

Efficient Region Tracking With Parametric Models of Geometry and Illumination

Gregory D. Hager, *Member, IEEE*, and Peter N. Belhumeur, *Member, IEEE*,

Abstract—As an object moves through the field of view of a camera, the images of the object may change dramatically. This is not simply due to the translation of the object across the image plane. Rather, complications arise due to the fact that the object undergoes changes in pose relative to the viewing camera, changes in illumination relative to light sources, and may even become partially or fully occluded. In this paper, we develop an efficient, general framework for object tracking—one which addresses each of these complications. We first develop a computationally efficient method for handling the geometric distortions produced by changes in pose. We then combine geometry and illumination into an algorithm that tracks large image regions using no more computation than would be required to track with no accommodation for illumination changes. Finally, we augment these methods with techniques from robust statistics and treat occluded regions on the object as statistical outliers. Throughout, we present experimental results performed on live video sequences demonstrating the effectiveness and efficiency of our methods.

Index Terms—Visual tracking, real-time vision, illumination, motion estimation, robust statistics.



1 INTRODUCTION

VISUAL tracking has emerged as an important component of systems in several application areas including vision-based control [1], [2], [3], [4], human-computer interfaces [5], [6], [7], surveillance [8], [9], agricultural automation [10], [11], medical imaging [12], [13], and visual reconstruction [14], [15], [16]. The central challenge in visual tracking is to determine the image configuration of a target region (or features) of an object as it moves through a camera's field of view. This is done by solving what is known as the temporal correspondence problem: the problem of matching the target region in successive frames of a sequence of images taken at closely-spaced time intervals. The correspondence problem for visual tracking has, of course, much in common with the correspondence problems which arise in stereopsis and motion estimation. It differs, however, in that the goal is not to determine the exact correspondence for every image location in a pair of images, but rather to determine, in a global sense, the movement of an entire target region over a long sequence of images.

What makes tracking difficult is the potential variability in the images of an object over time. This variability arises from three principle sources: variation in target pose or target deformations, variation in illumination, and partial or full occlusion of the target. When ignored, any one of these three sources of variability is enough to cause a tracking algorithm to lose its target. Thus, the two prin-

cipal challenges for visual tracking are to develop accurate models of image variability and to design effective and computationally efficient tracking algorithms which use these models.

In this article, we develop a framework for modeling image variability due to motion and illumination. In the case of motion, all points in the target region are presumed to be part of the same object allowing us the luxury—at least for most applications—of assuming that these points move coherently in space. This permits us to develop low-order parametric models for the image motion of points within a target region—models that can be used to predict the movement of the points and track the target through an image sequence. In the case of illumination, we exploit the observations of [17], [18], [19] to model image variation due to changing illumination by low-dimensional linear subspaces. We then show that these models can be incorporated into an efficient estimation algorithm which establishes temporal correspondence of the target region by simultaneously determining both motion and illumination parameters. Finally, in the case of partial occlusion, we apply results from robust statistics [20] to develop automatic methods of rejecting occluded pixels in a computationally efficient manner. The result is a family of region-tracking algorithms which can easily track large image regions (for example the face of a user at a workstation) at a 30 Hz frame rate using no special hardware other than a standard digitizer.

The tracking algorithms developed in this paper are based on minimizing the sum-of-squared differences (SSD) between two regions. Although this idea has been successfully employed in many contexts including stereo matching [21], optical flow computation [22], and visual motion analysis [23], previous SSD-based tracking algorithms have suffered from a variety of limitations. Many algorithms have modeled the motion of the target region as pure

• G.D. Hager is with the Departments of Computer Science and Electrical Engineering, Yale University, New Haven, CT 06520-8285. E-mail: hager@cs.yale.edu.

• P.N. Belhumeur is with the Departments of Electrical Engineering and Computer Science, Yale University, P.O. Box 208267, New Haven, CT, 06520-8267. E-mail: belhumeur@yale.edu.

Manuscript received 4 Mar. 1997; revised 16 July 1998. Recommended for acceptance by J. Connell.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number 107165.

translation in the image plane [16], [3]. This implicitly assumes that the underlying object is translating parallel to the image plane and is being viewed orthographically. While computationally efficient, over a long sequence these assumptions are often violated [23]. More elaborate tracking algorithms have included parametrized models for articulation [24], [25] or nonrigid deformations [26], [27] as well as linear image subspaces [28], [29]. However, the resulting algorithms rely on nonlinear optimization techniques which require from several seconds to several minutes per frame to compute. Furthermore, none explicitly address the problem of illumination changes. In fact, many algorithms avoid issues related to illumination by estimating and accumulating changes from frame to frame. As a result any error in motion estimation between any two frames is subsequently propagated through the entire sequence.

Another well-established route toward efficient tracking is to detect and track only a sparse collection of features (or contours) [30], [11], [31], [32], [33]. As such methods use local detection of areas of high contrast change, they tend to be insensitive to global changes in the intensity and/or composition of the incident illumination. However, in many situations persistent, strong edges are sparsely distributed throughout the image of the target. This sparseness makes it difficult to establish edge correspondences without strong geometric constraints [33], [31] or an accurate predictive model [11], [30]. In contrast, region-based methods such as those developed in this article make direct and complete use of all available image intensity information, thereby eliminating the need to identify and model a special set of features to track. By incorporating illumination models and robust estimation methods into an efficient correspondence algorithm, the performance of our region tracking algorithms appears to be comparable to that achieved by edge-based methods, thereby making region-based methods an effective complement to local feature-based algorithms.

The remainder of this article is organized as follows. Section 2 establishes a framework for posing the problem of region tracking for parametric motion models and describes conditions under which an efficient tracking algorithm can be developed. Section 3 then shows how models of illumination can be incorporated with no loss of computational efficiency. Section 4 details modifications for handling partial target occlusion via robust estimation techniques. Section 5 presents experimental results from an implementation of the algorithms. Finally, Section 6 presents a short discussion of performance improving extensions to our tracking algorithm.

2 TRACKING MOVING OBJECTS

In this section, we describe a framework for the efficient tracking of a target region through an image sequence. We first write down a general parametric model for the set of allowable image motions and deformations of the target region. We then pose the tracking problem as the problem of finding the best (in a least squares sense) set of parameter values describing the motions and defor-

mations of the target through the sequence. Finally, we describe how the best set of parameters can be efficiently computed.

2.1 On Recovering Structured Motion

Let $I(\mathbf{x}, t)$ denote the brightness value at the location $\mathbf{x} = (x, y)^t$ in an image acquired at time t and let $\nabla_{\mathbf{x}}I(\mathbf{x}, t)$ denote the spatial gradient at that location and time. The symbol t_0 denotes an identified "initial" time and we refer to the image at time t_0 as the *reference image*. Let the set $\mathcal{R} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ be a set of N image locations which define a *target region*. We refer to the brightness values of the target region in the reference image as the *reference template*.

Over time, the relative motion between the target object and the camera causes the image of the target to shift and to deform. Let us model the image motion of the target region of the object by a parametric *motion model* $\mathbf{f}(\mathbf{x}; \boldsymbol{\mu})$ parameterized by $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_n)^t$, with $\mathbf{f}(\mathbf{x}; 0) = \mathbf{x}$ and $N > n$. We assume that \mathbf{f} is differentiable in both $\boldsymbol{\mu}$ and \mathbf{x} . We call $\boldsymbol{\mu}$ the *motion parameter vector*. We consider recovering the motion parameter vector for each image in the tracking sequence as "tracking the object." We write $\boldsymbol{\mu}^*(t)$ to denote the ground truth values of these parameters at time t , and $\boldsymbol{\mu}(t)$ to denote the corresponding estimate. The argument t will be suppressed when it is obvious from its context.

Suppose that a reference template is acquired at time t_0 and that initially $\boldsymbol{\mu}^*(t_0) = \boldsymbol{\mu}(t_0) = 0$. Let us assume for now that the only changes in subsequent images of the target are completely described by \mathbf{f} , i.e., there are no changes in the illumination of the target. It follows that for any time $t > t_0$, there is a parameter vector $\boldsymbol{\mu}^*(t)$ such that

$$I(\mathbf{x}, t_0) = I(\mathbf{f}(\mathbf{x}; \boldsymbol{\mu}^*(t)), t) \text{ for all } \mathbf{x} \in \mathcal{R}. \quad (1)$$

This is a generalization of the so-called *image constancy assumption* [34]. Thus, the motion parameter vector of the target region can be estimated at time t by minimizing the following least squares objective function

$$O(\boldsymbol{\mu}) = \sum_{\mathbf{x} \in \mathcal{R}} \left(I(\mathbf{f}(\mathbf{x}; \boldsymbol{\mu}), t) - I(\mathbf{x}, t_0) \right)^2. \quad (2)$$

For later developments, it is convenient to rewrite this optimization problem in vector notation. To this end, let us consider images of the target region as vectors in an N -dimensional space. The image of the target region at time t , under the change of coordinates with parameters $\boldsymbol{\mu}$, is written as

$$\mathbf{I}(\boldsymbol{\mu}, t) = \begin{bmatrix} I(\mathbf{f}(\mathbf{x}_1, \boldsymbol{\mu}), t) \\ I(\mathbf{f}(\mathbf{x}_2, \boldsymbol{\mu}), t) \\ \vdots \\ I(\mathbf{f}(\mathbf{x}_N, \boldsymbol{\mu}), t) \end{bmatrix}. \quad (3)$$

This vector is subsequently referred to as the *rectified image* at time t with parameters $\boldsymbol{\mu}$. We also make use of the partial derivatives of \mathbf{I} with respect to the components of $\boldsymbol{\mu}$ and the time parameter t . These are written as

$$\mathbf{I}_{\mu_i}(\boldsymbol{\mu}, t) = \frac{\partial \mathbf{I}}{\partial \mu_i} = \begin{bmatrix} I_{\mu_i}(\mathbf{f}(\mathbf{x}_1, \boldsymbol{\mu}), t) \\ I_{\mu_i}(\mathbf{f}(\mathbf{x}_2, \boldsymbol{\mu}), t) \\ \vdots \\ I_{\mu_i}(\mathbf{f}(\mathbf{x}_N, \boldsymbol{\mu}), t) \end{bmatrix} \quad (4)$$

and

$$\mathbf{I}_t(\boldsymbol{\mu}, t) = \frac{\partial \mathbf{I}}{\partial t} = \begin{bmatrix} I_t(\mathbf{f}(\mathbf{x}_1, \boldsymbol{\mu}), t) \\ I_t(\mathbf{f}(\mathbf{x}_2, \boldsymbol{\mu}), t) \\ \vdots \\ I_t(\mathbf{f}(\mathbf{x}_N, \boldsymbol{\mu}), t) \end{bmatrix}, \quad (5)$$

where $1 \leq i \leq n$.

Using this vector notation, the image constancy assumption (1) can be rewritten as

$$\mathbf{I}(\boldsymbol{\mu}^*(t), t) = \mathbf{I}(0, t_0)$$

and (2) becomes

$$O(\boldsymbol{\mu}) = \|\mathbf{I}(\boldsymbol{\mu}, t) - \mathbf{I}(0, t_0)\|^2. \quad (6)$$

In general, (6) is a nonconvex objective function. Thus, in the absence of a good starting point, this problem will usually require some type of costly global optimization procedure to solve [35].

In the case of visual tracking, the continuity of motion provides such a starting point. Suppose that, at some arbitrary time $t > t_0$, the geometry of the target region is described by $\boldsymbol{\mu}(t)$. We recast the tracking problem as one of determining a vector of offsets, $\delta\boldsymbol{\mu}$, such that $\boldsymbol{\mu}(t + \tau) = \boldsymbol{\mu}(t) + \delta\boldsymbol{\mu}$ from an image acquired at $t + \tau$. Incorporating this modification into (6), we redefine the objective function as a function on $\delta\boldsymbol{\mu}$

$$O(\delta\boldsymbol{\mu}) = \|\mathbf{I}(\boldsymbol{\mu}(t) + \delta\boldsymbol{\mu}, t + \tau) - \mathbf{I}(0, t_0)\|^2. \quad (7)$$

If the magnitude of the components of $\delta\boldsymbol{\mu}$ are small, then it is possible to apply continuous optimization procedures to a linearized version of the problem [29], [34], [21], [36], [23]. The linearization is carried out by expanding $\mathbf{I}(\boldsymbol{\mu} + \delta\boldsymbol{\mu}, t + \tau)$ in a Taylor series about $\boldsymbol{\mu}$ and t ,

$$\mathbf{I}(\boldsymbol{\mu} + \delta\boldsymbol{\mu}, t + \tau) = \mathbf{I}(\boldsymbol{\mu}, t) + \mathbf{M}(\boldsymbol{\mu}, t) \delta\boldsymbol{\mu} + \tau \mathbf{I}_t(\boldsymbol{\mu}, t) + h.o.t, \quad (8)$$

where *h.o.t* denotes higher-order terms of the expansion, and \mathbf{M} is the *Jacobian matrix* of \mathbf{I} with respect to $\boldsymbol{\mu}$, i.e., the $N \times n$ matrix of partial derivatives which can be written in column form as

$$\mathbf{M}(\boldsymbol{\mu}, t) = [\mathbf{I}_{\mu_1}(\boldsymbol{\mu}, t) \mid \mathbf{I}_{\mu_2}(\boldsymbol{\mu}, t) \mid \dots \mid \mathbf{I}_{\mu_n}(\boldsymbol{\mu}, t)]. \quad (9)$$

As the expression above indicates, the values of the partial derivatives are a function of the evaluation point $(\boldsymbol{\mu}, t)$. These arguments will be suppressed when obvious from their context.

By substituting (8) into (7) and ignoring the higher-order terms, we have

$$O(\delta\boldsymbol{\mu}) \approx \|\mathbf{I}(\boldsymbol{\mu}, t) + \mathbf{M} \delta\boldsymbol{\mu} + \tau \mathbf{I}_t - \mathbf{I}(0, t_0)\|^2. \quad (10)$$

With the additional approximation

$$\tau \mathbf{I}_t(\boldsymbol{\mu}, t) \approx \mathbf{I}(\boldsymbol{\mu}, t + \tau) - \mathbf{I}(\boldsymbol{\mu}, t),$$

(10) becomes

$$O(\delta\boldsymbol{\mu}) \approx \|\mathbf{M} \delta\boldsymbol{\mu} + \mathbf{I}(\boldsymbol{\mu}, t + \tau) - \mathbf{I}(0, t_0)\|^2. \quad (11)$$

Solving the set of equations $\nabla O = 0$ yields the solution

$$\delta\boldsymbol{\mu} = -(\mathbf{M}^t \mathbf{M})^{-1} \mathbf{M}^t [\mathbf{I}(\boldsymbol{\mu}, t + \tau) - \mathbf{I}(0, t_0)], \quad (12)$$

provided the matrix $\mathbf{M}^t \mathbf{M}$ evaluated at $(\boldsymbol{\mu}, t)$ has full rank. When this is not the case, we are faced with a generalization of the aperture problem, i.e., the target region does not have sufficient structure to determine all of the elements of $\boldsymbol{\mu}$ uniquely. Further discussion of this point can be found in Section 2.4.

In subsequent developments, it will be convenient to define the *error vector*

$$\mathbf{e}(t + \tau) = \mathbf{I}(\boldsymbol{\mu}(t), t + \tau) - \mathbf{I}(0, t_0).$$

Incorporating this definition into (12), we see that the solution of (6) at time $t + \tau$ given a solution at time t is

$$\boldsymbol{\mu}(t + \tau) = \boldsymbol{\mu}(t) - (\mathbf{M}^t \mathbf{M})^{-1} \mathbf{M}^t \mathbf{e}(t + \tau). \quad (13)$$

It is important to note at this point that the solution for $\delta\boldsymbol{\mu}$ is homogeneous in \mathbf{e} . Thus, while errors in calculating \mathbf{M} may affect stability or speed of convergence, they do not affect the stationary points of (13).

2.2 An Efficient Tracking Algorithm

From (13), we see that to track the target region through the image sequence, we must compute the Jacobian matrix $\mathbf{M}(\boldsymbol{\mu}, t)$. Each element of this matrix is given by

$$\begin{aligned} m_{ij} &= I_{\mu_j}(\mathbf{f}(\mathbf{x}_i; \boldsymbol{\mu}), t) \\ &= \nabla_{\mathbf{f}} I(\mathbf{f}(\mathbf{x}_i; \boldsymbol{\mu}), t)^t \mathbf{f}_{\mu_j}(\mathbf{x}_i; \boldsymbol{\mu}) \end{aligned} \quad (14)$$

where $\nabla_{\mathbf{f}} I$ is the gradient of I with respect to the components of the vector \mathbf{f} . Recall that the Jacobian matrix of the transformation \mathbf{f} regarded as a function of $\boldsymbol{\mu}$ is the $2 \times n$ matrix

$$\mathbf{f}_{\mu}(\mathbf{x}; \boldsymbol{\mu}) = \begin{bmatrix} \frac{\partial \mathbf{f}(\mathbf{x}; \boldsymbol{\mu})}{\partial \mu_1} & \frac{\partial \mathbf{f}(\mathbf{x}; \boldsymbol{\mu})}{\partial \mu_2} & \dots & \frac{\partial \mathbf{f}(\mathbf{x}; \boldsymbol{\mu})}{\partial \mu_n} \end{bmatrix}. \quad (15)$$

By making use of (15), \mathbf{M} can be written compactly in row form as

$$\mathbf{M}(\boldsymbol{\mu}, t) = \begin{bmatrix} \nabla_{\mathbf{f}} I(\mathbf{f}(\mathbf{x}_1; \boldsymbol{\mu}), t)^t \mathbf{f}_{\mu}(\mathbf{x}_1; \boldsymbol{\mu}) \\ \nabla_{\mathbf{f}} I(\mathbf{f}(\mathbf{x}_2; \boldsymbol{\mu}), t)^t \mathbf{f}_{\mu}(\mathbf{x}_2; \boldsymbol{\mu}) \\ \vdots \\ \nabla_{\mathbf{f}} I(\mathbf{f}(\mathbf{x}_N; \boldsymbol{\mu}), t)^t \mathbf{f}_{\mu}(\mathbf{x}_N; \boldsymbol{\mu}) \end{bmatrix}. \quad (16)$$

Because \mathbf{M} depends on time-varying quantities, it may appear that it must be completely recomputed at each time step—a computationally expensive procedure involving the calculation of the image gradient vector, the calculation of a $2 \times n$ Jacobian matrix, and $n \cdot 2 \times 1$ vector inner products for each of the N pixels of the target region. However, we now show that it is possible to reduce this computation by both eliminating the need to recompute image gradients and by factoring \mathbf{M} .

First, we eliminate the need to compute image gradients. To do so, let us assume that our estimate is *exact*,

i.e., $\mu(t) = \mu^*(t)$. By differentiating both sides of (1), we obtain

$$\nabla_{\mathbf{x}} I(\mathbf{x}, t_0) = \mathbf{f}_{\mathbf{x}}(\mathbf{x}; \mu)^t \nabla_{\mathbf{x}} I(\mathbf{f}(\mathbf{x}; \mu), t), \quad (17)$$

where $\mathbf{f}_{\mathbf{x}}$ is the 2×2 Jacobian matrix of \mathbf{f} treated as a function of $\mathbf{x} = (x, y)^t$,

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}; \mu) = \left[\frac{\partial \mathbf{f}(\mathbf{x}; \mu)}{\partial x} \mid \frac{\partial \mathbf{f}(\mathbf{x}; \mu)}{\partial y} \right]. \quad (18)$$

Combining (17) with (16), we see that \mathbf{M} can be written as

$$\mathbf{M}(\mu) = \begin{bmatrix} \nabla_{\mathbf{x}} I(\mathbf{x}_1; t_0)^t \mathbf{f}_{\mathbf{x}}(\mathbf{x}_1; \mu)^{-1} \mathbf{f}_{\mu}(\mathbf{x}_1; \mu) \\ \nabla_{\mathbf{x}} I(\mathbf{x}_2; t_0)^t \mathbf{f}_{\mathbf{x}}(\mathbf{x}_2; \mu)^{-1} \mathbf{f}_{\mu}(\mathbf{x}_2; \mu) \\ \vdots \\ \nabla_{\mathbf{x}} I(\mathbf{x}_N; t_0)^t \mathbf{f}_{\mathbf{x}}(\mathbf{x}_N; \mu)^{-1} \mathbf{f}_{\mu}(\mathbf{x}_N; \mu) \end{bmatrix}. \quad (19)$$

It follows that for *any choice* of image deformations, the image spatial gradients need only be calculated once on the reference template. This is not surprising given that the target at time $t > t_0$ is only a geometric distortion of the target at time t_0 , and so its image gradients are also a distortion of those at t_0 . This transformation also allows us to drop the time argument of \mathbf{M} and regard it solely as a function of μ .

The remaining nonconstant factor in \mathbf{M} is a consequence of the fact that, in general, $\mathbf{f}_{\mathbf{x}}$ and \mathbf{f}_{μ} involve components of μ and, hence, implicitly vary with time. However, suppose that we choose \mathbf{f} so that $\mathbf{f}_{\mathbf{x}}^{-1} \mathbf{f}_{\mu}$ can be factored into the product of a $2 \times k$ matrix Γ which depends only on image coordinates, and a $k \times n$ matrix Σ which depends only on μ as

$$\mathbf{f}_{\mathbf{x}}(\mathbf{x}; \mu)^{-1} \mathbf{f}_{\mu}(\mathbf{x}; \mu) = \Gamma(\mathbf{x}) \Sigma(\mu). \quad (20)$$

For example, as discussed in more detail below, one family of such factorizations results when \mathbf{f} is a linear function of the image coordinate vector \mathbf{x} .

Combining (19) with (20), we have

$$\mathbf{M}(\mu) = \begin{bmatrix} \nabla_{\mathbf{x}} I(\mathbf{x}_1; t_0)^t \Gamma(\mathbf{x}_1) \\ \nabla_{\mathbf{x}} I(\mathbf{x}_2; t_0)^t \Gamma(\mathbf{x}_2) \\ \vdots \\ \nabla_{\mathbf{x}} I(\mathbf{x}_N; t_0)^t \Gamma(\mathbf{x}_N) \end{bmatrix} \Sigma(\mu) = \mathbf{M}_0 \Sigma(\mu). \quad (21)$$

As a result, we have shown that \mathbf{M} can be written as a product of an *constant* $N \times k$ matrix \mathbf{M}_0 and a time-varying $k \times n$ matrix Σ .

We can now exploit this factoring to define an efficient tracking algorithm which operates as follows:

Offline:

- Define the target region.
- Acquire and store the reference template.
- Compute and store \mathbf{M}_0 and $\Lambda = \mathbf{M}_0^t \mathbf{M}_0$.

Online:

- Use the most recent motion parameter estimate $\mu(t)$ to rectify the target region in the current image.

- Compute $\mathbf{e}(t + \tau)$ by taking the difference between the rectified image and the reference template.
- Solve the system $\Sigma^t \Lambda \Sigma \delta \mu = \Sigma^t \mathbf{M}_0^t \mathbf{e}(t + \tau)$ for $\delta \mu$, where Σ is evaluated at $\mu(t)$.
- Compute $\mu(t + \tau) = \mu(t) + \delta \mu$.

The online computation performed by this algorithm is quite small and consists of two $n \times k$ matrix multiplies, k N -vector inner products, n k -vector inner products, and an $n \times n$ linear system solution, where k and n are typically far smaller than N .

We note that the computation can be further reduced if Σ is invertible. In this case, the solution to the linear system can be expressed as

$$\delta \mu = -\Sigma^{-t} (\mathbf{M}_0^t \mathbf{M}_0)^{-1} \mathbf{M}_0^t \mathbf{e}(t + \tau), \quad (22)$$

where $\Sigma^{-t} = (\Sigma^{-1})^t$ is evaluated at $\mu(t)$. The factor $(\mathbf{M}_0^t \mathbf{M}_0)^{-1} \mathbf{M}_0^t$ can be computed offline, so the online computation is reduced to n N -vector inner products and n n -vector inner products.

2.3 Some Examples

2.3.1 Linear Models

Let us assume that $\mathbf{f}(\mathbf{x}; \mu)$ is linear in \mathbf{x} . Then we have

$$\mathbf{f}(\mathbf{x}; \mu) = \mathbf{A}(\mu) \mathbf{x} + \mathbf{u}(\mu) \quad (23)$$

and, hence, $\mathbf{f}_{\mathbf{x}} = \mathbf{A}$. It follows that $\mathbf{f}_{\mathbf{x}}^{-1} \mathbf{f}_{\mu}$ is linear in the components of \mathbf{x} and the factoring defined in (20) applies. We now present three examples illustrating these concepts.

Pure Translation. In the case of pure translation, the allowed image motions are parameterized by the vector $\mathbf{u} = (u, v)$ giving

$$\mathbf{f}(\mathbf{x}; \mathbf{u}) = \mathbf{x} + \mathbf{u}. \quad (24)$$

It follows immediately that $\mathbf{f}_{\mathbf{x}}$ and \mathbf{f}_{μ} are both the 2×2 identity matrix and, therefore,

$$\mathbf{M}_0 = [\mathbf{I}_x(t_0) \mid \mathbf{I}_y(t_0)], \quad (25)$$

and Σ is the 2×2 identity matrix.

The resulting linear system is nonsingular if the image gradients in the template region are not all collinear, in which case the solution at each time step is just

$$\delta \mathbf{u} = -(\mathbf{M}_0^t \mathbf{M}_0)^{-1} \mathbf{M}_0^t \mathbf{e}(t + \tau). \quad (26)$$

Note that in this case,

$$\Lambda = -(\mathbf{M}_0^t \mathbf{M}_0)^{-1} \mathbf{M}_0^t,$$

a constant matrix which can be computed offline.

Translation, Rotation, and Scale. The motion of objects which are viewed under scaled orthography and which do not undergo out-of-plane rotation can be modeled in the image plane by a planar rigid motion consisting of a translation \mathbf{u} and a rotation through an angle θ , plus scaling by a factor s . We subsequently refer to this as the RM+S model. The change of coordinates is given by

$$\mathbf{f}(\mathbf{x}; \mathbf{u}, \theta, s) = s \mathbf{R}(\theta) \mathbf{x} + \mathbf{u}, \quad (27)$$

where $\mathbf{R}(\theta)$ is a 2×2 rotation matrix. After some minor algebraic manipulations, we obtain

$$\Gamma(\mathbf{x}) = \begin{bmatrix} 1 & 0 & -y & x \\ 0 & 1 & x & y \end{bmatrix} \quad (28)$$

and

$$\Sigma(\theta, s) = \begin{bmatrix} \frac{1}{s}\mathbf{R}(-\theta) & 0 & 0 \\ \mathbf{0} & 1 & 0 \\ \mathbf{0} & 0 & \frac{1}{s} \end{bmatrix}. \quad (29)$$

From this \mathbf{M}_0 can be computed using (21) and, since Σ is invertible, the solution to the linear system becomes

$$\delta\mu = -\Sigma^t(\mathbf{M}_0^t \mathbf{M}_0)^{-1} \mathbf{M}_0^t \mathbf{e}(t + \tau). \quad (30)$$

This result can be explained as follows. The matrix \mathbf{M}_0 is the linearization of the system about $\theta = 0$ and $s = 1$. At time t , the target has orientation $\theta(t)$ and $s(t)$. Image rectification effectively rotates the target by $-\theta$ and scales by $\frac{1}{s}$, so the displacements of the target are computed *in the original target coordinate system*. Σ^{-t} then applies a change of coordinates to rotate and scale the computed displacements from the original target coordinate system back to the actual target coordinates.

Affine Motion. The image distortions of planar objects viewed under orthographic projection are described by a six-parameter linear change of coordinates. Suppose that we define

$$\begin{aligned} \mu &= (u, v, a, b, c, d)^t \\ \mathbf{f}(\mathbf{x}; \mu) &= \begin{bmatrix} a & c \\ b & d \end{bmatrix} \mathbf{x} + \begin{bmatrix} u \\ v \end{bmatrix} = \mathbf{A}\mathbf{x} + \mathbf{u} \end{aligned} \quad (31)$$

After some minor algebraic manipulations, we obtain

$$\Gamma(\mathbf{x}) = \begin{bmatrix} 1 & 0 & x & 0 & y & 0 \\ 0 & 1 & 0 & x & 0 & y \end{bmatrix} \quad (32)$$

and

$$\Sigma(\mu) = \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{A}^{-1} \end{bmatrix}. \quad (33)$$

Note that Σ is once again invertible which allows for additional computational savings as before.

2.3.2 Nonlinear Motion Models

The separability property needed for factoring does not hold for any type of nonlinear motion. However, consider a motion model of the form

$$\mathbf{f}(\mathbf{x}; u, v, a) = \mathbf{x} + \begin{bmatrix} u \\ v + 1 / 2ax^2 \end{bmatrix}, \quad (34)$$

where $\mathbf{x} = (x, y)^t$. Intuitively, this model performs a quadratic distortion of the image according to the equation $y = 1/2ax^2$. For example, a polynomial model of this form was used in [27] to model the motions of lips and eyebrows on a face. Again, after several algebraic steps we arrive at

$$\Gamma(\mathbf{x}) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & x & \frac{x^2}{2} \end{bmatrix} \text{ and } \Sigma(\mu) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -a & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (35)$$

Note that this general result holds for any distortion which can be expressed exclusively as either $y = f(x)$ or $x = g(y)$. However, adding more freedom to the motion model, for example combining affine and polynomial distortion, often makes factoring impossible. One possibility in such cases is to use a cascaded model in which the image is first rectified using an affine distortion model, and then the resulting rectified image is further rectified for polynomial distortion.

2.4 On the Structure of Image Change

The Jacobian matrix \mathbf{M} plays a central role in the algorithms described above, so it is informative to digress briefly on its structure. If we consider the rectified image as a continuous time-varying quantity, then its total derivative with respect to time is

$$\frac{d\mathbf{I}}{dt} = \mathbf{M} \frac{d\mu}{dt} + \frac{\partial \mathbf{I}}{\partial t} \text{ or } \dot{\mathbf{I}} = \mathbf{M}\dot{\mu} + \mathbf{I}_t. \quad (36)$$

Note that this is simply a differential form of (8). Due to the image constancy assumption (1), it follows that $\dot{\mathbf{I}} = 0$ when $\mu = \mu^*$. This is, of course, a parameterized version of Horn's optical flow constraint equation [34].

In this form, it is clear that the role of \mathbf{M} is to relate variations in motion parameters to variations in brightness values in the target region. The solution given in (13) effectively reverses this relationship and provides a method for interpreting observed *changes* in brightness as motion. In this sense, we can think of the algorithm as performing correlation on temporal changes (as opposed to spatial structure) to compute motion.

To better understand the structure of \mathbf{M} , recall that in column form, it can be written in terms of the partial derivatives of the rectified image:

$$\mathbf{M} = [\mathbf{I}_{\mu 1} \mid \mathbf{I}_{\mu 2} \mid \dots \mid \mathbf{I}_{\mu n}]. \quad (37)$$

Thus, the model states that the temporal variation in image brightness in the target region is a weighted combination of the vectors $\mathbf{I}_{\mu i}$. We can think of each of these columns (which have an entry for every pixel in the target region) as a "motion template" which directly represents the changes in brightness induced by the motion represented by the corresponding motion parameter. For example, in Fig. 1 we have shown these templates for four canonical motions of an image of a human face.

The development in this section has assumed that we start with a given parametric motion model from which these templates are derived. Based on that model, the structure of each entry of \mathbf{M} is given by (15) which states that

$$m_{i,j} = \nabla_{\mathbf{f}} I \cdot \mathbf{f}_{\mu_j} \big|_{\mathbf{x}=\mathbf{x}_i}. \quad (38)$$

The image gradient $\nabla_{\mathbf{f}} I$ defines, at each point in the image, the direction of strongest intensity change. The vector \mathbf{f}_{μ_j} evaluated at \mathbf{x}_i is the instantaneous direction and magnitude of motion of that image location captured by the parameter

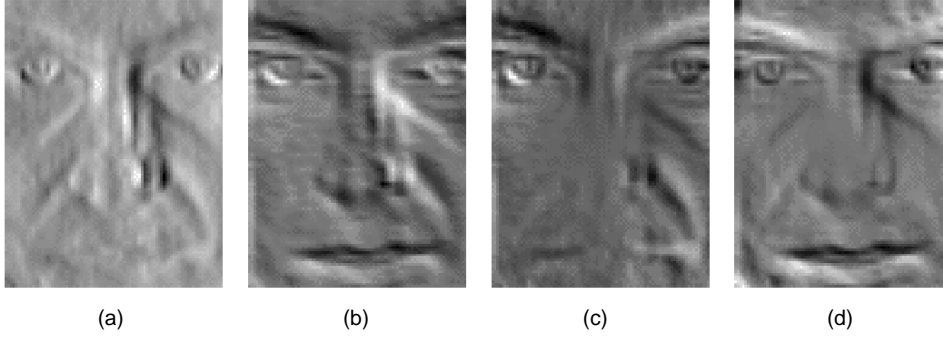


Fig. 1. The motion templates of a human face for four canonical motions. (a) X translation. (b) Y translation. (c) Rotation. (d) Scale.

μ_j . The collection of the latter for all pixels in the region represents the *motion field* defined by the motion parameter μ_j . Thus, the change in the brightness of the image location \mathbf{x}_i due to the motion parameter μ_j is the projection of the image gradient onto the motion vector. This also explains why each pixel in the image contributes only one constraint to the parameter computation.

More importantly, the methods described above assume that $\mathbf{M}^t \mathbf{M}$ is full rank. Although, in general, this condition depends on both the structure of the motion to be computed and the structure of the image itself, the form of (38) provides some insight into the rank structure of \mathbf{M} . In particular, it follows that for $\mathbf{M}^t \mathbf{M}$ to be rank deficient, there must exist a $\gamma \in \mathbb{R}^n$ such that

$$(\nabla_{\mathbf{f}} I_{\mu}^t \big|_{\mathbf{x}=\mathbf{x}_i}) \gamma = 0, \quad 1 \leq i \leq N. \quad (39)$$

Geometrically, this condition corresponds to a motion γ such that the displacement of every pixel in the image is orthogonal to the local image gradient.¹ Thus, we can view the rank deficiency of \mathbf{M} as a generalization of the well-known aperture problem [34] in optical flow.

Finally, (38) suggests how our techniques can be used to perform structured motion estimation without an explicit parametric motion model. First, if the changes in images due to motion can be observed directly (for example, by computing the differences of images taken before and after small reference motions are performed), then these can be used as the motion templates which comprise \mathbf{M} . Second, if a one or more *motion fields* can be observed (for example, by tracking a set of fiducial points in a series of training images), then projecting each element of the motion field onto the corresponding image gradient yields motion templates for those motion fields. The linear estimation process described above can be used to interpret time-varying images in terms of those basis motions.

3 ILLUMINATION-INSENSITIVE TRACKING

The systems described above are inherently sensitive to changes in illumination of the target region. This is not surprising, as the incremental estimation step is effectively computing a structured optical flow, and optical flow

methods are well-known to be sensitive to illumination changes [34]. Thus, shadowing or shading changes of the target object over time lead to bias, or, in the worst case, complete loss of the target.

Recently, it has been shown that a relatively small number of “basis” images can often be used to account for large changes in illumination [19], [18], [17], [37]. Briefly, the reason for this is as follows. Consider a point p on a Lambertian surface and a collimated light source characterized by a vector $\mathbf{s} \in \mathbb{R}^3$, such that the direction of \mathbf{s} gives the direction of the light rays and $\|\mathbf{s}\|$ gives the intensity of the light source. The irradiance at the point p is given by

$$E = \mathbf{a} \mathbf{n} \cdot \mathbf{s}, \quad (40)$$

where \mathbf{n} is the unit inward normal vector to the surface at p and a is the nonnegative absorption coefficient (albedo) of the surface at the point p [34]. This shows that the irradiance at the point p , and hence the gray level measured by a camera, is linear on $\mathbf{s} \in \mathbb{R}^3$.

Therefore, in the absence of self-shadowing, given three images of a Lambertian surface from the same viewpoint taken under three linearly independent light source directions, one can reconstruct the image of the surface under a novel lighting direction by a linear combination of the three original images [37], [38]. In other words, if the surface is purely Lambertian and there is no shadowing, then all images under varying illumination lie within a 3D linear subspace of \mathbb{R}^N , the space of all possible images (where N is the number of pixels in the images).

A complication comes when handling shadowing: All images are no longer guaranteed to lie in a linear subspace [19]. Nevertheless, as done in [17], we can still use a linear model as an approximation: A small set of basis images can account for much of the shading changes that occur on patches of nonspecular surfaces. Naturally, we need more than three images (we use between eight and 15) and a higher than three-dimensional linear subspace (we use four or five) if we hope to provide good approximation to these effects.

Returning to the problem of region tracking, suppose now that we have a basis of image vectors $\mathbf{B}_1, \mathbf{B}_2, \dots, \mathbf{B}_m$ where the i th element of each of the basis vectors corresponds to the image location $\mathbf{x}_i \in \mathcal{R}$. To accommodate changes in contrast, we choose the first basis vector to be the template image itself, i.e., $\mathbf{B}_1 = \mathbf{I}(0, t_0)$. To model bright-

1. Note that one possibility is that the gradient at a point is zero, in which case this is true of any motion.

ness changes, we choose the second basis vector to be a column of ones, i.e., $\mathbf{B}_2 = (1, 1, \dots, 1)^t$.² Let us choose the remaining basis vectors by performing SVD (singular value decomposition) on a set of training images of the target, taken under varying illumination. We denote the collection of basis vectors by the matrix $\mathbf{B} = [\mathbf{B}_1 | \mathbf{B}_2 | \dots | \mathbf{B}_m]$ and the corresponding parameters by the vector $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)^t$.

Combining motion with illumination, the image constancy constraint, (1), can now be rewritten as

$$\mathbf{I}(\mu^*(t), t) = \mathbf{I}(\mathbf{0}, t_0) + \mathbf{B}\lambda(t), \quad (41)$$

and (2) becomes

$$O(\mu, \lambda) = \|\mathbf{I}(\mu, t) - \mathbf{I}(\mathbf{0}, t_0) - \mathbf{B}\lambda\|^2. \quad (42)$$

In short, we now have expressions which simultaneously model both geometric and photometric image changes. By rewriting this optimization as

$$O(\delta\mu, \lambda) = \|\mathbf{I}(\mu(t) + \delta\mu, t + \tau) + \mathbf{B}\lambda - \mathbf{I}(\mathbf{0}, t_0)\|^2, \quad (43)$$

and substituting in (8) we arrive at

$$O(\delta\mu, \lambda) = \|\mathbf{M}\delta\mu + \mathbf{B}\lambda + \mathbf{I}(\mu(t), t + \tau) - \mathbf{I}(\mathbf{0}, t_0)\|^2. \quad (44)$$

Solving $\nabla O(\delta\mu, \lambda) = \mathbf{0}$ yields

$$\begin{bmatrix} \delta\mu \\ \lambda \end{bmatrix} = -\begin{bmatrix} \mathbf{M}'\mathbf{M} & \mathbf{M}'\mathbf{B} \\ \mathbf{B}'\mathbf{M} & \mathbf{B}'\mathbf{B} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{M}' \\ \mathbf{B}' \end{bmatrix} \mathbf{e}(t + \tau), \quad (45)$$

where $\mathbf{e}(t + \tau) = \mathbf{I}(\mu(t), t + \tau) - \mathbf{I}(\mathbf{0}, t_0)$ as before.

We would now like to apply the factoring methods of the previous section to reduce the online computation needed for estimation. However, letting $\mathbf{B}(\mathbf{x}; \lambda)$ denote the value for pixel location \mathbf{x} of $\mathbf{B}\lambda$, from (41) we have

$$\nabla_{\mathbf{x}} \mathbf{I}(\mathbf{x}, t_0) = \mathbf{f}_{\mathbf{x}}(\mathbf{x}; \mu)^t \nabla_{\mathbf{x}} \mathbf{I}(\mathbf{f}(\mathbf{x}; \mu), t) - \nabla_{\mathbf{x}} \mathbf{B}(\mathbf{x}; \lambda). \quad (46)$$

If we follow the same steps as before in factoring \mathbf{M} , we find that $\nabla_{\mathbf{x}} \mathbf{B}(\mathbf{x}; \lambda)$ will appear in \mathbf{M}_0 , thus requiring recomputation of that form. In practice, we have found that, for the specific case of illumination, these terms are small and can be safely ignored without seriously affecting the stability of the resulting tracking system.³ Ignoring these terms, \mathbf{M} factors as before and, since \mathbf{B} is constant, the system can be efficiently computed.

Further efficiencies can be realized if we are only interested in the motion parameters and hence we only need to compute the portions of (45) pertaining to those parameters. We can compute an explicit form of this expression by first optimizing over λ as a function of $\delta\mu$ in (44) and substituting the solution back into (44). Doing so, solving the resulting expression for $\delta\mu$, and writing \mathbf{M} in factored form, we arrive at

$$\delta\mu = -\Sigma^{-t} (\mathbf{M}_0^t \mathbf{N} \mathbf{M}_0)^{-1} \mathbf{M}_0^t \mathbf{N} \mathbf{e}(t + \tau), \quad (47)$$

$$\mathbf{N} = (\mathbf{I} - \mathbf{B}(\mathbf{B}'\mathbf{B})^{-1}\mathbf{B}'). \quad (48)$$

2. In practice, choosing a value close to the mean of the brightness of the image produces a better conditioned linear system.

3. Note that this may not hold true for other subspace decompositions such as those used by [29].

Since \mathbf{N} is constant, the computation needed to realize (47) depends only on the number of motion fields to be computed, not on the illumination model. As a result, we can compute motion parameters while accounting for variations in illumination *using no more online computation than would be required to compute pure motion*.

4 MAKING TRACKING RESISTANT TO OCCLUSION

As a system tracks objects over a large space, it is not uncommon that other objects “intrude” into the picture. For example, the system may be in the process of tracking a target region which is the side of a building when, due to observer motion, a parked car begins to occlude a portion of that region. Similarly the target object may rotate, causing the tracked region to “slide off” and pick up a portion of the background. Such intrusions will bias the motion parameter estimates and, in the long term can potentially cause mistracking. In this section, we describe how to avoid such problems. For the sake of simplicity, we develop a solution for the case where we are only recovering motion parameters; the modifications for combined motion and illumination models are straightforward.

A common approach to this problem is to assume that occlusions create large image differences which can be viewed as “outliers” by the estimation process [29]. The error metric is then modified to reduce sensitivity to outliers by solving a robust optimization problem of the form

$$O_R(\mu) = \sum_{\mathbf{x} \in \mathcal{R}} \rho(I(\mathbf{f}(\mathbf{x}; \mu), t) - I(\mathbf{x}, t_0)), \quad (49)$$

where ρ is one of a variety of “robust” regression metrics [39].

It is well-known that optimization of (49) is closely related to another approach to robust estimation—iteratively reweighted least squares (IRLS). We have chosen to implement the optimization using a somewhat unusual form of IRLS due to Dutter and Huber [20]. In order to formulate the algorithm, we introduce the notation of an “inner iteration” which is performed one or more times at each time step. We will use a superscript to denote these iterations, and refer each time step in the estimation as an “outer iteration.”

Let $\delta\mu^i$ denote the value of $\delta\mu$ computed by the i th inner iteration with $\delta\mu^0 = \mathbf{0}$. Define the vector of residuals in the i th iteration \mathbf{r}^i as

$$\mathbf{r}^i = \mathbf{e}(t + \tau) - \mathbf{M}(\mu)\delta\mu^i. \quad (50)$$

We introduce a diagonal weighting matrix $\mathbf{W}^i = \mathbf{W}(\mathbf{r}^i)$ which has entries

$$w_{k,k}^i = \eta(r_k^i) = \rho'(r_k^i) / r_k^i, \quad 1 \leq k \leq N. \quad (51)$$

The inner iteration cycle at time $t + \tau$ consists of performing an estimation step by solving the linear system

$$\Sigma^t \Lambda \Sigma \delta\mu^{i+1} = \Sigma^t \mathbf{M}_0^t \mathbf{W}^i \mathbf{r}^i, \quad (52)$$

where Σ is evaluated at $\mu(t)$ and \mathbf{r}^i and \mathbf{W}^i are given by (50) and (51), respectively. This process is repeated for k iterations.

This form of IRLS is particularly efficient for our problem. It does not require recomputation of Λ or Σ and, since the weighting matrix is diagonal, does not add significantly

to the overall computation time needed to solve the linear system. In addition, the error vector \mathbf{e} is fixed over all inner iterations, so these iterations do not involve acquiring or warping images.

As discussed in [20], on linear problems this procedure is guaranteed to converge to a unique global minimum for a large variety of choices of ρ . In this article, ρ is taken to be a so-called “windsorizing” function [39] which is of the form

$$\rho(r) = \begin{cases} r^2 / 2 & \text{if } |r| \leq \tau \\ c|r| - c^2 / 2 & \text{if } |r| > \tau \end{cases} \quad (53)$$

where r is normalized to have unit variance. The parameter τ is a user-defined threshold which places a limit on the variations of the residuals before they are considered outliers. This function has the advantage of guaranteeing global convergence of the IRLS method while being cheap to compute. The updating function for matrix entries is

$$\eta(r) = \begin{cases} 1 & \text{if } |r| \leq \tau \\ c / |r| & \text{if } |r| > \tau \end{cases} \quad (54)$$

As stated, the weighting matrix is computed anew at each outer iteration, a process which can require several inner iterations. However, given that tracking is a continuous process, it is natural to start each outer iteration with a weighting matrix which is closely related to that computed at the end of the previous outer iteration. In doing so, two issues arise. First, the fact that the linear system we are solving is a local linearization of a nonlinear system means that, in cases when interframe motion is large, the effect of higher-order terms of the Taylor series expansion will cause areas of the image to masquerade as outliers. Second, if we assume that areas of the image with low weights correspond to intruders, it makes sense to add a “buffer zone” around those areas before the next outer iteration to proactively cancel the effects of intruder motion.

Both of these problems can be dealt with by noting that the diagonal elements of \mathbf{W} themselves form an image where “dark areas” (those locations with low value) are areas of occlusion or intrusion, while “bright areas” (those with value one) are the expected target. Let $Q(\mathbf{x})$ to be the pixel values in the eight-neighborhood of the image coordinate \mathbf{x} plus the value at \mathbf{x} itself. We use two common morphological operators [40]

$$\text{erode}(\mathbf{x}) = \max_{v \in Q(\mathbf{x})} v \quad (55)$$

$$\text{dilate}(\mathbf{x}) = \min_{v \in Q(\mathbf{x})} v. \quad (56)$$

When applied to a weighting matrix image, *erode* has the effect of removing small areas of outlier pixels, while *dilate* increases their size. Between frames of the sequence we propagate the weighting matrix forward after applying one step of *erode* to remove small areas of outliers followed by two or three steps of *dilate* to provide a “buffer” about previously detected intruders.

5 IMPLEMENTATION AND EXPERIMENTS

This section illustrates the performance of the tracking algorithm under a variety of circumstances, noting particularly

the effects of image warping, illumination compensation, and outlier detection. All experiments were performed on live video sequences by an SGI Indy equipped with a 175 Mhz R4400 SC processor and VINO image acquisition system.

5.1 Implementation

We have implemented the methods described above within the X Vision environment [41]. The implemented system incorporates all of the linear motion models described in Section 2, nonorthonormal illumination bases as described in Section 3, and outlier rejection using the algorithm described in Section 4.

The image warping required to support the algorithm is implemented by factoring linear transformations into a rotation matrix and a positive-definite upper-diagonal matrix. This factoring allows image warping to be implemented by first acquiring a rotated rectangular image region surrounding the target, and then scaling and shearing the region using bilinear interpolation. The resolution of the region is then reduced by averaging neighboring pixels. Spatial and temporal derivatives are computed by applying Prewitt operators on the reduced scale images. More details on this level of the implementation can be found in [41].

In these experiments, the algorithm is initialized by interactively selecting a region to track in a live video stream. The algorithm immediately acquires the selected region as the reference template and performs tracking on all subsequent images of the stream. When an illumination basis is used, care was taken to select the reference template to correspond to the basis, but no automatic registration was performed.

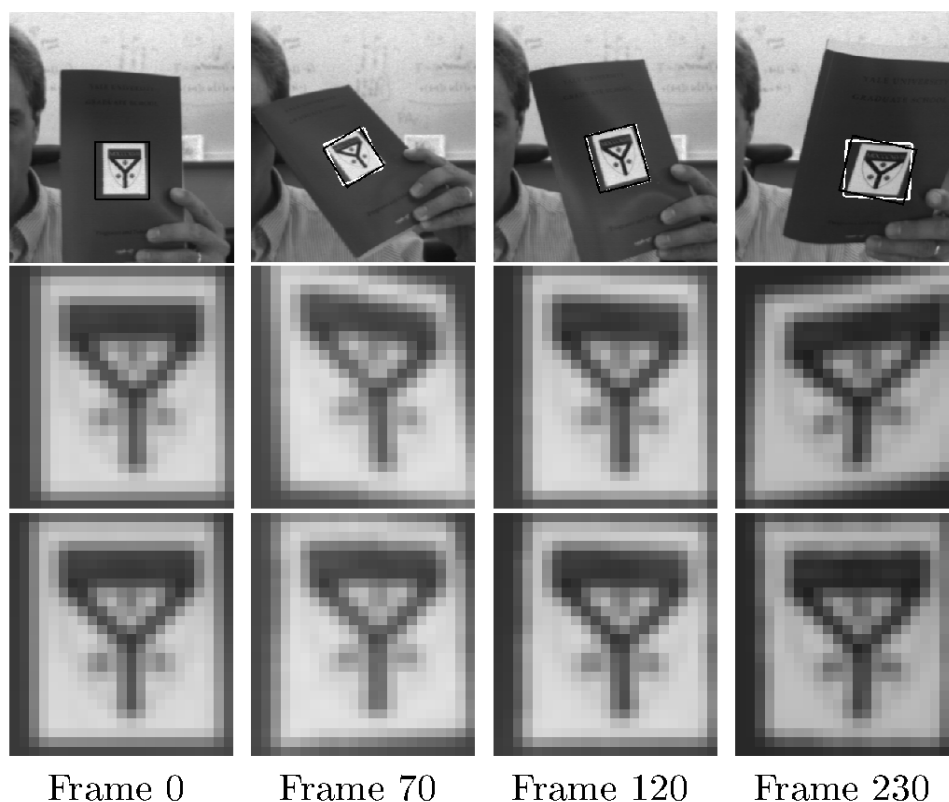
Timings of the algorithm⁴ indicate that it can perform frame rate (30 Hz) tracking of image regions of up to 100×100 pixels undergoing affine distortions and illumination changes at one-half resolution. Similar performance has been achieved on a 120 Mhz Pentium processor and 70 Mhz Sun Sparc-Station. Higher performance is achieved for smaller regions, lower resolutions, or fewer parameters. For example, tracking the same size region while computing just translation at one-fourth resolution takes just four milliseconds per cycle.

5.2 Planar Tracking

As a baseline, we first consider tracking a non-specular planar object—the cover of a book. Affine warping augmented with brightness and contrast compensation is a good approximation in this case (it is exact for an orthographic camera model and purely Lambertian surface). As a point of comparison, recent work by Black and Jepson [29] used the rigid motion plus scaling model for SSD-based region tracking. Their reduced model is more efficient and may be more stable since fewer parameters must be computed, but it does ignore the effects of changing aspect ratio and shear.

We tested both the rigid motion plus scale (RM+S) and full affine (FA) motion models on the same live video sequence of the book cover in motion. Fig. 2 shows the set of

4. Because of additional data collection overhead, the tracking performance in the experiments presented here is slower than the stated figures.



Residuals: Planar Test

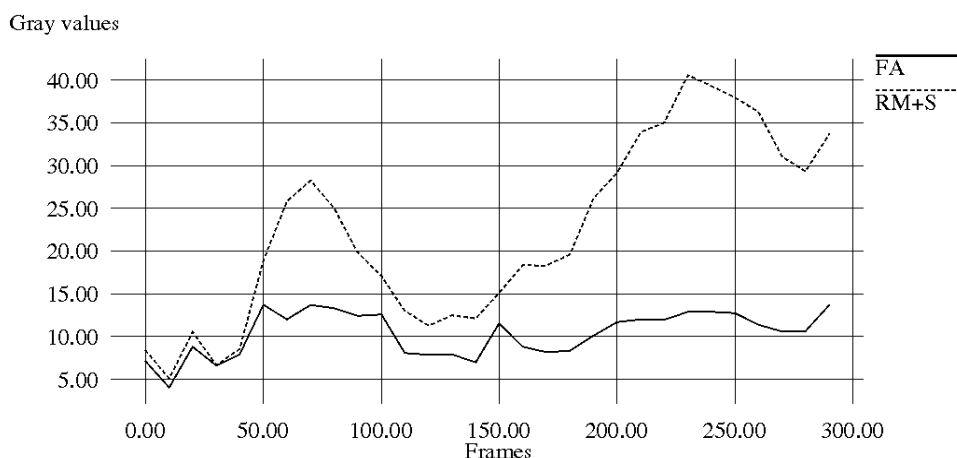


Fig. 2. Top, several images of a planar region and the corresponding warped image computed by a tracker computing position, orientation, and scale (RM+S), and one computing a full affine deformation (FA). The image at the left is the initial reference image. Bottom, the graph of the SSD residuals for both algorithms.

motion templates (the columns of the motion matrix) for an 81×72 region of a book cover tracked at one third resolution. The upper series of images shows several images of the object with the region tracked indicated with a black frame (the RM+S algorithm) and a white frame (the FA algorithm).⁵ The middle row of images shows the output of the warping operator from the RM+S algorithm. If the computed parameters were error-free, these images would be identical. However, because of the inability to correct for aspect ratio and skew, the best fit leads to a skewed image.

5. These annotations indicate the region acquired in the first stage of image warping and so do not indicate distortions due to image shear.

The bottom row shows the output of the warping operator for the FA algorithm. Here, we see that full affine warping is much better at accommodating the full range of image distortions. The graph at the bottom of the figure shows the least squares residual (in squared gray-values per pixel). Here, the difference between the two geometric models is clearly evident.

5.3 Human Face Tracking

There has been a great deal of recent interest in face tracking in the computer vision literature [27], [6], [42]. Although faces can produce images with significant variation due to

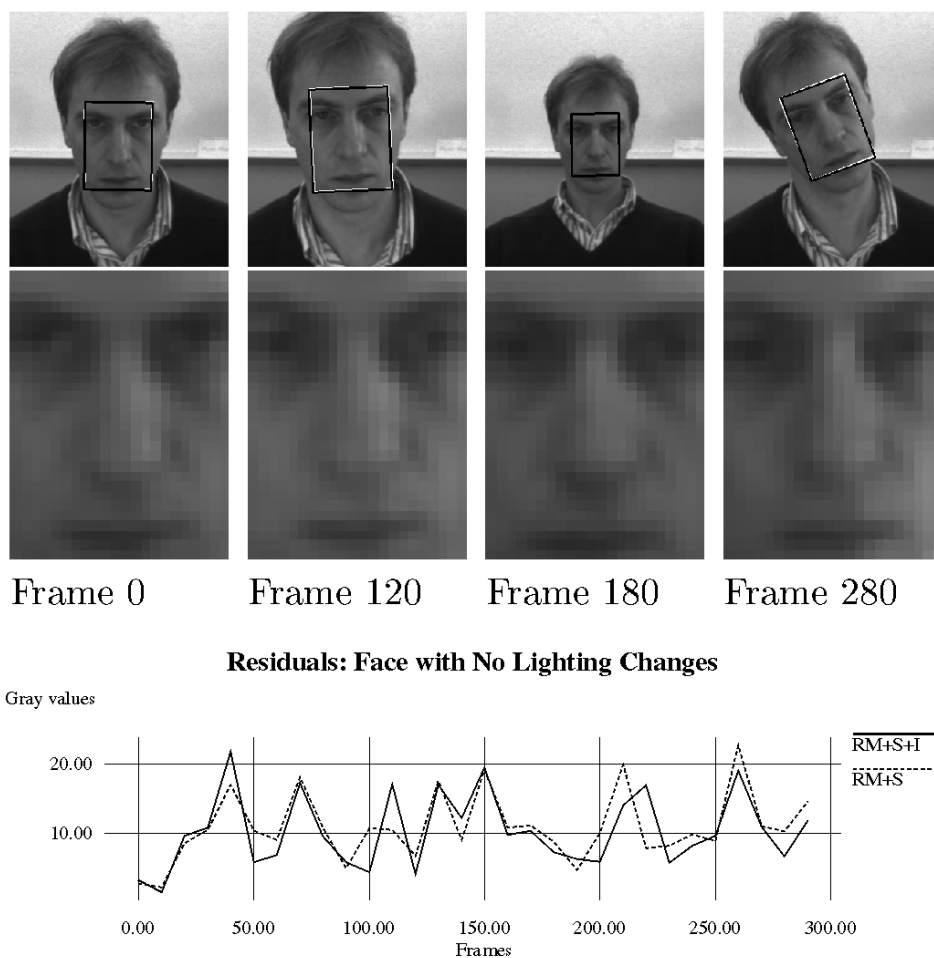


Fig. 3. Top row, excerpts from a sequence of tracked images of a face. The black frames represent the region tracked by an SSD algorithm using no illumination model (RM+S) and the white frames represent the regions tracked by an algorithm which includes an illumination model (RM+S+I). In some cases the estimates are so close that only one box is visible. Middle row, the region within the frame warped by the current motion estimate. Bottom, the residuals of the algorithms expressed in gray-scale units per pixel as a function of time.

illumination, empirical results suggest that a small number of basis images of a face gathered under different illuminations is sufficient to accurately account for most gross shading and illumination effects [17]. At the same time, the depth variations exhibited by facial features are small enough to be well-approximated by an affine warping model. The following experiments demonstrate the ability of our algorithm to track a face as it undergoes changes in pose and illumination, and under partial occlusion. Throughout, we assume the subject is roughly looking toward the camera, so we use the rigid motion plus scaling (RM+S) motion model. Fig. 1 shows the columns of the motion matrix for this model.

5.3.1 Geometry

We first performed a test to determine the accuracy of the computed motion parameters for the face and to investigate the effect of the illumination basis on the sensitivity of those estimates. During this test, we simultaneously executed two tracking algorithms: one using the rigid motion plus scale model (RM+S) and one which additionally included an illumination model for the face (RM+S+I). The algorithms were executed on a sequence which did not contain large

changes in the illumination of the target. The top row of Fig. 3 shows images excerpted from the video sequence. In each image, the black frames denote the region selected as the best match by RM+S and the white frames correspond to the best match computed by RM+S+I. For this test, we would expect both algorithms to be quite accurate and to exhibit similar performance unless the illumination basis significantly affected the sensitivity of the computation. As is apparent from the figures, the computed motion parameters of both algorithms are extremely similar for the entire run—so close that in many cases one frame is obscured by the other.

In order to demonstrate the absolute accuracy of the tracking solution, below each live image in Fig. 3 we have included the corresponding rectified image computed by RM+S+I. The rectified image at time zero is the reference template. If the motion of the target fit the RM+S motion model, and the computed parameters were exact, then we would expect each subsequent rectified image to be identical to the reference template. Despite the fact that the face is nonplanar and we are using a reduced motion model, we see that the algorithm is quite effective at computing an accurate geometric match.

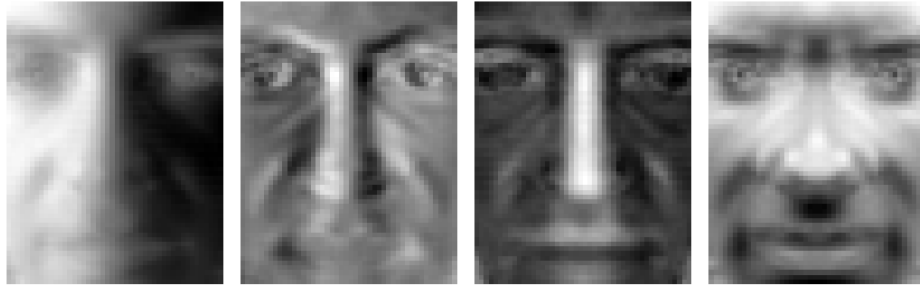


Fig. 4. The illumination basis for the face (contrast and brightness components not shown).

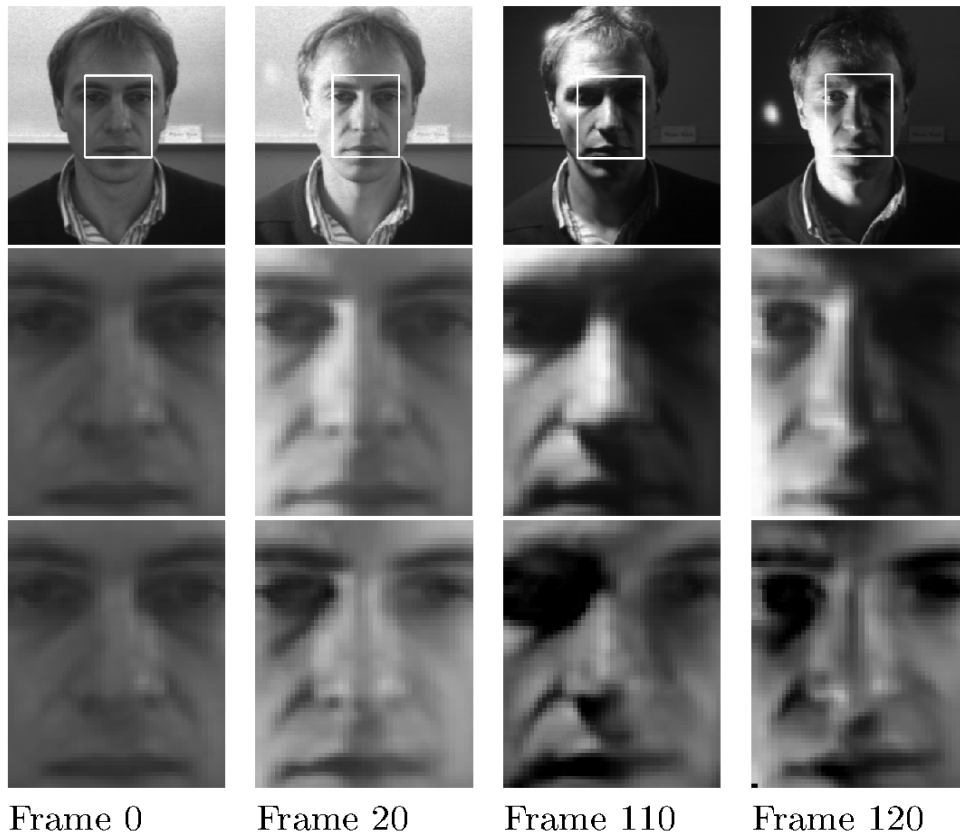


Fig. 5. The first row of images shows excerpts of a tracking sequence. The second row is a magnified view of the region in the white frame. The third row contains the images in the second row after adjustment for illumination using the illumination basis shown in Fig. 4 (for the sake of comparison, we have not adjusted for brightness and contrast across the sequence).

Finally, the graph in Fig. 3 shows the residuals of the linearized SSD computation at each time step. As is apparent from the figures, the residuals of both algorithms are also extremely similar for the entire run. From this experiment we conclude that, in the absence of illumination changes, the performance of both algorithms is quite similar—including illumination models does not appear to reduce accuracy.

5.3.2 Illumination

In a second set of experiments, we kept the face nearly motionless and varied the illumination. We used an illumination basis of four orthogonal image vectors. This basis was computed offline by acquiring ten images of the face under various lighting conditions. A singular value decomposi-

tion was applied to the resulting image vectors and the vectors with the maximum singular values were chosen to be included in the basis. The illumination basis is shown in Fig. 4.

Fig. 5 shows the effects of illumination compensation for the illumination situations depicted in the first row. As with warping, if the compensation were perfect, the images of the bottom row would appear to be identical up to brightness and contrast. In particular, note how the strong shading effects of frames 110 and 120 have been “corrected” by the illumination basis.

5.3.3 Combining Illumination and Geometry

Next, we present a set of experiments illustrating the interaction of geometry and illumination. In these experiments,

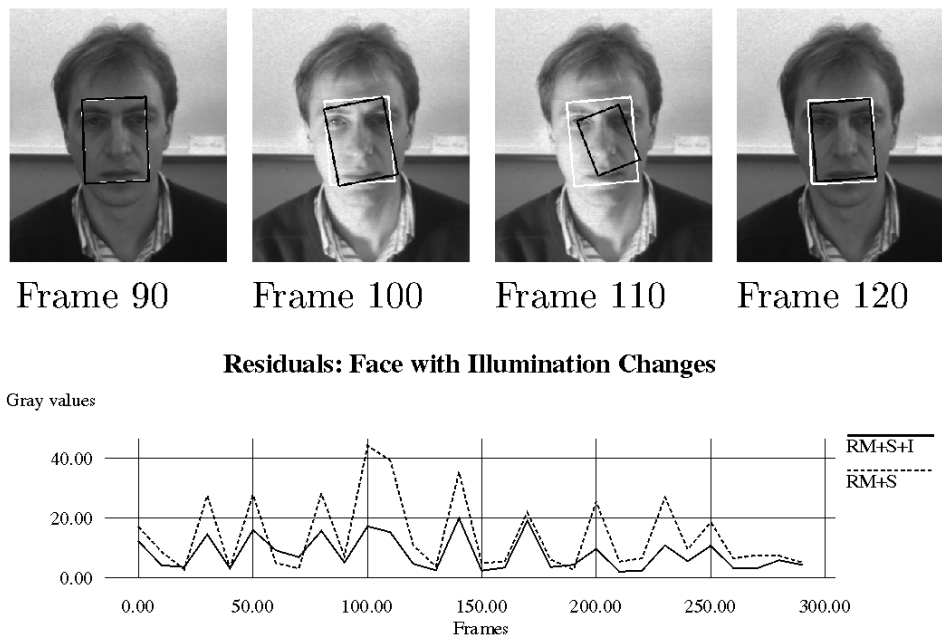


Fig. 6. Top, an excerpt from a tracking sequence containing changes in both geometry and illumination. The black frame corresponds to the algorithm without illumination (RM+S) and the white frame corresponds to the algorithm with an illumination basis (RM+S+I). Note that the algorithm which does not use illumination completely loses the target until the original lighting is restored. Bottom, the residuals, in gray-scale units per pixel, of the two algorithms as a light is turned on and off.



Fig. 7. A run combining illumination and geometry in which the algorithm without illumination compensation (black frame) loses the target while the algorithm with illumination compensation (white frame) does not.

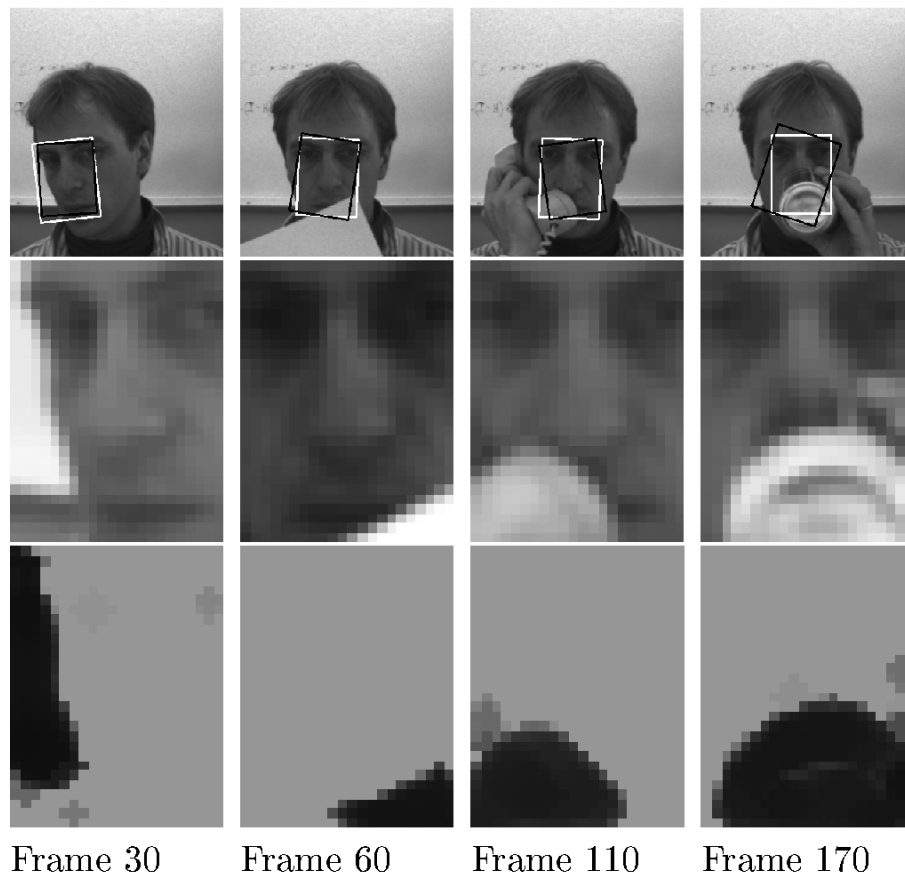
we again executed two algorithms labeled RM+S and RM+S+I. As the algorithms were operating, a light was periodically switched on and off and the face moved slightly. The results appear in Fig. 6. In the residual graph, we see that the illumination basis clearly “accounts” for the shading on the face quite well, leading to a much lower fluctuation of the residuals. The sequence of images shows an excerpt near the middle of the sequence where the RM+S algorithm (which could not compensate for illumination changes) completely lost the target for several frames, only regaining it after the original lighting was restored. Since the target was effectively motionless during this period, this can be completely attributed to biases due to illumination effects. Similar sequences with larger target motions often cause the purely geometric algorithm to lose the target completely as shown in Fig. 7.

5.3.4 Tracking With Outliers

Finally, we illustrate the performance of the method when the image of the target becomes partially occluded. We again track a face. The motion and illumination basis are

the same as before. In the weighting matrix calculations, the pixel gray-scale variance was set to five (about what is observed in our camera) and the outlier threshold was set to a conservative value of five variance units.

The sequence is an “office” sequence which includes several “intrusions” including the background, a piece of paper, a telephone and a soda can. As before we executed two versions of the tracker, the nonrobust algorithm from the previous experiment (RM+S+I) and a robust version (RM+S+I+O). Fig. 8 shows the results. The upper series of images shows the region acquired by both algorithms (the black frame corresponds to RM+S+I, the white to RM+S+I+O). As is clear from the sequence, the nonrobust algorithm is disturbed significantly by the occlusion, whereas the robust algorithm is much more stable. In fact, a slight motion of the head while the soda can is in the image caused the nonrobust algorithm to mistrack completely. The middle series of images shows the output of the warping operation for the robust algorithm. The lower row of images depicts the weighting values attached to each pixel in the warped image. Dark areas correspond to “outliers.”



Residuals: Face with Partial Occlusion

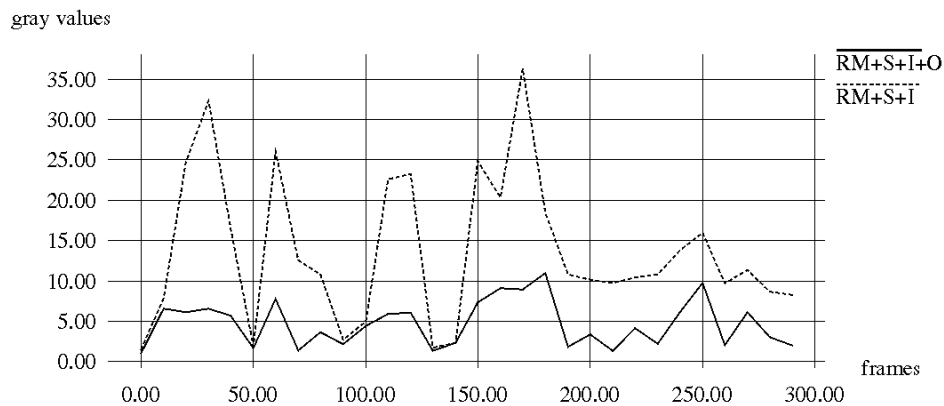


Fig. 8. The first row of images shows excerpts of a tracking sequence with occurrences of partial occlusion. The black frame corresponds to the algorithm without outlier rejection (RM+S+I) and the white frame corresponds to the algorithm with outlier rejection (RM+S+I+O). The second row is a magnified view of the region in the white frame. The third row contains the corresponding outlier images where darker areas mark outliers. The graph at the bottom compares the residual values for both algorithms.

Note that, although the occluded region is clearly identified by the algorithm, there are some small regions away from the occlusion which received a slightly reduced weight. This is due to the fact that the robust metric used introduces some small bias into the computed parameters. In areas where the spatial gradient is large (e.g., near the eyes and mouth), this introduces some false rejection of pixels. At the same time, intruding regions of a similar intensity as the face are not rejected as seen in the lower left of the left-most column of images.

It is also important to note that the dynamical performance of the tracker is reduced by including outliers. Large, fast motions tend to cause the algorithm to “turn off” areas of the image where there are large gradients, slowing convergence. At the same time, performing outlier rejection is more computationally intensive as it requires explicit computation of both the motion and illumination parameters to calculate the residual values.

6 DISCUSSION AND CONCLUSIONS

We have shown a straightforward and efficient solution to the problem of tracking regions undergoing geometric distortion, changing illumination, and partial occlusion. The method is simple and efficient, yet robust to reasonable deviations from underlying motion and illumination models. For example, although we have modeled the face as a rigid object undergoing limited motion in our experiments, the algorithm can still track the subject as he or she is changing expression or, as illustrated in the previous section, performing out-of-plane rotations.

Although the focus in this article has been on parameter estimation techniques for tracking using image rectification, the same estimation methods can be used for directly controlling devices. For example, instead of computing a parameter estimate μ , the incremental solutions $\delta\mu$ can be used to control the position and orientation of a camera so to stabilize the target image by active motion. Hybrid combinations of camera control and image warping are also possible.

One possible objection to the methods is the requirement that the change from frame to frame is small (generally within a few pixels), limiting the speed at which objects can move. Luckily, there are several means for improving the dynamical performance of the algorithms. One possibility is to include a model for the motion of the underlying object and to incorporate prediction into the tracking algorithm. Likewise, if a model of the noise characteristics of images is available, the updating method can be modified to incorporate this model. In fact, the form of the solution makes it straightforward to incorporate the estimation algorithm into a Kalman filter or similar iterative estimation procedure.

Performance can also be improved by operating the tracking algorithm at multiple levels of resolution. One possibility, as is used by many authors [29], [23], is to perform a complete coarse to fine progression of estimation steps on each image in the sequence. Another possibility, which we have used successfully in prior work [41], is to dynamically adapt resolution based on the motion of the target. That is, when the target moves quickly estimation is performed at a coarse resolution, and when it moves slowly the algorithm changes to a higher resolution. The advantage of this approach is that it not only increases the range over which the linearized problem is valid, but it also reduces the computation time required on each image when motion is fast.

We are actively continuing to evaluate the performance of these methods, and to extend their theoretical underpinnings. One area that still needs attention is the problem of determining an illumination basis online, i.e. while tracking the object. Initial experiments in this direction have shown that online determination of the illumination basis can be achieved, although we have not included such results in this paper. As in [29], we are also exploring the use of basis images to handle changes of view or aspect not well addressed by warping.

We are also looking at the problem of extending the method to utilize shape information on the target when such information is available [43]. In particular, it is well known [44] that under orthographic projection, the image

deformations of a surface due to motion can be described with a linear motion model. This suggests that our methods can be extended to handle such models. Furthermore, as with the illumination basis, it may be possible to estimate the deformation models online, thereby making it possible to efficiently track arbitrary objects under changes in illumination, pose, and partial occlusion.

ACKNOWLEDGMENTS

G.D. Hager was supported by ARO grant DAAG55-98-1-0168, U.S. National Science Foundation grant IRI-9420982, and by funds provided by Yale University. P.N. Belhumeur was supported by a Presidential Early Career Award IIS-9703134, a U.S. National Science Foundation Career Award IRI-9703134, and ARO grant DAAH04-95-1-0494. The authors would like to thank David Mumford, Alan Yuille, David Kriegman, Peter Hallinan, and Jorgen Karlholm for contributing to the ideas in this paper.

REFERENCES

- [1] P. Allen, B. Yoshimi, and A. Timcenko, "Hand-Eye Coordination for Robotics Tracking and Grasping," K. Hashimoto, ed., *Visual Servicing*, pp. 33-70. World Scientific, 1994.
- [2] S. Hutchinson, G.D. Hager, and P. Corke, "A Tutorial Introduction to Visual Servo Control," *IEEE Trans. Robotics and Automation*, vol. 12, no. 5, pp. 651-670, 1996.
- [3] N. Papanikolopoulos, P. Khosla, and T. Kanade, "Visual Tracking of a Moving Target by a Camera Mounted on a Robot: A Combination of Control and Vision," *IEEE Trans. Robotics and Automation*, vol. 9, no. 1, pp. 14-35, 1993.
- [4] E. Dickmanns and V. Graefe, "Dynamic Monocular Machine Vision," *Machine Vision and Applications*, vol. 1, pp. 223-240, 1988.
- [5] A.F. Bobick and A.D. Wilson, "A State-Based Technique for the Summarization of Recognition of Gesture," *Proc. Int'l Conf. Computer Vision*, pp. 382-388, 1995.
- [6] T. Darrell, B. Moghaddam, and A. Pentland, "Active Face Tracking and Pose Estimation in an Interactive Room," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 67-72, 1996.
- [7] D. Gavrilu and L. Davis, "Tracking Humans in Action: A 3D Model-Based Approach," *Proc. Image Understanding Workshop*, pp. 737-746, 1996.
- [8] R. Howarth and H. Buxton, "Visual Surveillance Monitoring and Watching," *Proc. European Conf. Computer Vision*, vol. 2, pp. 321-334, 1996.
- [9] T. Frank, M. Haag, H. Kollnig, and H.-H. Nagel, "Tracking of Occluded Vehicles in Traffic Scenes," *Proc. European Conf. Computer Vision*, vol. 2, pp. 485-494, 1996.
- [10] R.C. Harrell, D.C. Slaughter, and P.D. Adsit, "A Fruit-Tracking System for Robotic Harvesting," *Machine Vision and Applications*, vol. 2, pp. 69-80, 1989.
- [11] D. Reynard, A. Wildenberg, A. Blake, and J. Marchant, "Learning Dynamics of Complex Motions From Image Sequences," *Proc. European Conf. Computer Vision*, vol. 1, pp. 357-368, 1996.
- [12] E. Bardinet, L. Cohen, and N. Ayache, "Tracking Medical 3D Data With a Deformable Parametric Model," *Proc. European Conf. Computer Vision*, vol. 1, pp. 317-328, 1996.
- [13] P. Shi, G. Robinson, T. Constable, A. Sinusas, and J. Duncan, "A Model-Based Integrated Approach to Track Myocardial Deformation Using Displacement and Velocity Constraints," *Proc. Int'l Conf. Computer Vision*, pp. 687-692, 1995.
- [14] E. Boyer, "Object Models From Contour Sequences," *Proc. European Conf. Computer Vision*, vol. 2, pp. 109-118, 1996.
- [15] L. Shapiro, *Affine Analysis of Image Sequences*. Cambridge, England: Cambridge Univ. Press, 1995.
- [16] C. Tomasi and T. Kanade, "Shape and Motion From Image Streams Under Orthography: A Factorization Method," *Int'l J. Computer Vision*, vol. 9, no. 2, pp. 137-154, 1992.

- [17] P. Hallinan, "A Low-Dimensional Representation of Human Faces for Arbitrary Lighting Conditions," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 995-999, 1994.
- [18] R. Epstein, P. Hallinan, and A. Yuille, " 5 ± 2 Eigenimages Suffice: An Empirical Investigation of Low-Dimensional Lighting Models," Technical Report 94-11, Harvard Univ., 1994.
- [19] P.N. Belhumeur and D.J. Kriegman, "What Is the Set of Images of an Object Under All Possible Lighting Conditions," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 270-277, 1996.
- [20] R. Dutter and P. Huber, "Numerical Methods for the Nonlinear Robust Regression Problem," *J. Statistical Computer Simulation*, vol. 13, no. 2, pp. 79-113, 1981.
- [21] B.D. Lucas and T. Kanade, "An Iterative Image Registration Technique With an Application to Stereo Vision," *Proc. Int'l Joint Conf. Artificial Intelligence*, pp. 674-679, 1981.
- [22] P. Anandan, "A Computational Framework and an Algorithm for the Measurement of Structure From Motion," *Int'l J. Computer Vision*, vol. 2, pp. 283-310, 1989.
- [23] J. Shi and C. Tomasi, "Good Features to Track," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 593-600, IEEE CS Press, 1994.
- [24] J. Rehag and T. Kanade, "Visual Tracking of High DOF Articulated Structures: An Application to Human Hand Tracking," *Proc. European Conf. Computer Vision*, vol. B, pp. 35-46, 1994.
- [25] C. Bregler, "Learning and Recognizing Human Dynamics in Video Sequences," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 568-574, 1997.
- [26] J. Rehag and A. Witkin, "Visual Tracking With Deformation Models," *Proc. IEEE Int'l Conf. Robotics and Automation*, pp. 844-850, 1991.
- [27] M. Black and Y. Yacoob, "Tracking and Recognizing Rigid and Non-Rigid Facial Motions Using Local Parametric Models of Image Motion," *Proc. Int'l Conf. Computer Vision*, pp. 374-381, 1995.
- [28] H. Murase and S. Nayar, "Visual Learning and Recognition of 3-D Objects From Appearance," *Int'l J. Computer Vision*, vol. 14, pp. 5-24, 1995.
- [29] M. Black and A. Jepson, "Eigentracking: Robust Matching and Tracking of Articulated Objects Using a View-Based Representation," *Proc. European Conf. Computer Vision*, pp. 329-342, 1996.
- [30] M. Isard and A. Blake, "Contour Tracking by Stochastic Propagation of Conditional Density," *European Conf. on Computer Vision*, vol. 1, pp. 343-356, 1996.
- [31] D.G. Lowe, "Robust Model-Based Motion Tracking Through the Integration of Search and Estimation," *Int'l J. Computer Vision*, vol. 8, no. 2, pp. 113-122, 1992.
- [32] D.B. Gennery, "Visual Tracking of Known Three-Dimensional Objects," *Int'l J. Computer Vision*, vol. 7, no. 3, pp. 243-270, 1992.
- [33] A. Blake, R. Curwen, and A. Zisserman, "A Framework for Spatio-Temporal Control in the Tracking of Visual Contour," *Int'l J. Computer Vision*, vol. 11, no. 2, pp. 127-145, 1993.
- [34] B. Horn, *Computer Vision*. Cambridge, Mass.: MIT Press, 1986.
- [35] M. Betke and N. Makris, "Fast Object Recognition in Noisy Images Using Simulated Annealing," *Proc. Int'l Conf. Computer Vision*, pp. 523-530, 1995.
- [36] R. Szeliski, "Image Mosaicing for Tele-Reality Applications," *Proc. Workshop Applications of Computer Vision*, pp. 44-53, 1994.
- [37] A. Shashua, "Geometry and Photometry in 3D Visual Recognition," PhD thesis, Massachusetts Institute of Technology, 1992.
- [38] R. Woodham, "Analysing Images of Curved Surfaces," *Artificial Intelligence*, vol. 17, pp. 117-140, 1981.
- [39] P. Huber, *Robust Statistics*. New York: John Wiley & Sons, 1981.
- [40] R.M. Haralick and L.G. Shapiro, *Computer and Robot Vision*. Reading, Mass.: Addison Wesley, 1993.
- [41] G.D. Hager and K. Toyama, "XVision: A Portable Substrate for Real-Time Vision Applications," *Computer Vision and Image Understanding*, vol. 69, no. 1, pp. 23-37, 1998.
- [42] S. McKenna, S. Gong, and J. Collins, "Face Tracking and Pose Representation," *British Machine Vision Conf.*, 1996.
- [43] P. Belhumeur and G.D. Hager, "Tracking in 3D: Image Variability Decomposition for Recovering Object Pose and Illumination," *Proc. Int'l Conf. Pattern Analysis Applications*, 1998. Also available as Yale Computer Science #1141.
- [44] S. Ullman and R. Basri, "Recognition by a Linear Combination of Models," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 13, pp. 992-1,006, 1991.



Gregory D. Hager received his BA degree in computer science and mathematics from Luther College in 1983 and his MS and PhD in computer science from the University of Pennsylvania in 1985 and 1988, respectively. From 1988 to 1990, he was a Fulbright junior research fellow at the University of Karlsruhe and the Fraunhofer Institute IITB in Karlsruhe, Germany. Upon returning to the United States, he joined the Computer Science Department at Yale University, where he is currently an associate professor. He is a member of IEEE and AAAI and is currently cochairman of the Robotics and Automation Society Technical Committee on Computer and Robot Vision. His research interests include visual tracking, hand-eye coordination, sensor data fusion, and sensor planning. A book on his dissertation work entitled *Task-Directed Sensor Fusion and Planning* has been published by Kluwer Academic Publishers, Inc.



Peter N. Belhumeur graduated in 1985 from Brown University with Highest Honors, receiving an ScB degree in computer and information engineering. He received an SM in 1991 and a PhD in 1993 from Harvard University, where he studied under a Harvard Fellowship. In 1993, he was a postdoctoral fellow at the University of Cambridge's Sir Isaac Newton Institute for Mathematical Sciences. He was appointed assistant professor of electrical engineering at Yale University in 1994 and was given a joint appointment with the Department of Computer Science in 1998. He is a recipient of the Presidential Early Career Award for Scientists and Engineers, the U.S. National Science Foundation Career Award, and a Yale University Junior Faculty Fellowship for Natural Sciences. He won the Best Paper Award at the 1996 *IEEE Conference on Computer Vision and Pattern Recognition* and an Outstanding Paper Award at the 1998 *European Conference on Computer Vision*.