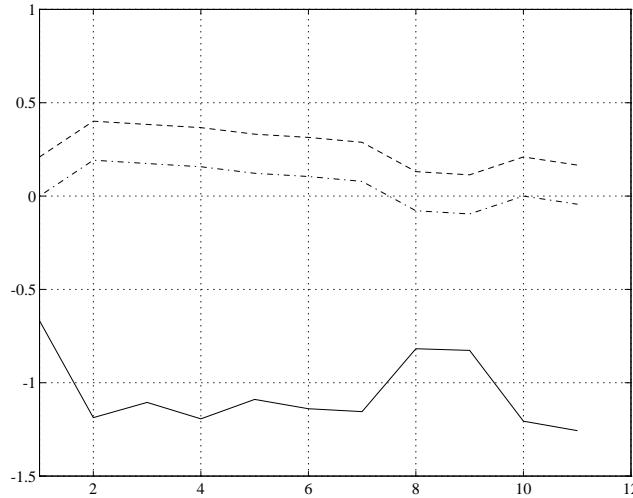


Figure 11: Angles α , β , and θ for hitting with a shovel. α is given by a dashed line, β is given by a solid line, and θ is given by a dash-dot line.

Tightening with the wrench is type of motion in which the rotational angle θ and the angle β (which corresponds to the translational direction in the object coordinate frame) have opposite signs. This is equivalent to saying that the head of the tool is fixed while the handle is moving. Figure 12 shows the flow vectors taken from the 6th sample and a composite image of the wrench taken from the 1st, 6th and 11th samples of the tightening with a wrench experiment. Figure 13 shows a plot of the triple (α, β, θ) with respect to time (frame numbers). We can see that the values of α are decreasing (this is equivalent to $\theta < 0$) while β is close to $\pi/2$.

Hammering with the wrench is a type of motion in which the translational part of the

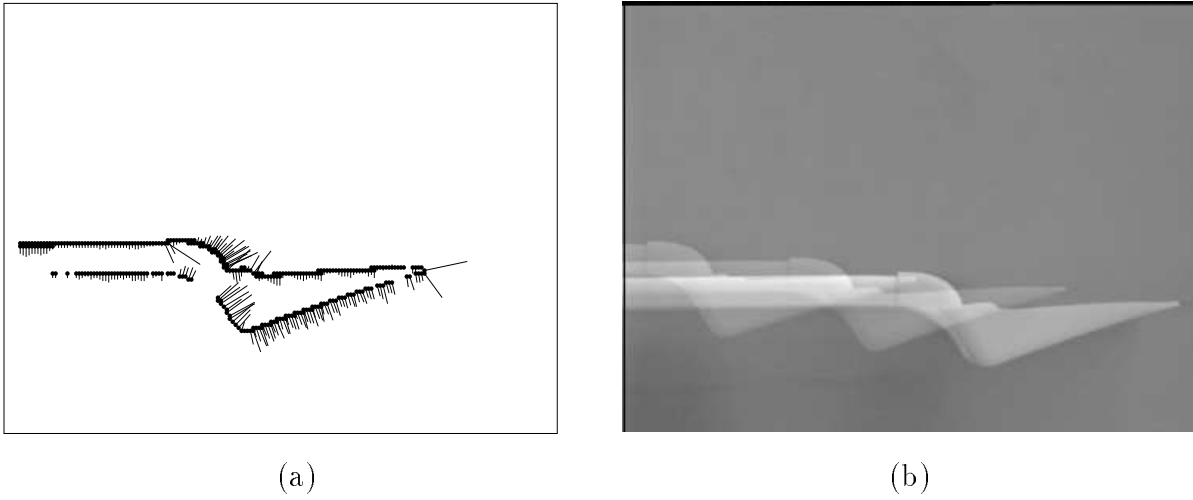


Figure 8: (a) Flow vectors for scooping with a Shovel. (b) Scooping motion.

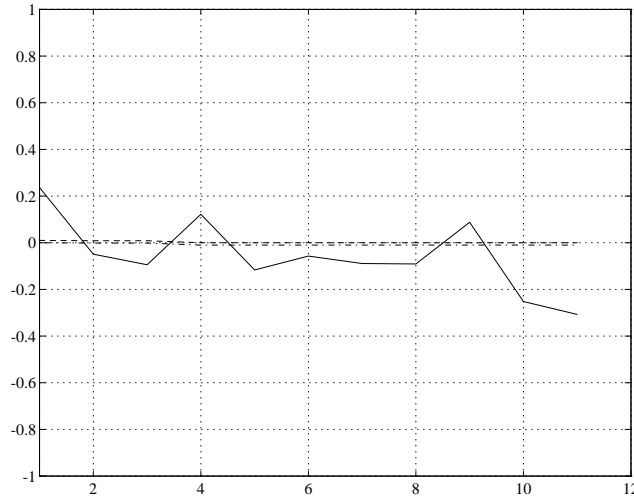


Figure 9: Angles α , β , and θ for scooping with a Shovel. α is given by a dashed line, β is given by a solid line, and θ is given by a dash-dot line.

triple (α, β, θ) with respect to time (frame numbers). We can see that the values of α are small and that $\theta \approx 0$ while β is close to $-\pi/2$.

6.3.2 Wrench

Two actions using a wrench were examined. In one experiment, the wrench was used to tighten a bolt; in the other sequence, the wrench was used as a hammer. In these cases the same tool is being used for multiple, inherently different functions. Motion analysis enables us to differentiate between the two.

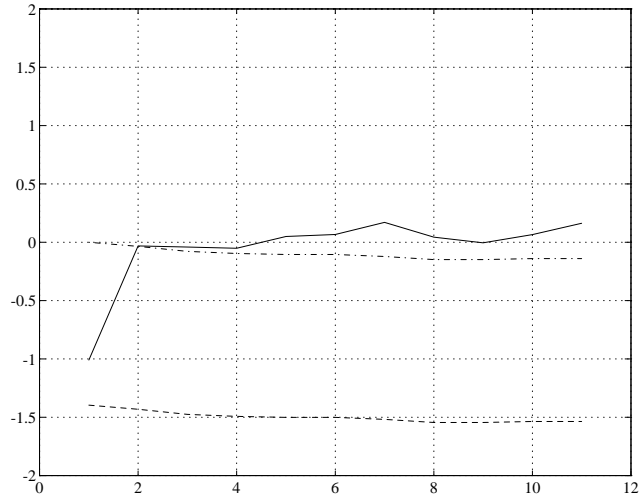


Figure 7: Angles α , β , and θ for Stabbing. α is given by a dashed line, β is given by a solid line, and θ is given by a dash-dot line.

We can see that the values of α are very close to $-\pi/2$, as was expected, β is close to 0 and θ is around 0.

6.3 Multi-Usage Objects

In this section we have two examples of multiple use of objects. We examine two actions using a shovel and two actions using a wrench.

6.3.1 Shovel

Two actions using a shovel were examined. In one experiment, the shovel was used in a scooping action; in the other sequence, it was used in a hitting action. In these cases the same tool is being used for two inherently different functions. Motion analysis enables us to differentiate between the two.

Scooping with a shovel is a type of motion in which the rotational angle θ is small and the angle β (which corresponds to the translational direction in the object coordinate frame) is small. Figure 8 shows the flow vectors taken from the 6th sample and a composite image of the shovel taken from the 1st, 6th and 11th samples of the scooping with a shovel experiment. Figure 9 shows a plot of the triple (α, β, θ) with respect to time (frame numbers). We can see that the values of θ are small while α and β are close to 0.

Hitting with the shovel is a type of motion in which the translational part of the motion dominates over the rotational part of the motion and the direction of translation is approximately orthogonal to the direction of the medial axis of the tool. Figure 10 shows the flow vectors taken from the 6th sample and a composite image of the shovel taken from the 1st, 6th and 11th samples of the hitting with a shovel experiment. Figure 11 shows a plot of the

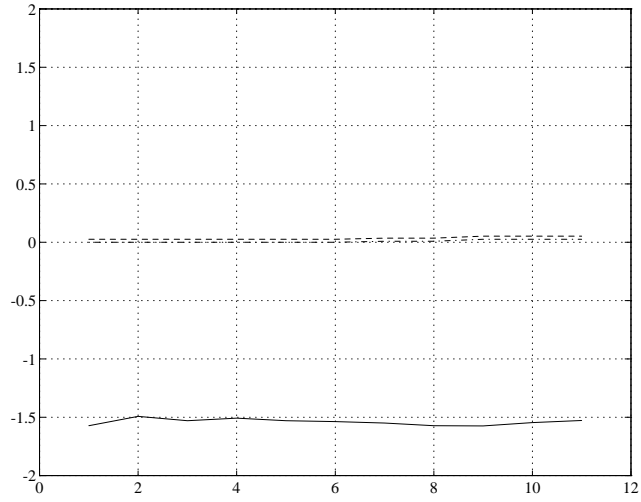


Figure 5: Angles α , β , and θ for Chopping. α is given by a dashed line, β is given by a solid line, and θ is given by a dash-dot line.

6.2.3 Stabbing

Stabbing is defined as the cutting motion of a knife in which α (the angle between the projection of l_c onto the plane $Z = Z_c$ and the Ox axis) is close to either $-\pi/2$ or $\pi/2$, β is approximately 0, and θ is small and approximately constant. The difference between jabbing and stabbing is in α . Figure 6 shows the flow vectors taken from the 6th sample



Figure 6: (a) Flow vectors for Stabbing. (b) Stabbing motion.

and a composite image of the knife taken from the 1st, 6th and 11th samples of the stabbing experiment. Figure 7 shows a plot of the triple (α, β, θ) with respect to time (frame numbers).

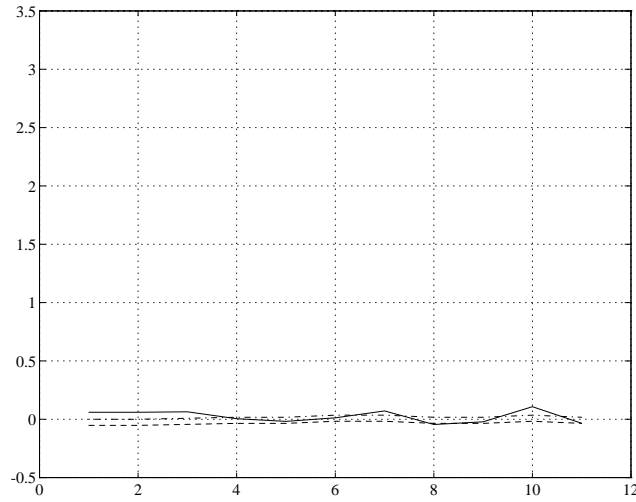


Figure 3: Angles α , β , and θ for Jabbing. α is given by a dashed line, β is given by a solid line, and θ is given by a dash-dot line.

6.2.2 Chopping

Chopping is defined as the cutting motion of a knife in which α (the angle between the projection of l_c onto the plane $Z = Z_c$ and the Ox axis) is close to either 0 or π , β is close to $\pi/2$ ($\alpha \approx \pi$) or $-\pi/2$ (when $\alpha \approx 0$), and θ is small and approximately constant. Figure 4

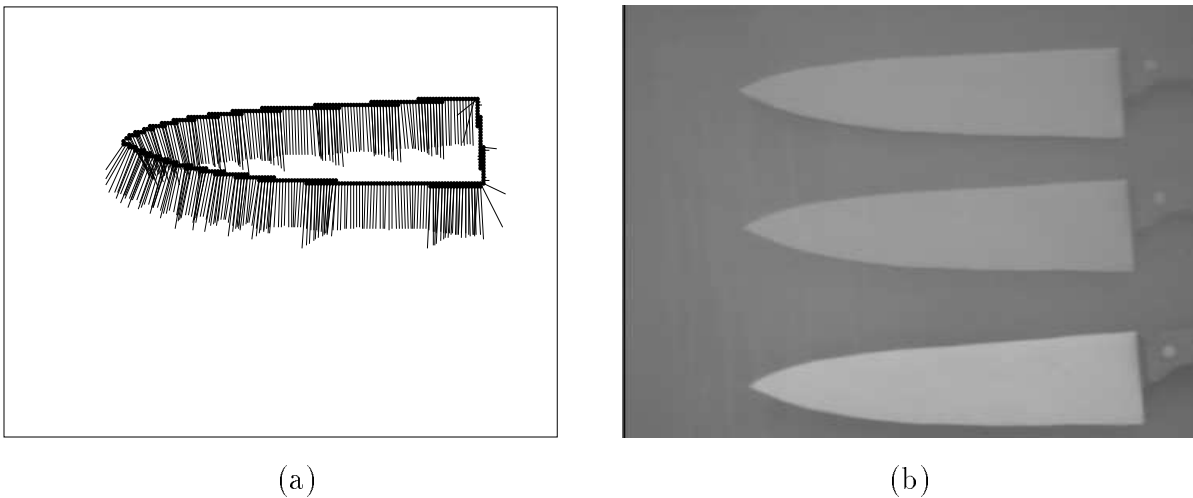


Figure 4: (a) Flow vectors for Chopping. (b) Chopping motion.

shows the flow vectors taken from the 6th sample and a composite image of the knife taken from the 1st, 6th and 11th samples of the chopping experiment. Figure 5 shows a plot of the triple (α, β, θ) with respect to time (frame numbers). We can see that the values of α are very close to 0, as was expected, β is close to $-\pi/2$ and θ is around 0.

For the following experiments we use these approximations to compute the object motion from images.

Let β be the angle between the vector $(U_1 \ V_1 \ 0)^T$ (the projection of \vec{T}_1 onto the plane $Z = Z_c$) and \vec{i}_c (the unit vector along the projection of the medial axis l_c onto the plane $Z = Z_c$). We have

$$\beta = \arctan \frac{V_1}{U_1}. \quad (19)$$

Let θ be the total rotation angle as a function of time. For a fronto-parallel surface the total rotation angle is approximately equal to the change in α and we have

$$\theta = \int_0^t C_1 dt \approx \alpha - \alpha_0. \quad (20)$$

We use the triples (α, β, θ) to recognize the functionalities of simple objects.

6.2 Action recognition for a class of manipulation tasks: Cutting

We start with three examples of simple functions performed by knives: chopping, jabbing and stabbing. In what follows we demonstrate how motion is used to differentiate between the three.

6.2.1 Jabbing

Jabbing is defined as the cutting motion of a knife in which α (the angle between the projection of l_c onto the plane $Z = Z_c$ and the Ox axis) is close to either 0 or π , β is approximately 0, and θ is small and approximately constant.

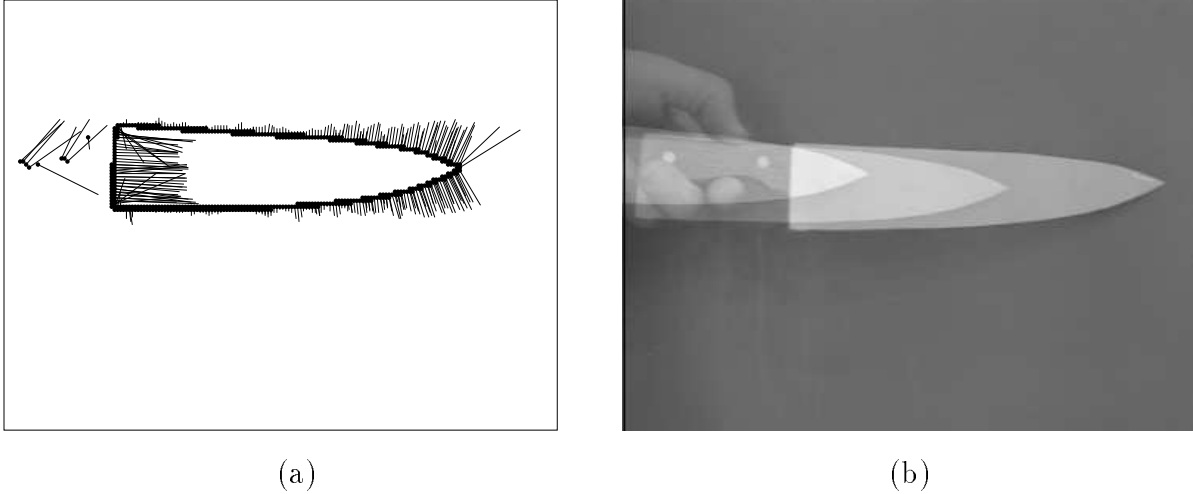


Figure 2: (a) Flow vectors for Jabbing. (b) Jabbing motion.

Figure 2 shows the flow vectors taken from the 6th sample and a composite image of the knife taken from the 1st, 6th and 11th samples of the jabbing experiment. Figure 3 shows a plot of the triple (α, β, θ) with respect to time (frame numbers). We can see that the values of α are very close to 0, as was expected, β is close to 0 and θ is around 0.

6.1 Motion Estimation from Normal Flow

In what follows we show how the different motion parameters defined in Section 5 can be estimated based on normal flow data computed from an image sequence.

Let $g_1(x, y) = (x - x_c) \cos \alpha - (y - y_c) \sin \alpha$. We can then define the vectors $\vec{\mathbf{a}}$ and $\vec{\mathbf{d}}$:

$$\vec{\mathbf{a}} = \begin{pmatrix} f n_x \\ f n_y \\ -x n_x - y n_y \\ n_y g_1(x, y) \\ -n_x g_1(x, y) \\ -n_x(y - y_c) + n_y(x - x_c) \end{pmatrix}, \quad \vec{\mathbf{d}} = \begin{pmatrix} U/Z_c \\ V/Z_c \\ W/Z_c \\ A \tan \varphi \\ B \tan \varphi \\ C \end{pmatrix}.$$

For a given $\vec{n}_r = n_x \vec{i} + n_y \vec{j}$ we then have from (10-11)

$$\dot{x} n_x + \dot{y} n_y = \vec{\mathbf{a}} \cdot \vec{\mathbf{d}}. \quad (17)$$

If we use the spatial image gradient as the normal direction $\vec{n}_r \equiv \nabla I / \|\nabla I\| = n_x \vec{i} + n_y \vec{j}$ and $\dot{r}_n \approx \dot{u}_n$ we can obtain an approximate equation by replacing the left hand side of (17) by normal flow $-I_t / \|\nabla I\|$. In this way we obtain one approximate equation in the six unknown elements of $\vec{\mathbf{d}}$. For each point (x_i, y_i) , $i = 1, \dots, N$ of the image at which $\|\nabla I(x_i, y_i, t)\|$ is large we can write one equation. If we have more than six points we have an over-determined system of equations $A \vec{\mathbf{d}} = \vec{\mathbf{b}}$; the rows of the $N \times 6$ matrix A are the vectors $\vec{\mathbf{a}}_i$. and the elements of the N -vector $\vec{\mathbf{b}}$ are $-(\partial I(x_i, y_i, t) / \partial t) / \|\nabla I(x_i, y_i, t)\|$.

We seek the solution for which $\|\vec{\mathbf{b}} - A \vec{\mathbf{d}}\|$ is minimal. This solution is the same as the solution of the system

$$A^T A \vec{\mathbf{d}} = A^T \vec{\mathbf{b}} \equiv \vec{\mathbf{e}}.$$

We solve the system $A^T A \vec{\mathbf{d}} = \vec{\mathbf{e}}$ using the Cholesky decomposition. Since the matrix $A^T A$ is a positive definite 6×6 matrix there exists a lower triangular matrix L such that $L L^T = A^T A$. We solve two triangular systems $L \vec{\mathbf{f}} = \vec{\mathbf{d}}$ and $L^T \vec{\mathbf{d}} = \vec{\mathbf{f}}$ to obtain the parameter vector $\vec{\mathbf{d}}$.

In the case when $\varphi \approx 0$ (fronto-parallel case) and the rotation B_1 around the Cy_1 axis is small the equations (10-11) become

$$\begin{aligned} \dot{x} &\approx \frac{Uf - xW}{Z_c} - C(y - y_c), \\ \dot{y} &\approx \frac{Vf - yW}{Z_c} + C(x - x_c). \end{aligned}$$

In this case we need to estimate only four parameters U/Z_c , V/Z_c , W/Z_c , and C in the parameter vector $\vec{\mathbf{d}}$; thus, $A^T A$ is a 4×4 matrix. We then have from (14-16) that $C_1 \approx C$ and from (12-13) we have

$$\begin{pmatrix} U_1/Z_c \\ V_1/Z_c \end{pmatrix} \approx \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix} \begin{pmatrix} U/Z_c \\ V/Z_c \end{pmatrix}, \quad \frac{W_1}{Z_c} = \frac{W}{Z_c}. \quad (18)$$

If $\text{sgn } c_4 = \text{sgn } c_6$ we have

$$C_1 = \text{sgn}(c_6)\sqrt{c_4c_6 + c_6^2}, \quad \varphi = \pm \arctan \sqrt{c_4/c_6}, \quad B_1 = \pm \frac{c_5}{\sqrt{c_4/c_6}} \quad (16)$$

where sgn is the sign function. Note that if $\text{sgn } c_5 = 1$ the signs of B_1 and φ must be equal, otherwise they must be different.

If translation is non-zero we can combine (13) and (16) to estimate φ and the motion parameters. When $U_1 \equiv 0$ we have $\tan \varphi = c_1/c_3$; when $W_1 \equiv 0$ we have $\tan \varphi = -c_3/c_1$. In both of these cases B_1 can be determined directly, as well as φ and C_1 . If $\varphi \approx 0$ (or $\varphi < 0.1$) then $C_1 \approx c_6$, and $|c_5| = |B_1| \tan \varphi \approx |B_1| \varphi < 0.1|B_1|$ and thus $|B_1| > 10|c_5|$. If (16) is used to compute φ we have two solutions for the $(U_1/Z_c, W_1/Z_c)$ pair.

6 Experiments

This section illustrates how our methods can be applied to real image sequences. In each sequence, we captured the motion of an object performing a task. The vision system used images at 25 frames per second for 5 seconds, yielding 125 images per experiment. After each image sequence was recorded, a representative sampling of the 125 images was used for further processing. Eleven evenly spaced samples, each composed of three consecutive images, were used.³ This resulted in 33 images for each experiment.

In our experiments we assumed a table-top scenario, with a stationary observer on one side of the table. Based on this assumption we used a coordinate system that was fixed to the center of the image, with the X axis horizontal and pointing toward the right side of the image, the Y axis pointing upward, and the Z axis chosen to yield a right-handed coordinate frame (pointing toward the scene). All measurements were made relative to this coordinate system. The focal length f of the camera was 550.

Estimation of the medial axis of the object was done by taking the median of all edge orientations at those points for which the normal flow was computed. We estimated (x_c, y_c) —the image position of C (the reference point and the center of mass of the object)—as the average of the coordinates of all edge points for which the normal flow was computed.

In the following subsections we describe our method of motion estimation for sticks and strips. This motion estimation procedure was used in two scenarios. In the first we show how the motion can be used to discriminate between different functionalities of the same object. All the functionalities belong to the same family of manipulation tasks, namely cutting. In the second scenario we again use motion information to differentiate between two different functionalities of the same object, but this time for two different families of manipulation tasks. We give two examples of this scenario. In the first example, a shovel is used for scooping or for hitting. In the second example, a wrench is used for tightening or for hammering. In both examples, the first use is the normal one and the second use is an instance of improvisation. The motion gives clear information for a correct interpretation of the action that is taking place.

³For instance, samples 1 and 2 in any given experiment used images 0–2 and 10–12, respectively.

This is an exact formula for thin planar strips; in the case of sticks this formula is exact for an occluding contour.

From (4–5), (7), and (9) we obtain the equations of projected motion for points on \mathcal{B} under weak perspective:

$$\dot{x} = \frac{Uf - xW}{Z_c} - C(y - y_c) - B \tan \varphi [(x - x_c) \cos \alpha - (y - y_c) \sin \alpha], \quad (10)$$

$$\dot{y} = \frac{Vf - yW}{Z_c} + C(x - x_c) + A \tan \varphi [(x - x_c) \cos \alpha - (y - y_c) \sin \alpha]. \quad (11)$$

Equations (10-11) relate the projected motion field, α , and (x_c, y_c) to the scaled translational velocity $Z_c^{-1}\vec{T} = Z_c^{-1}(U \ V \ W)^T$ and the three parameters of rotation and slant $(A \tan \varphi, B \tan \varphi, C)$.

Now, from (8) we have

$$Z_c^{-1} \begin{pmatrix} U_1 \cos \varphi + W_1 \sin \varphi \\ V_1 \\ -U_1 \sin \varphi + W_1 \cos \varphi \end{pmatrix} = R_Z^T(\alpha) \begin{pmatrix} U/Z_c \\ V/Z_c \\ W/Z_c \end{pmatrix} \equiv \begin{pmatrix} c_1 \\ c_2 \\ c_3 \end{pmatrix} \quad (12)$$

and by rearrangement we obtain

$$\frac{V_1}{Z_c} = c_2, \quad \begin{pmatrix} U_1/Z_c \\ W_1/Z_c \end{pmatrix} = \begin{pmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{pmatrix} \begin{pmatrix} c_1 \\ c_3 \end{pmatrix} \quad (13)$$

When $Z_c^{-1}\vec{T} = Z_c^{-1}(U \ V \ W)^T$ and α are known c_1 , c_2 , and c_3 are computable. From (13) we then can compute V_1/Z_c . However, φ and $(U_1/Z_c, W_1/Z_c)$ cannot be estimated without some additional assumptions. If we assume a fronto-parallel surface, i.e., $\varphi = 0$, we obtain $U_1/Z_c = c_1$ and $W_1/Z_c = c_3$; similarly, a bound on φ (e.g. $\varphi < \pi/6$) gives bounds on the motion parameters too. Finally, if we assume that $W_1 = 0$ we obtain $\varphi = -\arctan(c_3/c_1)$ and then from (13) we have U_1/Z_c ; similarly, if we assume that $U_1 = 0$ we obtain $\varphi = \arctan(c_1/c_3)$ and then from (13) we have W_1/Z_c .

From (7) we have

$$R_Z(\alpha) \begin{pmatrix} C_1 \tan \varphi \sin \varphi \\ B_1 \tan \varphi \\ C_1 \cos \varphi \end{pmatrix} = I_\varphi R_Z(\alpha) \begin{pmatrix} C_1 \sin \varphi \\ B_1 \\ C_1 \cos \varphi \end{pmatrix} = \begin{pmatrix} A \tan \varphi \\ B \tan \varphi \\ C \end{pmatrix} \quad (14)$$

where I_φ is a diagonal matrix with the first two diagonal elements $\tan \varphi$ and the third diagonal element 1; we have used the fact that I_φ and $R_Z(\alpha)$ are commutative matrices. From (14) we have

$$\begin{pmatrix} C_1 \tan \varphi \sin \varphi \\ B_1 \tan \varphi \\ C_1 \cos \varphi \end{pmatrix} = R_Z^T(\alpha) \begin{pmatrix} A \tan \varphi \\ B \tan \varphi \\ C \end{pmatrix} = \begin{pmatrix} c_4 \\ c_5 \\ c_6 \end{pmatrix}. \quad (15)$$

When $(A \tan \varphi \ B \tan \varphi \ C)^T$ and α are known c_4 , c_5 , and c_6 are computable. Since $\varphi \in [0, \pi/2)$ both c_4 and c_6 should have the same sign, otherwise we can assume that $\varphi = 0$.

have the center of mass C at the middle of its medial axis l_c : in this case l_c corresponds to the longest principal axis of the ellipsoid of inertia; the other two principal axes are orthogonal to l_c and can be chosen arbitrarily. We assume that there is no rotational velocity around l_c .

We choose the center of mass C of a stick or a strip \mathcal{B} as the origin of the object coordinate system $Cx_1y_1z_1$; the coordinates of C expressed in the fixed frame are (X_c, Y_c, Z_c) . We choose the unit vector \vec{v}_1 along l_c with the orientation chosen to be in the direction of the acting part of the tool. Let Π_{l_c} be the plane orthogonal to the plane $Z = Z_c$ in which the line l_c lies (we can obtain Π_{l_c} by sliding the line parallel to the \vec{k} along l_c). We chose \vec{k}_1 to lie in the plane Π_{l_c} with the orientation of \vec{k}_1 chosen so that $\vec{k} \cdot \vec{k}_1 \geq 0$; the unit vector \vec{j}_1 is then normal to the Π_{l_c} plane. We assume that strips are orthogonal to the Π_{l_c} plane.

The orthographic image of l_c in the plane $Z = Z_c$ is the line l'_c which is the intersection of the planes $Z = Z_c$ and Π_{l_c} ; let the unit vector in the direction of l'_c be \vec{v}_c and let it be oriented so that $\vec{v}_c \cdot \vec{v}_1 \geq 0$; and let the angle between l_c and l'_c be φ . The rotation $R_{Y_c}(\varphi)$ through the angle φ around the normal \vec{j}_1 of Π_{l_c} transforms \vec{v}_1 into \vec{v}_c and \vec{k}_1 into \vec{k} . The rotation $R_Z(\alpha)$ through the angle α (this is the angle between \vec{v}_c and \vec{v}) around the Oz axis transforms \vec{v}_c into \vec{v} . The rotation matrix $R = R_Z(\alpha)R_{Y_c}(\varphi)$ in (1) is then given by

$$R = \begin{pmatrix} \cos \alpha & \sin \alpha & 0 \\ -\sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \cos \varphi & 0 & \sin \varphi \\ 0 & 1 & 0 \\ -\sin \varphi & 0 & \cos \varphi \end{pmatrix}.$$

By our assumption about the rotational velocity and the choice of the object coordinate system we have $\vec{\omega}_1 = B_1\vec{j}_1 + C_1\vec{k}_1$. The expression for the rotational velocity in the fixed frame is given by

$$\vec{\omega} = \begin{pmatrix} A \\ B \\ C \end{pmatrix} = R(B_1\vec{j}_1 + C_1\vec{k}_1) = R_Z(\alpha) \begin{pmatrix} C_1 \sin \varphi \\ B_1 \\ C_1 \cos \varphi \end{pmatrix}. \quad (7)$$

Similarly, since the translational velocity of the object is $\vec{T}_1 = (U_1 \ V_1 \ W_1)^T$ and $\vec{T} = R_Z(\alpha)R_{Y_1}(\varphi)\vec{T}_1$ we have

$$\vec{T} = \begin{pmatrix} U \\ V \\ W \end{pmatrix} = R_Z(\alpha) \begin{pmatrix} U_1 \cos \varphi + W_1 \sin \varphi \\ V_1 \\ -U_1 \sin \varphi + W_1 \cos \varphi \end{pmatrix}. \quad (8)$$

We now consider the term $(Z - Z_c)/Z_c$ for the points on the object \mathcal{B} . The equations we derive are valid for points in the plane in which l_c lies and is orthogonal to Π_{l_c} ; the unit vector \vec{k}_1 is normal to this plane. The expression for \vec{k}_1 in the fixed frame is $R\vec{k}_1 = (\cos \alpha \sin \varphi \ -\sin \alpha \sin \varphi \ \cos \varphi)^T$. The equation of the plane orthogonal to $R\vec{k}_1$ and in which the point (X_c, Y_c, Z_c) lies is given by

$$(X - X_c) \cos \alpha \sin \varphi - (Y - Y_c) \sin \alpha \sin \varphi + (Z - Z_c) \cos \varphi = 0.$$

Multiplying by $f(Z_c \cos \varphi)^{-1}$ and using (3) we obtain

$$f \frac{Z - Z_c}{Z_c} = -(x - x_c) \cos \alpha \tan \varphi + (y - y_c) \sin \alpha \tan \varphi. \quad (9)$$

are important. When pure rotation is involved we have a screwdriver or a rotor (long and short axis respectively). When pure translation along the main axis is performed we have stabbing or jabbing.² For a plate, translation along the direction of the normal and rotation around the normal are important. In other directions there is no special component because a plate is isotropic. For a strip there are two important axes (it can be regarded as a plate and stick combined). Note that in all of these examples the motion is in a plane.

In this work, we are interested in the mapping $f : M \mapsto F$ from motion to function. Given a moving object as seen by an observer we would like to infer the function being performed by the acting agent. We are interested in the object's motion over time in the object's coordinate system and its relation to the object it acts on (the actee). Both of these measurements are necessary for the mapping. The object's motion over time in the object coordinate system gives us the relationship between the main axis of the object and its direction of motion. Given an object, these relationships help to determine the intended function. For example, we would expect the motion of a knife that a person is using to "cut" to be parallel to the main axis of the knife, whereas if the person is "chopping" with the knife we would expect motion perpendicular to the main axis.

When determining function from motion, attention must be paid to the intended recipient. The relation to the actee is essential for establishing the mapping and creating a frame of reference. Once this frame is established, motion of a knife in one direction could signify murder while motion in the opposite direction could signify suicide. Humans usually employ reference frames in which one axis represents the gravity vector, but this is not necessary. We can slice bread on a wall as well as on a table; what matters is the motion of the knife relative to the actee.

In the next section, we develop the motion estimation machinery needed for this class of examples and we formalize our procedure for obtaining $f : M \mapsto F$. In the following section, we present experimental results.

5 Motion of Sticks and Strips

Consider a moving object \mathcal{B} . There is an *ellipsoid of inertia* associated with \mathcal{B} . The center of the ellipsoid is at the center of mass C of \mathcal{B} ; the axes of the ellipsoid are called the *principal axes*. We associate the coordinate system $Cx_1y_1z_1$ with the ellipsoid and choose the axes of $Cx_1y_1z_1$ to be parallel to the principal axes. Let \vec{v}_1 be the unit vector in the direction of the longest axis l_c (this axis corresponds to the smallest principal moment of inertia); let \vec{k}_1 be the unit vector in the direction of the shortest principal axis (this axis corresponds to the largest moment of inertia); and let \vec{j}_1 be the unit vector in the direction of the remaining principal axis with the direction chosen so that the vectors $(\vec{v}_1, \vec{j}_1, \vec{k}_1)$ form a right-handed coordinate system.

In this paper we consider only planar and approximately straight strips and sticks. For a planar strip the axis of the maximal moment of inertia is orthogonal to the plane of the strip; if the strip is approximately straight, the axis of the minimal moment of inertia is approximately parallel to the medial axis l_c of the strip. In the case of a straight stick we

²When torsion is also involved we have screwing, drilling, etc.

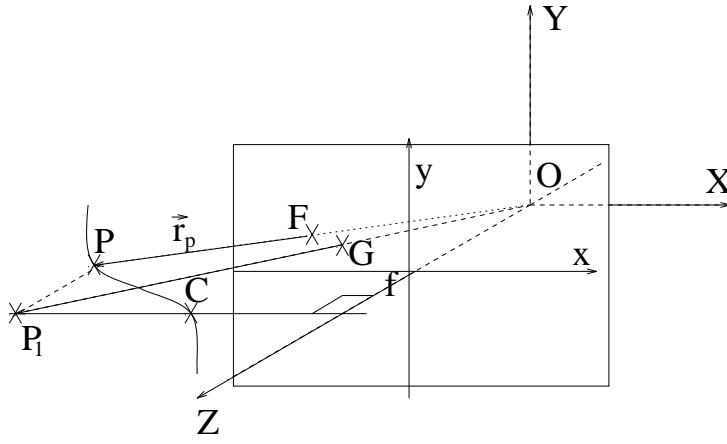


Figure 1: The plane perspective projection image of P is $F = f(X/Z, Y/Z, 1)$; the weak perspective projection image of P is obtained through the plane perspective projection of the intermediate point $P_1 = (X, Y, Z_c)$ and is given by $G = f(X/Z_c, Y/Z_c, 1)$.

where $(x_c, y_c) = (fX_c/Z_c, fY_c/Z_c)$ is the image of the point C . Let \vec{i} and \vec{j} be the unit vectors in the x and y directions, respectively; $\dot{\vec{r}} = \dot{x}\vec{i} + \dot{y}\vec{j}$ is the projected motion field at the point $\vec{r} = x\vec{i} + y\vec{j}$.

If we choose a unit direction vector \vec{n}_r in the image point \vec{r} and call it the normal direction, then the *normal motion field* at \vec{r} is $\dot{\vec{r}}_n = (\dot{\vec{r}} \cdot \vec{n}_r)\vec{n}_r$. \vec{n}_r can be chosen in various ways; the usual choice (as we shall now see) is the direction of the image intensity gradient.

Let $I(x, y, t)$ be the image intensity function. The time derivative of I can be written as

$$\frac{dI}{dt} = \frac{\partial I}{\partial x} \frac{dx}{dt} + \frac{\partial I}{\partial y} \frac{dy}{dt} + \frac{\partial I}{\partial t} = (I_x \vec{i} + I_y \vec{j}) \cdot (\dot{x}\vec{i} + \dot{y}\vec{j}) + I_t = \nabla I \cdot \dot{\vec{r}} + I_t$$

where ∇I is the image gradient and the subscripts denote partial derivatives.

If we assume $dI/dt = 0$, i.e. that the image intensity does not vary with time [5], then we have $\nabla I \cdot \vec{u} + I_t = 0$. The vector field \vec{u} in this expression is called the *optical flow*. If we choose the normal direction \vec{n}_r to be the image gradient direction, i.e. $\vec{n}_r \equiv \nabla I / \|\nabla I\|$, we then have

$$\vec{u}_n = (\vec{u} \cdot \vec{n}_r)\vec{n}_r = \frac{-I_t \nabla I}{\|\nabla I\|^2} \quad (6)$$

where \vec{u}_n is called the *normal flow*.

It was shown in [17] that the magnitude of the difference between \vec{u}_n and the normal motion field $\dot{\vec{r}}_n$ is inversely proportional to the magnitude of the image gradient. Hence $\dot{\vec{r}}_n \approx \vec{u}_n$ when $\|\nabla I\|$ is large. Equation (6) thus provides an approximate relationship between the 3-D motion and the image derivatives. We will use this approximation later in this paper.

4 Function from Motion

Basic or primitive motions (which can be rotational or translational) are motions relative to the main axes of a primitive object. For a stick, translation and rotation along the main axis

of P in $Cx_1y_1z_1$ we have the position \vec{r}_p of P in $Oxyz$

$$\vec{r}_p \equiv \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} = \begin{pmatrix} \vec{i} \cdot \vec{i}_1 & \vec{i} \cdot \vec{j}_1 & \vec{i} \cdot \vec{k}_1 \\ \vec{j} \cdot \vec{i}_1 & \vec{j} \cdot \vec{j}_1 & \vec{j} \cdot \vec{k}_1 \\ \vec{k} \cdot \vec{i}_1 & \vec{k} \cdot \vec{j}_1 & \vec{k} \cdot \vec{k}_1 \end{pmatrix} \begin{pmatrix} X_1 \\ Y_1 \\ Z_1 \end{pmatrix} + \begin{pmatrix} X_c \\ Y_c \\ Z_c \end{pmatrix} \equiv R\vec{p} + \vec{d}_c \quad (1)$$

where R is the matrix of the direction cosines (the frames are taken as right-handed so that $\det R = 1$). The velocity of \vec{r}_p is then given by

$$\dot{\vec{r}}_p = \vec{\omega} \times (\vec{r}_p - \vec{d}_c) + \vec{T}$$

where $\vec{\omega} = (A \ B \ C)^T$ is the rotational velocity of the moving frame; $\dot{\vec{d}}_c = (\dot{X}_c \ \dot{Y}_c \ \dot{Z}_c)^T \equiv (U \ V \ W)^T \equiv \vec{T}$ is the translational velocity of the point C . This can be written as

$$\begin{pmatrix} \dot{X} \\ \dot{Y} \\ \dot{Z} \end{pmatrix} = \begin{pmatrix} 0 & -C & B \\ C & 0 & -A \\ -B & A & 0 \end{pmatrix} \begin{pmatrix} X - X_c \\ Y - Y_c \\ Z - Z_c \end{pmatrix} + \begin{pmatrix} U \\ V \\ W \end{pmatrix}. \quad (2)$$

Let the rotational velocity in the moving frame be $\vec{\omega}_1 = (A_1 \ B_1 \ C_1)^T$; we can write $\vec{\omega} = R\vec{\omega}_1$ and $\vec{\omega}_1 = R^T\vec{\omega}$.

3.3 The Imaging Model

Let (X, Y, Z) denote the Cartesian coordinates of a scene point with respect to the fixed camera frame (see Figure 1), and let (x, y) denote the corresponding coordinates in the image plane. The equation of the image plane is $Z = f$, where f is the focal length of the camera. The perspective projection is given by $x = fX/Z$, and $y = fY/Z$. For weak perspective projection we need a reference point (X_c, Y_c, Z_c) . A scene point (X, Y, Z) is first projected onto the point (X, Y, Z_c) ; then, through plane perspective projection the point (X, Y, Z_c) is projected onto the image point (x, y) . The projection equations are then given by

$$x = \frac{X}{Z_c}f, \quad y = \frac{Y}{Z_c}f. \quad (3)$$

3.4 The Motion Field and the Optical Flow Field

The instantaneous velocity of the image point (x, y) under weak perspective projection can be obtained by taking derivatives of (3) with respect to time and using (2):

$$\begin{aligned} \dot{x} &= f \frac{\dot{X}Z_c - X\dot{Z}_c}{Z_c^2} = f \frac{[-C(Y - Y_c) + B(Z - Z_c) + U]Z_c - XW}{Z_c^2} \\ &= \frac{Uf - xW}{Z_c} - C(y - y_c) + fB \left(\frac{Z}{Z_c} - 1 \right), \end{aligned} \quad (4)$$

$$\begin{aligned} \dot{y} &= f \frac{\dot{Y}Z_c - Y\dot{Z}_c}{Z_c^2} = f \frac{[C(X - X_c) - A(Z - Z_c) + V]Z_c - YW}{Z_c^2} \\ &= \frac{Vf - yW}{Z_c} + C(x - x_c) - fA \left(\frac{Z}{Z_c} - 1 \right) \end{aligned} \quad (5)$$

3.1 Primitive shapes and primitive motions

Following [9, 10] we regard objects as composed of primitive parts. On the most coarse level we consider four types of primitive parts: sticks, strips, plates, and blobs, which differ in the values of their relative dimensions. If all three dimensions are about the same, we have a blob. If two are about the same, and the third is very different, we have two cases: if the two are bigger than the one, we have a plate, and in the reverse case we have a stick. When no two dimensions are about the same we have a strip. For example, a knife blade is a strip, because no two of its dimensions are similar.

These primitives can be combined to create compound objects. We can qualitatively describe the different ways in which primitives can be combined—for example, end to end, end to side, end to edge, etc. In addition to specifying the two attachment surfaces participating in the junction of two primitives, we could also consider the angles at which they join, and classify the joints as perpendicular, oblique, tangential, etc. Another refinement would be to describe qualitatively the position of the joint on each surface; an attachment can be near the middle, near a side, near a corner, or near an end of the surface. We can also specialize the primitives by adding qualitative features such as axis shape (straight or curved), cross-section size (constant or tapered), etc.

Functional recognition is based on compatibility with some action requirement. Some basic “actions” are static in nature (supporting, containing, etc.), but most actions involve using an object while it is moving. To illustrate the ways in which one can interact with a primitive, consider the action of “cutting” with a sharp strip or plate. Here a sharp edge is interacting with a surface. The interaction can be described from a kinematic point of view. The direction of motion of the primitive relative to its axis defines the action—for example, slicing or chopping. We define basic or primitive motions to be motions along, or perpendicular to, the main axes of a primitive object.¹ The motion can be a translation or a rotation.

3.2 Rigid Body Motion

To facilitate the derivation of the motion equations of a rigid body \mathcal{B} we use two rectangular coordinate frames, one $(Oxyz)$ fixed in space, the other $(Cx_1y_1z_1)$ fixed in the body and moving with it. The coordinates X_1, Y_1, Z_1 of any point P of the body with respect to the moving frame are constant with respect to time t , while the coordinates X, Y, Z of the same point P with respect to the fixed frame are functions of t . It is assumed that these functions are differentiable with respect to t . The position of the moving frame at any instant is given by the position $\vec{d}_c = (X_c \ Y_c \ Z_c)^T$ of the origin C , and by the nine direction cosines of the axes of the moving frame with respect to the fixed frame. Let $\vec{i}, \vec{j},$ and \vec{k} be the unit vectors in the directions of the $Ox, Oy,$ and Oz axes, respectively; and let $\vec{i}_1, \vec{j}_1,$ and \vec{k}_1 be the unit vectors in the directions of the $Cx_1, Cy_1,$ and Cz_1 axes, respectively. For a given position \vec{p}

¹It is interesting to note that motions along the main axis of a primitive preserve “degenerate views” [6].

in Section 5. In Section 6 we present experimental results demonstrating that motion analysis can indeed be used in determining functionality. In Section 7 we discuss planned future work in the area.

2 Related Work

Our research is concerned with the problem of determining the function of an object by analyzing its motion. Motion and functionality have appeared in the literature in several contexts. Early work on functional recognition can be found in [3, 11, 18]. More recently, Stark and Bowyer [12, 13, 14, 15] used these ideas to solve some of the problems presented by more traditional model-based methods of object recognition. In the so-called function-based approach, an object category is defined in terms of properties that an object must have in order to function as an instance of that category [14]. This work deals only with the stationary objects; no motion is involved.

Gould and Shah [4] use motion characteristics obtained from the extended trajectories followed by representative points on an object to identify important events corresponding to changes in direction, speed and acceleration. They believe that “in many cases where an object has a fixed and predefined motion, the trajectories of several points on the object may serve to uniquely identify the object”. This identification would be achieved by analyzing motion characteristics alone without requiring an object model; but no object identification results were given. We believe that since many objects display similar motion characteristics, motion alone is insufficient for function-based analysis, in cases where the function of an object is dependent not only on its motion, but also on its form.

Motion analysis for recognition of activities was described by Polana and Nelson [8]. They use Fourier analysis to detect and localize periodic activities such as walking or flying in a sequence of images. This work is similar in nature to our work in that both use motion as a basis for identifying activities. However, Polana and Nelson are concerned only with detecting the activities, without concern for the source of the motion. This is not adequate for function-based analysis since many objects can display similar motion characteristics. An object model is necessary to distinguish between the functions of objects that display similar motion characteristics.

Our work depends on segmenting the object into primitive parts and analyzing their motions. This kind of segmentation into functional parts was discussed by Rivlin et al. in [9]. They proposed a technique for functional recognition which extends the “Recognition by Parts” paradigm of object recognition to support “Recognition by Functional Parts”.

3 Preliminaries

In this section we begin with a discussion of primitive shapes and motions. Next, we derive equations of motion for both the observer-centered and the object-centered coordinate systems. We then derive projected motion equations for the plane perspective imaging model and show how these equations can be simplified by the use of weak perspective projection [16]. Finally, we derive the relationship between the image velocities and the projected motion.

1 Introduction

In the field of robotics, researchers have long pursued the goal of enabling a robot to act autonomously in its environment. For robots, as for humans, recognizing the functions of objects is a prerequisite to autonomous interaction with them. Functionality can be defined as the usability of an object for a particular purpose [1]. As an example, suppose we would like to open a letter. We seek a sharp object such as a knife or a pair of scissors that would be appropriate for opening the letter. Clearly the knife or scissors are functional in the context of opening a letter, and a robot given the task of opening a letter would at some point be required to recognize such objects as being functional for its task.

Recent research has focused on the problem of recognizing object functionality [1]. The goal of this research has been to determine functional capabilities of an object based on characteristics such as shape, physics and causation [14]. Little attention has been given to the problem of determining the functionality of an object from its motion. We believe that motion provides a strong indication of function. In particular, velocity, acceleration, and force of impact resulting from motion strongly constrain possible function. As in other approaches to functional recognition, the object (and in our case, its motion) should not be evaluated in isolation, but in context. The context includes the nature of the agent and the frame of reference it uses.

Information derived from motion can be useful in several ways. We expect a robot to take actions based on perceived events. In many instances, the events are perceived visually as motions in the environment. For example, a robot serving as a mechanic's mate [2] might "see" a person tightening a bolt with a pair of pliers and offer the person a wrench which would be more suitable for the task. Here, the robot determines the function of the pliers based on their motion (i.e. tightening) and determines that the wrench would be more suitable for tightening the bolt. This is an example of action perception. Additionally, a robot can learn object functionality by watching the object in use. This "visual learning" paradigm is a well-developed and vital component of biological visual systems [7]. As an example, a robot might "see" a knife being used to open a letter and learn the function of cutting and the context in which it can be used. In both of the above examples, the motion of the tool (i.e. tightening or cutting) is used to determine the function of the tool. Since the mapping between function and form is many to many, we need the information provided by motion to make the mapping more exact.

In this paper, we address the following problem: given a model of an object, how can we use the motion of the object, while it is being used to perform a task, to determine its function? Our method of answering this question takes into consideration the angular relationships between three vectors obtained from image sequence analysis. These vectors are compared with angular relationships that arise in known motion-to-function mappings. If a close enough match occurs, the functionality is identified; if not, a new functionality is learned.

In Section 2 we review literature that describes related work. In Section 3 we cover some preliminaries related to the problem. Section 4 considers the problem of determining the functionality of a known object by analyzing an image sequence showing that object performing the function. The motion estimation machinery needed for this task is developed

Function from Motion

Zoran Duric¹
Jeffrey Fayman²
Ehud Rivlin²

¹Computer Vision Laboratory
Center for Automation Research
University of Maryland
College Park, MD 20742-3275

²Department of Computer Science
Israel Institute of Technology – Technion
Haifa, Israel

Abstract

In order for a robot to operate autonomously in its environment, it must be able to perceive its environment and take actions based on these perceptions. Recognizing the functionalities of objects is an important component of this ability.

In this paper, we look into a new area of functionality recognition: determining the function of an object from its motion. Given a sequence of images of a known object performing some function, we attempt to determine what that function is. We show that the motion of an object, when combined with information about the object and its normal uses, provides us with strong constraints on possible functions that the object might be performing.

The support of the Air Force Office of Scientific Research under Grant F49620-93-1-0039 is gratefully acknowledged, as is the help of Sandy German in preparing this paper.