

Surface Orientation and Time to Contact from Image Divergence and Deformation

Roberto Cipolla* and Andrew Blake

Department of Engineering Science, University of Oxford, OX1 3PJ, England

Abstract. This paper describes a novel method to measure the differential invariants of the image velocity field robustly by computing average values from the integral of normal image velocities around image contours. This is equivalent to measuring the temporal changes in the area of a closed contour. This avoids having to recover a dense image velocity field and taking partial derivatives. It also does not require point or line correspondences. Moreover integration provides some immunity to image measurement noise.

It is shown how an *active* observer making small, deliberate motions can use the estimates of the divergence and deformation of the image velocity field to determine the object surface orientation and time to contact. The results of real-time experiments are presented in which arbitrary image shapes are tracked using B-spline snakes and the invariants are computed efficiently as closed-form functions of the B-spline snake control points. This information is used to guide a robot manipulator in obstacle collision avoidance, object manipulation and navigation.

1 Introduction

Relative motion between an observer and a scene induces deformation in image detail and shape. If these changes are smooth they can be economically described locally by the first order differential invariants of the image velocity field [16] – the curl (vorticity), divergence (dilatation), and shear (deformation) components. These invariants have simple geometrical meanings which do not depend on the particular choice of co-ordinate system. Moreover they are related to the three dimensional structure of the scene and the viewer's motion – in particular the surface orientation and the time to contact² – in a simple geometrically intuitive way. Better still, the divergence and deformation components of the image velocity field are unaffected by arbitrary viewer rotations about the viewer centre. They therefore provide an efficient, reliable way of recovering these parameters.

Although the analysis of the differential invariants of the image velocity field has attracted considerable attention [16, 14] their application to real tasks requiring visual inferences has been disappointingly limited [23, 9]. This is because existing methods have failed to deliver reliable estimates of the differential invariants when applied to real images. They have attempted the recovery of dense image velocity fields [4] or the accurate extraction of points or corner features [14]. Both methods have attendant problems concerning accuracy and numerical stability. An additional problem concerns the domain of

* Toshiba Fellow, Toshiba Research and Development Center, Kawasaki 210, Japan.

² The time duration before the observer and object collide if they continue with the same relative translational motion [10, 20]

applications to which estimates of differential invariants can be usefully applied. First order invariants of the image velocity field at a single point in the image cannot be used to provide a *complete* description of shape and motion as attempted in numerous structure from motion algorithms [27]. This in fact requires second order spatial derivatives of the image velocity field [21, 29]. Their power lies in their ability to efficiently recover reliable but incomplete solutions to the structure from motion problem which can be augmented with other information to accomplish useful visual tasks.

The reliable, real-time extraction of these invariants from image data and their application to visual tasks will be addressed in this paper. First we present a novel method to measure the differential invariants of the image velocity field robustly by computing average values from the integral of simple functions of the normal image velocities around image contours. This is equivalent to measuring the temporal changes in the area of a closed contour and avoids having to recover a dense image velocity field and taking partial derivatives. It also does not require point or line correspondences. Moreover integration provides some immunity to image measurement noise.

Second we show that the 3D interpretation of the differential invariants of the image velocity field is especially suited to the domain of *active vision* in which the viewer makes deliberate (although sometimes imprecise) motions, or in stereo vision, where the relative positions of the two cameras (eyes) are constrained while the cameras (eyes) are free to make arbitrary rotations (eye movements). Estimates of the divergence and deformation of the image velocity field, augmented with constraints on the direction of translation, are then sufficient to efficiently determine the object surface orientation and time to contact.

The results of preliminary real-time experiments in which arbitrary image shapes are tracked using B-spline snakes [6] are presented. The invariants are computed as closed-form functions of the B-spline snake control points. This information is used to guide a robot manipulator in obstacle collision avoidance, object manipulation and navigation.

2 Differential Invariants of the Image Velocity Field

2.1 Review

For a sufficiently small field of view (defined precisely in [26, 5]) and smooth change in viewpoint the image velocity field and the change in apparent image shape is well approximated by a linear (*affine*) transformation [16]. The latter can be decomposed into independent components which have simple geometric interpretations. These are an image translation (specifying the change in image position of the centroid of the shape); a 2D rigid rotation (vorticity), specifying the change in orientation, $\text{curl}\mathbf{v}$; an isotropic expansion (divergence) specifying a change in scale, $\text{div}\mathbf{v}$; and a pure shear or deformation which describes the distortion of the image shape (expansion in a specified direction with contraction in a perpendicular direction in such a way that area is unchanged) described by a magnitude, $\text{def}\mathbf{v}$, and the orientation of the axis of expansion (maximum extension), μ . These quantities can be defined as combinations of the partial derivatives of the image velocity field, $\mathbf{v} = (u, v)$, at an image point (x, y) :

$$\text{div}\mathbf{v} = (u_x + v_y) \tag{1}$$

$$\text{curl}\mathbf{v} = -(u_y - v_x) \tag{2}$$

$$(\text{def}\mathbf{v}) \cos 2\mu = (u_x - v_y) \tag{3}$$

$$(\text{def}\mathbf{v}) \sin 2\mu = (u_y + v_x) \tag{4}$$

where subscripts denote differentiation with respect to the subscript parameter. The curl, divergence and the magnitude of the deformation are scalar invariants and do not depend on the particular choice of image co-ordinate system [16, 14]. The axes of maximum extension and contraction change with rotations of the image plane axes.

2.2 Relation to 3D Shape and Viewer Motion

The differential invariants depend on the viewer motion (translational velocity, \mathbf{U} , and rotational velocity, $\mathbf{\Omega}$), depth, λ and the relation between the viewing direction (ray direction \mathbf{Q}) and the surface orientation in a simple and geometrically intuitive way. Before summarising these relationships let us define two 2D vector quantities: the component of translational velocity parallel to the image plane scaled by depth, λ , \mathbf{A} where:

$$\mathbf{A} = \frac{\mathbf{U} - (\mathbf{U} \cdot \mathbf{Q})\mathbf{Q}}{\lambda} \quad (5)$$

and the *depth gradient* scaled by depth³, \mathbf{F} , to represent the surface orientation and which we define in terms of the 2D vector gradient:

$$\mathbf{F} = \frac{\mathbf{grad}\lambda}{\lambda} \quad (6)$$

The magnitude of the depth gradient, $|\mathbf{F}|$, determines the tangent of the *slant* of the surface (angle between the surface normal and the visual direction). It vanishes for a frontal view and is infinite when the viewer is in the tangent plane of the surface. Its direction, $\angle\mathbf{F}$, specifies the direction in the image of increasing distance. This is equal to the *tilt* of the surface tangent plane. The exact relationship between the magnitude and direction of \mathbf{F} and the slant and tilt of the surface (σ, τ) is given by:

$$|\mathbf{F}| = \tan \sigma \quad (7)$$

$$\angle\mathbf{F} = \tau \quad (8)$$

With this new notation the relations between the differential invariants, the motion parameters and the surface position and orientation are given by [15]:

$$\text{curlv} = -2\mathbf{\Omega} \cdot \mathbf{Q} + |\mathbf{F} \wedge \mathbf{A}| \quad (9)$$

$$\text{divv} = \frac{2\mathbf{U} \cdot \mathbf{Q}}{\lambda} + \mathbf{F} \cdot \mathbf{A} \quad (10)$$

$$\text{defv} = |\mathbf{F}||\mathbf{A}| \quad (11)$$

where μ (which specifies the axis of maximum extension) bisects \mathbf{A} and \mathbf{F} :

$$\mu = \frac{\angle\mathbf{A} + \angle\mathbf{F}}{2} \quad (12)$$

The geometric significance of these equations is easily seen with a few examples. For example, a translation towards the surface patch leads to a uniform expansion in the

³ Koenderink [15] defines \mathbf{F} as a “nearness gradient”, $\mathbf{grad}(\log(1/\lambda))$. In this paper \mathbf{F} is defined as a scaled depth gradient. These two quantities differ by a sign.

image, i.e. positive divergence. This encodes the distance to the object which due to the speed-scale ambiguity⁴ is more conveniently expressed as a time to contact, t_c :

$$t_c = \frac{\lambda}{\mathbf{U} \cdot \mathbf{Q}} . \quad (13)$$

Translational motion perpendicular to the visual direction results in image deformation with a magnitude which is determined by the slant of the surface, σ and with an axis depending on the tilt of the surface, τ and the direction of the viewer translation. Divergence (due to foreshortening) and curl components may also be present.

Note that divergence and deformation are unaffected by (and hence insensitive to errors in) viewer rotations such as panning or tilting of the camera whereas these lead to considerable changes in point image velocities or disparities⁵. As a consequence the deformation component efficiently encodes the orientation of the surface while the divergence component can be used to provide an estimate of the time to contact or collision.

This formulation clearly exposes both the speed-scale ambiguity and the *bas-relief* ambiguity [11]. The latter manifests itself in the appearance of surface orientation, \mathbf{F} , with \mathbf{A} . Increasing the slant of the surface \mathbf{F} while scaling the movement by the same amount will leave the local image velocity field unchanged. Thus, from two weak perspective views and with no knowledge of the viewer translation, it is impossible to determine whether the deformation in the image is due to a large $|\mathbf{A}|$ (large “turn” of the object or “vergence angle”) and a small slant or a large slant and a small rotation around the object. Equivalently a nearby “shallow” object will produce the same effect as a far away “deep” structure. We can only recover the depth gradient \mathbf{F} up to an unknown scale. These ambiguities are clearly exposed with this analysis whereas this insight is sometimes lost in the purely algorithmic approaches to solving the equations of motion from the observed point image velocities. A consequence of the latter is the numerically ill-conditioned nature of structure from motion solutions when perspective effects are small.

3 Extraction of Differential Invariants

The analysis above treated the differential invariants as observables of the image. There are a number of ways of extracting the differential invariants from the image. These are summarised below before presenting a novel method based on the temporal derivatives of the moments of the area enclosed by a closed curve.

3.1 Summary of Existing Methods

1. Partial derivatives of image velocity field

⁴ Translational velocities appear scaled by depth making it impossible to determine whether the effects are due to a nearby object moving slowly or a far-away object moving quickly.

⁵ This is somewhat related to the reliable estimation of relative depth from the relative image velocities of two nearby points – motion parallax [21, 24, 6]. Both motion parallax and the deformation of the image velocity field relate local measurements of relative image velocities to scene structure in a simple way which is uncorrupted by the rotational image velocity component. In the case of parallax, the depths are discontinuous and differences of discrete velocities are related to the difference of inverse depths. Equation (11) on the otherhand assumes a smooth and continuous surface and derivatives of image velocities are related to derivatives of inverse depth.

This is the most commonly stressed approach. It is based on recovering a dense field of image velocities and computing the partial derivatives using discrete approximation to derivatives [17] or a least squares estimation of the affine transformation parameters from the image velocities estimated by spatio-temporal methods [23, 4]. The recovery of the image velocity field is usually computationally expensive and ill-conditioned [12].

2. Point velocities in a small neighbourhood

The image velocities of a minimum of three points in a small neighbourhood are sufficient, in principle, to estimate the components of the affine transformation and hence the differential invariants [14, 18]. In fact it is only necessary to measure the change in area of the triangle formed by the three points and the orientations of its sides [7]. There is, however, no redundancy in the data and hence this method requires very accurate image positions and velocities. In [7] this is attempted by tracking large numbers of "corner" features [28] and using Delaunay triangulation [3] in the image to approximate the physical world by planar facets. Preliminary results showed that the localisation of "corner" features was insufficient for reliable estimation of the differential invariants.

3. Relative Orientation of Line Segments

Koenderink [15] showed how temporal texture density changes can yield estimates of the divergence. He also presented a method for recovering the curl and shear components that employs the orientations of texture elements. Orientations are not affected by the divergence term. They are only affected by the curl and deformation components. In particular the curl component changes all the orientations by the same amount. It does not affect the angles between the image edges. These are only affected by the deformation component. The relative changes in orientation can be used to recover deformation in a simple way since the effects of the curl component are cancelled out. Measurement at three oriented line segments is sufficient to completely specify the deformation components. The main advantage is that point velocities or partial derivatives are not required.

4. Curves and Closed Contours

The methods described above require point and line correspondences. Sometimes these are not available or are poorly localised. Often we can only reliably extract portions of curves (although we can not always rely on the end points) or closed contours.

Image shapes or contours only "sample" the image velocity field. At contour edges it is only possible to measure the normal component of image velocity. This information can in certain cases be used to recover the image velocity field. Waxman and Wohn [30] showed how to recover the full velocity field from the normal components at image contours. In principle, measurement of eight normal velocities around a contour allow the characterisation of the full velocity field for a planar surface. Kanatani [13] related line integrals of image velocities around closed contours to the motion and orientation parameters of a planar contour. In the following we will not attempt to solve for these structure and motion parameters directly but only to recover the divergence and deformation.

3.2 Recovery of Invariants from Area Moments of Closed Contours

It has been shown that the differential invariants of the image velocity field conveniently characterise the changes in apparent shape due to relative motion between the viewer and scene. Contours in the image sample this image velocity field. It is usually only possible,

however, to recover the normal image velocity component from local measurements at a curve [27, 12]. It is now shown that this information is often sufficient to estimate the differential invariants within closed curves.

Our approach is based on relating the temporal derivative of the area of a closed contour and its moments to the invariants of the image velocity field. This is a generalisation of the result derived by Maybank [22] in which the rate of change of area scaled by area is used to estimate the divergence of the image velocity field. The advantage is that point or line correspondences are not used. Only the correspondence between shapes is required. The computationally difficult, ill-conditioned and poorly defined process of making explicit the full image velocity field [12] is avoided. Moreover, since taking temporal derivatives of area (and its moments) is equivalent to the integration of normal image velocities (scaled by simple functions) around closed contours our approach is effectively computing average values of the differential invariants (not point properties) and has better immunity to image noise leading to reliable estimates. Areas can also be estimated accurately, even when the full set of first order derivatives can not be obtained.

The moments of area of a contour, I_f , are defined in terms of an area integral with boundaries defined by the contour in the image plane:

$$I_f = \int_{a(t)} f dx dy \quad (14)$$

where $a(t)$ is the area of a contour of interest at time t and f is a scalar function of image position (x, y) that defines the moment of interest. For instance setting $f = 1$ gives us area. Setting $f = x$ or $f = y$ gives the first-order moments about the image x and y axes respectively.

The moments of area can be measured directly from the image (see below for a novel method involving the control points of a B-spline snake). Better still, their temporal derivatives can also be measured. Differentiating (14) with respect to time and using a result from calculus⁶ we can relate the temporal derivative of the moment of area to an integral of the normal component of image velocities at an image contour, $\mathbf{v} \cdot \mathbf{n}$, weighted by a scalar $f(x, y)$. By Green's theorem, this integral around the contour $c(t)$, can be re-expressed as an integral over the area enclosed by the contour, $a(t)$.

$$\frac{d}{dt}(I_f) = \oint_{c(t)} [f \mathbf{v} \cdot \mathbf{n}] ds \quad (15)$$

$$= \int_{a(t)} [\text{div}(f \mathbf{v})] dx dy \quad (16)$$

$$= \int_{a(t)} [f \text{div} \mathbf{v} + (\mathbf{v} \cdot \text{grad} f)] dx dy \quad (17)$$

If the image velocity field, \mathbf{v} , can be represented by constant partial derivatives in the area of interest, substituting the coefficients of the affine transformation for the velocity field into (17) leads to a linear equation in which the left hand side is the temporal derivative of the moment of area described by f (which can be measured, see below) while the integrals on the right-hand side are moments of area (also directly measurable). The coefficients of each term are the required parameters of the affine transformation. In summary, the

⁶ This equation can be derived by considering the *flux* linking the area of the contour. This changes with time since the contour is carried by the velocity field. The *flux* field, f , in our example does not change with time. Similar integrals appear in fluid mechanics, e.g. the *flux transport theorem* [8].

image velocity field deforms the shape of contours in the image. Shape can be described by moments of area. Hence measuring the change in the moments of area is an alternative way characterising the transformation. In this way the change in the moments of area have been expressed in terms of the parameters of the affine transformation.

If we initially set up the $x - y$ co-ordinate system at the centroid of the image contour of interest so that the first moments are zero, (17) with $f = x$ and $f = y$ shows that the centroid of the deformed shape specifies the mean translation $[u_0, v_0]$. Setting $f = 1$ leads to the simple and useful result that the divergence of the image velocity field can be estimated as the derivative of area scaled by area:

$$\frac{d}{dt}a(t) = a(t)\text{div}\mathbf{v} \quad (18)$$

Increasing the order of the moments, i.e. different values of $f(x, y)$, generates new equations and additional constraints. In principle, if it is possible to find six linearly independent equations, we can solve for the affine transformation parameters and combine the co-efficients to recover the differential invariants. The validity of the affine approximation can be checked by looking at the error between the transformed and observed image contours. The choice of which moments to use is a subject for further work. Listed below are some of the simplest equations which have been useful in the experiments presented here.

$$\frac{d}{dt} \begin{bmatrix} a \\ I_x \\ I_y \\ I_{x^2} \\ I_{y^2} \\ I_{x^3y} \end{bmatrix} = \begin{bmatrix} 0 & 0 & a & 0 & 0 & a \\ a & 0 & 2I_x & I_y & 0 & I_x \\ 0 & a & I_y & 0 & I_x & 2I_y \\ 2I_x & 0 & 3I_{x^2} & 2I_{xy} & 0 & I_{x^2} \\ 0 & 2I_y & I_{y^2} & 0 & 2I_{xy} & 3I_{y^2} \\ 3I_{x^2y} & I_{x^3} & 4I_{x^3y} & 3I_{x^2y^2} & I_{x^4} & 2I_{x^3y} \end{bmatrix} \begin{bmatrix} u_0 \\ v_0 \\ u_x \\ u_y \\ v_x \\ v_y \end{bmatrix} \quad (19)$$

(Note that in this equation subscripts are used to label the moments of area. The left-hand side represents the temporal derivative of the moments in the column vector.) In practice certain contours may lead to equations which are not independent and their solution is ill-conditioned. The interpretation of this is that the normal components of image velocity are insufficient to recover the true image velocity field globally, e.g. a fronto-parallel circle rotating about the optical axis. This was termed the ‘‘aperture problem in the large’’ by Waxman and Wohn [30] and investigated by Berghom and Carlsson [2]. Note however, that it is always possible to recover the divergence from a closed contour.

4 Recovery of Surface Orientation and Time to Contact

Applications of the estimates of the image divergence and deformation of the image velocity field are summarised below. It has already been noted that measurement of the differential invariants in a single neighbourhood is insufficient to completely solve for the structure and motion since (9,10,11,12) are four equations in the six unknowns of scene structure and motion. In a single neighbourhood a complete solution would require the computation of second order derivatives [21, 29] to generate sufficient equations to solve for the unknowns. Even then the solution of the resulting set of non-linear equations is non-trivial.

In the following, the information available from the first-order differential invariants alone is investigated. It will be seen that the differential invariants are sufficient to constrain surface position and orientation and that this partial solution can be used to

perform useful visual tasks when augmented with additional information. Useful applications include providing information which is used by pilots when landing aircraft [10], estimating time to contact in braking reactions [20] and in the recovery of 3D shape up to a relief transformation [18, 19]. We now show how surface orientation and position (expressed as a time to contact) can be recovered from the estimates of image divergence and the magnitude and axis of the deformation.

1. With knowledge of translation but arbitrary rotation

An estimate of the direction of translation is usually available when the viewer is making deliberate movements (in the case of active vision) or in the case of binocular vision (where the camera or eye positions are constrained). It can also be estimated from image measurements by motion parallax [21, 24].

If the viewer translation is known, (10), (11) and (12) are sufficient to unambiguously recover the surface orientation and the distance to the object in temporal units. Due to the speed-scale ambiguity the latter is expressed as a time to contact. A solution can be obtained in the following way.

- (a) The axis of expansion (μ) of the deformation component and the projection in the image of the direction of translation ($\angle A$) allow the recovery of the tilt of the surface from (12).
- (b) We can then subtract the contribution due to the surface orientation and viewer translation parallel to the image axis from the image divergence (10). This is equal to $|\text{defv}| \cos(\tau - \angle A)$. The remaining component of divergence is due to movement towards or away from the object. This can be used to recover the time to contact, t_c . This can be recovered despite the fact that the viewer translation may not be parallel to the visual direction.
- (c) The time to contact fixes the viewer translation in temporal units. It allows the specification of the magnitude of the translation parallel to the image plane (up to the same speed-scale ambiguity), A . The magnitude of the deformation can then be used to recover the slant, σ , of the surface from (11).

The advantage of this formulation is that camera rotations do not affect the estimation of shape and distance. The effects of errors in the direction of translation are clearly evident as scalings in depth or by a relief transformation [15].

2. With fixation

If the cameras or eyes rotate to keep the object of interest in the middle of the image (null the effect of image translation) the magnitude of the rotations needed to bring the object back to the centre of the image determines A and hence allows us to solve for surface orientation, as above. Again the major effect of any error in the estimate of rotation is to scale depth and orientations.

3. With no additional information – constraints on motion

Even without any additional assumptions it is still possible to obtain useful information from the first-order differential invariants. The information obtained is best expressed as bounds. For example inspection of (10) and (11) shows that the time to contact must lie in an interval given by:

$$\frac{1}{t_c} = \frac{\text{divv}}{2} \pm \frac{\text{defv}}{2} . \quad (20)$$

The upper bound on time to contact occurs when the component of viewer translation parallel to the image plane is in the opposite direction to the depth gradient. The lower bound occurs when the translation is parallel to the depth gradient. The upper and lower estimates of time to contact are equal when there is no deformation

component. This is the case in which the viewer translation is along the ray or when viewing a fronto-parallel surface (zero depth gradient locally). The estimate of time to contact is then exact. A similar equation was recently described by Subbarao [25].

4. With no additional information – the constraints on 3D shape

Koenderink and Van Doorn [18] showed that surface shape information can be obtained by considering the variation of the deformation component alone in small field of view when weak perspective is a valid approximation. This allows the recovery of 3D shape up to a scale and relief transformation. That is they effectively recover the axis of rotation of the object but not the magnitude of the turn. This yields a family of solutions depending on the magnitude of the turn. Fixing the latter determines the slants and tilts of the surface. This has recently been extended in the affine structure from motion theorem [19].

The solutions presented above use knowledge of a single viewer translation and measurement of the divergence and deformation of the image velocity field. An alternative solution exists if the observer is free to translate along the ray and also in two orthogonal directions parallel to the image plane. In this case measurement of divergence alone is sufficient to recover the surface orientation and the time to contact.

5 Implementation and Experimental Results

5.1 Tracking Closed Loop Contours

The implementation and results follow. Multi-span closed loop B-spline snakes [6] are used to localise and track closed image contours. The B-spline is a curve in the image plane

$$\mathbf{x}(s) = \sum_i f_i(s) \mathbf{q}_i \quad (21)$$

where f_i are the spline basis functions with coefficients \mathbf{q}_i (control points of the curve) and s is a curve parameter (not necessarily arc length)[1]. The snakes are initialised as points in the centre of the image and are forced to expand radially outwards until they were in the vicinity of an edge where image “forces” make the snake stabilise close to a high contrast closed contour. Subsequent image motion is automatically tracked by the snake [5].

B-spline snakes have useful properties such as local control and continuity. They also compactly represent image curves. In our applications they have the additional advantage that the area enclosed is a simple function of the control points. This also applies to the other area moments. From Green’s theorem in the plane it is easy to show that the area enclosed by a curve with parameterisation $x(s)$ and $y(s)$ is given by:

$$a = \int_{s_0}^{s_N} x(s)y'(s)ds \quad (22)$$

where $x(s)$ and $y(s)$ are the two components of the image curve and $y'(s)$ is the derivative with respect to the curve parameter s . For a B-spline, substituting (21) and its derivative:

$$a(t) = \int_{s_0}^{s_N} \sum_i \sum_j (q_{x_i} q_{y_j}) f_i f'_j ds \quad (23)$$

$$= \sum_i \sum_j (q_{x_i} q_{y_j}) \int_{s_0}^{s_N} f_i f'_j ds. \quad (24)$$

Note that for each span of the B-spline and at each time instant the basis functions remain unchanged. The integrals can thus be computed off-line in closed form. (At most 16 coefficients need be stored. In fact due to symmetry there are only 10 possible values for a cubic B-spline). At each time instant multiplication with the control point positions gives the area enclosed by the contour. This is extremely efficient, giving the exact area enclosed by the contour. The same method can be used for higher moments of area as well. The temporal derivatives of the area and its moments is then used to estimate image divergence and deformation.

5.2 Applications

Here we present the results of a preliminary implementation of the theory. The examples are based on a camera mounted on a robot arm whose translations are deliberate while the rotations around the camera centre are performed to keep the target of interest in the centre of its field of view. The camera intrinsic parameters (image centre, scaling factors and focal length) and orientation are unknown. The direction of translation is assumed known and expressed with bounds due to uncertainty.

Braking Figure 1 shows four samples from a sequence of images taken by a moving observer approaching the rear windscreen of a stationary car in front. In the first frame (time $t = 0$) the relative distance between the two cars is approximately 7m. The velocity of approach is uniform and approximately 1m/time unit.

A B-spline snake is initialised in the centre of the windscreen, and expands out until it localises the closed contour of the edge of the windscreen. The snake can then automatically track the windscreen over the sequence. Figure 2 plots the apparent area, $a(t)$ (relative to the initial area, $a(0)$) as a function of time, t . For uniform translation along the optical axis the relationship between area and time can be derived from (10) and (18) by solving the first-order partial differential equation:

$$\frac{d}{dt}(a(t)) = \left(\frac{2\mathbf{U} \cdot \mathbf{Q}}{\lambda} \right) a(t) . \quad (25)$$

Its solution is given by:

$$a(t) = \frac{a(0)}{\left[1 - \frac{t}{t_c(0)} \right]^2} \quad (26)$$

where $t_c(0)$ is the time to contact at time $t = 0$:

$$t_c(0) = \frac{\lambda(0)}{\mathbf{U} \cdot \mathbf{Q}} . \quad (27)$$

This is in close agreement with the data (Fig. 2a). This is more easily seen if we look at the variation of the time to contact with time. For uniform motion this should decrease linearly. The experimental results are plotted in Fig. 2b. These are obtained by dividing the area of the contour at a given time by its temporal derivative (estimated by finite differences). The variation is linear, as predicted. These results are of useful accuracy, predicting the collision time to the nearest half time unit (corresponding to 50cm in this example).

For non-uniform motion the profile of the time to contact as a function of time is a very important cue for braking and landing reactions [20].

Collision avoidance It is well known that image divergence can be used in obstacle collision avoidance. Nelson and Aloimonos [23] demonstrated a robotics system which computed divergence by spatio-temporal techniques applied to the images of highly textured visible surfaces. We describe a real-time implementation based on image contours and “act” on the visually derived information.

Figure 3 shows the results of a camera mounted on an Adept robot manipulator and pointing in the direction of a target contour. (We hope to extend this so that the robot initially searches by rotation for a contour of interest. In the present implementation, however, the target object is placed in the centre of the field of view.) The closed contour is then localised automatically by initialising a closed loop B-spline snake in the centre of the image. The snake “explodes” outwards and deforms under the influence of image forces which cause it to be attracted to high contrast edges.

The robot manipulator then makes a deliberate motion towards the target. Tracking the area of the contour and computing its rate of change allows us to estimate the divergence. For motion along the visual ray this is sufficient information to estimate the time to contact or impact. The estimate of time to contact – decreased by the uncertainty in the measurement and any image deformation (20) – can be used to guide the manipulator so that it stops just before collision (Fig. 3d). The manipulator in fact, travels “blindly” after its sensing actions (above) and at a uniform speed for the time remaining until contact. In repeated trials image divergences measured at distances of 0.5m to 1.0m were estimated accurately to the nearest half of a time unit. This corresponds to a positional accuracy of 20mm for a manipulator translational velocity of 40mm/s.

The affine transformation approximation breaks down at close proximity to the target. This may lead to a degradation in the estimate of time to contact when very close to the target.

Landing reactions and object manipulation If the translational motion has a component parallel to the image plane, the image divergence is composed of two components. The first is the component which determines immediacy or time to contact. The other term is due to image foreshortening when the surface has a non-zero slant. The two effects can be computed separately by measuring the deformation. The deformation also allows us to recover the surface orientation.

Note that unlike stereo vision, the magnitude of the translation is not needed. Nor are the camera parameters (focal length and aspect ratio is not needed for divergence) known or calibrated. Nor are the magnitudes and directions of the camera rotations needed to keep the target in the field of view. Simple measurements of area and its moments – obtained in closed form as a function of the B-spline snake control points – were used to estimate divergence and deformation. The only assumption was of uniform motion and known direction of translation.

Figure 3 shows an example in which a robot manipulator uses these estimates of time to contact and surface orientation to approach the object surface perpendicularly so as to position a suction gripper for manipulation. The image contours are shown in Fig. 3a and 3b highlighting the effect of deformation due to the sideways component of translation. The successful execution is shown in Fig. 3c and 3d.

Qualitative visual navigation Existing techniques for visual navigation have typically used stereo or the analysis of image sequences to determine the camera ego-motion and then the 3D positions of feature points. The 3D data are then analysed to determine, for example, navigable regions, obstacles or doors. An example of an alternative approach

is presented. This computes qualitative information about the orientation of surfaces and times to contact from estimates of image divergence and deformation. The only requirement is that the viewer can make deliberate movements or has stereoscopic vision. Figure 4a shows the image of a door and an object of interest, a pallet. Movement towards the door and pallet produce a deformation in the image. This is seen as an expansion in the apparent area of the door and pallet in Fig. 4b. This can be used to determine the distance to these objects, expressed as a time to contact – the time needed for the viewer to reach the object if the viewer continued with the same speed. The image deformation is not significant. Any component of deformation can, anyhow, be absorbed by (20) as a bound on the time to contact. A movement to the left (Fig. 4c) produces image deformation, divergence and rotation. This is immediately evident from both the door (positive deformation and a shear with a horizontal axis of expansion) and the pallet (clockwise rotation with shear with diagonal axis of expansion). These effects with the knowledge of the direction of translation between the images taken at figure 4a and 4c are consistent with the door having zero tilt, i.e. horizontal direction of increasing depth, while the pallet has a tilt of approximately 90° , i.e. vertical direction of increasing depth. These are the effects predicted by (9, 10, 11 and 12) even though there are also strong perspective effects in the images. They are sufficient to determine the orientation of the surface qualitatively (Fig. 4d). This has been done without knowledge of the intrinsic properties of the cameras (camera calibration), the orientations of the cameras, their rotations or translational velocities. No knowledge of epipolar geometry is used to determine exact image velocities or disparities. The solution is incomplete. It can, however, be easily augmented into a complete solution by adding additional information. Knowing the magnitude of the sideways translational velocity, for example, can determine the exact quantitative orientations of the visible surfaces.

6 Conclusions

We have presented a simple and efficient method for estimating image divergence and deformation by tracking closed image contours with B-spline snakes. This information has been successfully used to estimate surface orientation and time to contact.

Acknowledgements

The authors acknowledge discussions with Mike Brady, Kenichi Kanatani, Christopher Longuet-Higgins, and Andrew Zisserman. This work was partially funded by Esprit BRA 3274 (FIRST) and the SERC. Roberto Cipolla also gratefully acknowledges the support of the IBM UK Scientific Centre, St. Hugh's College, Oxford and the Toshiba Research and Development Centre, Japan.

References

1. R.H. Bartels, J.C. Beatty, and B.A. Barsky. *An Introduction to Splines for use in Computer Graphics and Geometric Modeling*. Morgan Kaufmann, 1987.
2. F. Bergholm. Motion from flow along contours: a note on robustness and ambiguous case. *Int. Journal of Computer Vision*, 3:395–415, 1989.
3. J.D. Boissonat. Representing solids with the delaunay triangulation. In *Proc. ICPR*, pages 745–748, 1984.

4. M. Campani and A. Verri. Computing optical flow from an overconstrained system of linear algebraic equations. In *Proc. 3rd Int. Conf. on Computer Vision*, pages 22–26, 1990.
5. R. Cipolla. *Active Visual Inference of Surface Shape*. PhD thesis, University of Oxford, 1991.
6. R. Cipolla and A. Blake. The dynamic analysis of apparent contours. In *Proc. 3rd Int. Conf. on Computer Vision*, pages 616–623, 1990.
7. R. Cipolla and P. Kovesi. Determining object surface orientation and time to impact from image divergence and deformation. (University of Oxford (Memo)), 1991.
8. H.F. Davis and A.D. Snider. *Introduction to vector analysis*. Allyn and Bacon, 1979.
9. E. Francois and P. Bouthemy. Derivation of qualitative information in motion analysis. *Image and Vision Computing*, 8(4):279–288, 1990.
10. J.J. Gibson. *The Ecological Approach to Visual Perception*. Houghton Mifflin, 1979.
11. C.G. Harris. Structure from motion under orthographic projection. In O. Faugeras, editor, *Proc. 1st European Conference on Computer Vision*, pages 118–123. Springer-Verlag, 1990.
12. E.C. Hildreth. *The measurement of visual motion*. The MIT press, Cambridge Massachusetts, 1984.
13. K. Kanatani. Detecting the motion of a planar surface by line and surface integrals. *Computer Vision, Graphics and Image Processing*, 29:13–22, 1985.
14. K. Kanatani. Structure and motion from optical flow under orthographic projection. *Computer Vision, Graphics and Image Processing*, 35:181–199, 1986.
15. J.J. Koenderink. Optic flow. *Vision Research*, 26(1):161–179, 1986.
16. J.J. Koenderink and A.J. Van Doorn. Invariant properties of the motion parallax field due to the movement of rigid bodies relative to an observer. *Optica Acta*, 22(9):773–791, 1975.
17. J.J. Koenderink and A.J. Van Doorn. How an ambulant observer can construct a model of the environment from the geometrical structure of the visual inflow. In G. Hauske and E. Butenandt, editors, *Kybernetik*. Oldenburg, Munchen, 1978.
18. J.J. Koenderink and A.J. Van Doorn. Depth and shape from differential perspective in the presence of bending deformations. *J. Opt. Soc. Am.*, 3(2):242–249, 1986.
19. J.J. Koenderink and A.J. van Doorn. Affine structure from motion. *Journal of Optical Society of America*, 1991.
20. D.N. Lee. The optic flow field: the foundation of vision. *Phil. Trans. R. Soc. Lond.*, 290, 1980.
21. H.C. Longuet-Higgins and K. Pradzny. The interpretation of a moving retinal image. *Proc. R. Soc. Lond.*, B208:385–397, 1980.
22. S. J. Maybank. Apparent area of a rigid moving body. *Image and Vision Computing*, 5(2):111–113, 1987.
23. R.C. Nelson and J. Aloimonos. Using flow field divergence for obstacle avoidance: towards qualitative vision. In *Proc. 2nd Int. Conf. on Computer Vision*, pages 188–196, 1988.
24. J.H. Rieger and D.L. Lawton. Processing differential image motion. *J. Optical Soc. of America*, A2(2), 1985.
25. M. Subbarao. Bounds on time-to-collision and rotational component from first-order derivatives of image flow. *Computer Vision, Graphics and Image Processing*, 50:329–341, 1990.
26. D.W. Thompson and J.L. Mundy. Three-dimensional model matching from an unconstrained viewpoint. In *Proceedings of IEEE Conference on Robotics and Automation*, 1987.
27. S. Ullman. *The interpretation of visual motion*. MIT Press, Cambridge, USA, 1979.
28. H. Wang, C. Bowman, M. Brady, and C. Harris. A parallel implementation of a structure from motion algorithm. In *Proc. 2nd European Conference on Computer Vision*, 1992.
29. A.M. Waxman and S. Ullman. Surface structure and three-dimensional motion from image flow kinematics. *Int. Journal of Robotics Research*, 4(3):72–94, 1985.
30. A.M. Waxman and K. Wohn. Contour evolution, neighbourhood deformation and global image flow: planar surfaces in motion. *Int. Journal of Robotics Research*, 4(3):95–108, 1985.



Fig. 1. Using image divergence to estimate time to contact.

Four samples of a video sequence taken from a moving observer approaching a stationary car at a uniform velocity (approximately 1m per time unit). A B-spline snake automatically tracks the area of the rear windscreen (Fig. 2a). The image divergence is used to estimate the time to contact (Fig. 2b). The next image in the sequence corresponds to collision!

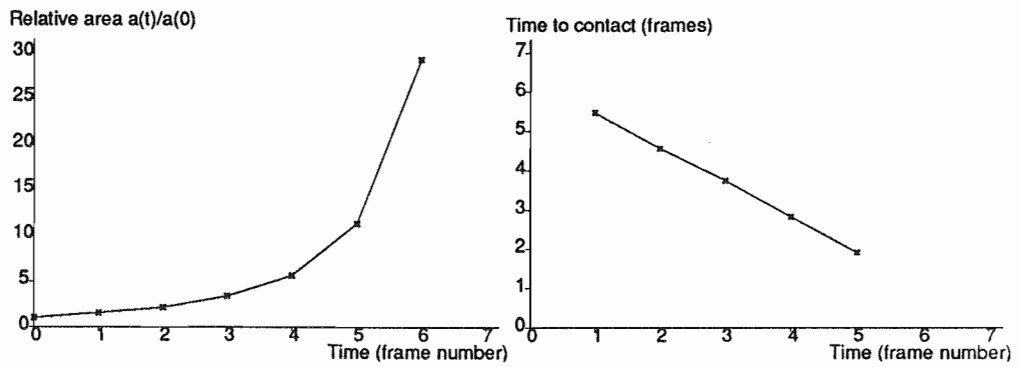


Fig. 2. Apparent area of windscreen for approaching observer and the estimated time to contact.

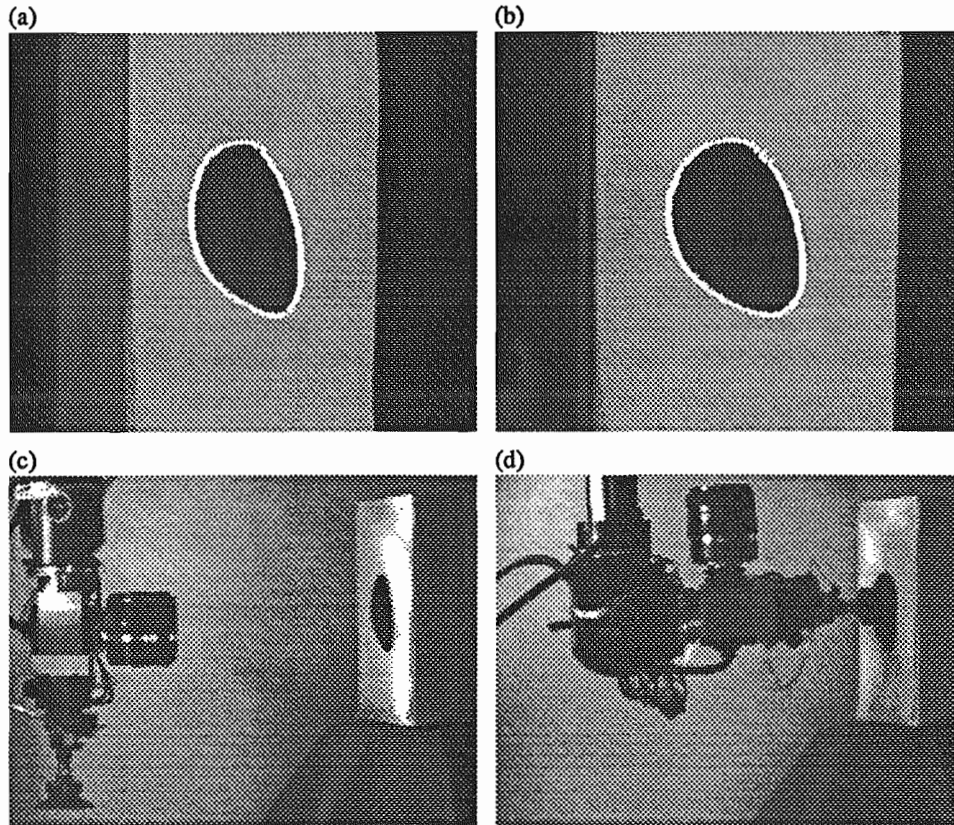


Fig. 3. Visually guided object manipulation using image divergence and deformation.

(a) The image of a planar contour (zero tilt and positive slant, i.e. the direction of increasing depth, F , is horizontal and from left to right). The image contour is localised automatically by a B-spline snake initialised in the centre of the field of view. (b) The effect on apparent shape when the viewer translates to the right while fixating on the target (i.e. A is horizontal, left to right). The apparent shape undergoes an isotropic expansion (positive divergence which increases the area) and a deformation in which the axis of expansion is horizontal. Measurement of the divergence and deformation can be used to estimate the time to contact and surface orientation. This is used to guide the manipulator so that it comes to rest perpendicular to the surface with a pre-determined clearance. Estimates of divergence and deformation made approximately 1m away were sufficient to estimate the target object position and orientation to the nearest 2cm in position and 1° in orientation. This information is used to position a suction gripper in the vicinity of the surface. A contact sensor and small probing motions can then be used to refine the estimate of position and guide the suction gripper before manipulation (d).

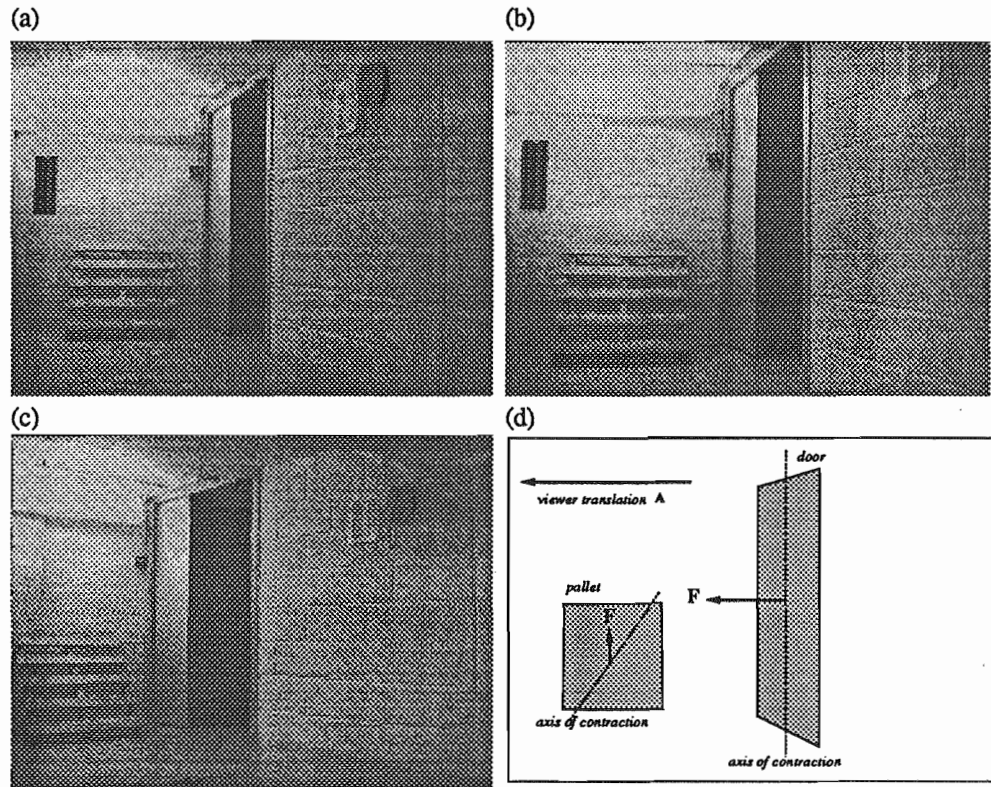


Fig. 4. Qualitative visual navigation using image divergence and deformation.

(a) The image of a door and an object of interest, a pallet. (b) Movement towards the door and pallet produces a deformation in the image seen as an expansion in the apparent area of the door and pallet. This can be used to determine the distance to these objects, expressed as a time to contact – the time needed for the viewer to reach the object if it continued with the same speed. (c) A movement to the left produces combinations of image deformation, divergence and rotation. This is immediately evident from both the door (positive deformation and a shear with a horizontal axis of expansion) and the pallet (clockwise rotation with shear with diagonal axis of expansion). These effects, combined with the knowledge that the movement between the images, are consistent with the door having zero tilt, i.e. horizontal direction of increasing depth, while the pallet has a tilt of approximately 90° , i.e. vertical direction of increasing depth. They are sufficient to determine the orientation of the surface qualitatively (d). This has been done with no knowledge of the intrinsic properties of the camera (camera calibration), its orientations or the translational velocities. Estimation of divergence and deformation can also be recovered by comparison of apparent areas and the orientation of edge segments.