

**RuCCS TR-51**

August, 1999

## Adjusting Model Parameters using Model-Based Optical Flow Residuals

**Doug DeCarlo**

decarlo@cs.rutgers.edu  
Rutger's University  
Department of Computer Science

Center for Cognitive Science

**Dimitris Metaxas**

dnm@central.cis.upenn.edu  
University of Pennsylvania  
Department of Computer and Information  
Science

Technical Report TR-51  
Center for Cognitive Science  
Psych Bldg Addition, Busch Campus  
Rutgers University - New Brunswick  
152 Frelinghuysen Road  
Piscataway, NJ 08854-8020



# Adjusting Model Parameters using Model-Based Optical Flow Residuals

Douglas DeCarlo  
Department of Computer Science  
and Center for Cognitive Science  
Rutgers University  
decarlo@cs.rutgers.edu

Dimitris Metaxas  
Department of Computer and Information Science  
University of Pennsylvania  
dnm@central.cis.upenn.edu

## Abstract

We present a method for estimating the shape and motion of a deformable model using the least-squares residuals from a model-based optical flow computation. This method is built on top of an estimation framework using optical flow and image features, where optical flow affects only the motion parameters of the model. Using the results of this computation, our new method adjusts all of the parameters (both shape and motion) so that the residuals from the flow computation are minimized. We present experiments that demonstrate that this method is a considerable improvement over a framework using only optical flow and features, especially in the estimation of the shape.

## 1 Introduction

Analyzing the error in model-based tracking frameworks can lead to further improvements of the parameter estimates. Typical model-based tracking methods choose a least-squares solution to set the value of those parameters that the measurements directly determine. In over-determined situations, this approach can leave behind a significant residual. Inevitably, this residual reflects deviations that can be caused by idealizations inherent in the model as well as noisy measurements.

But importantly, while parameters affected by the measurements are set as well as possible, parameters that remain unaffected by the measurements may be inaccurate, and so also contribute to the residual. Within this modeling framework, this is a problem we can actually correct.

Consider a model-based optical flow computation, for example. The model parameters are split in those which model motion, and those which model the underlying and unchanging shape. Observed motion directly determines only the motion parameters; the true values of the shape parameters do not change over time. However, inaccurate shape estimates make it impossible for the model to explain the observed motion using the motion parameters—this results in an increased flow residual. In this case, we can adjust the shape parameters to improve the entire estimate.

The challenge is to be faithful to the distinction between shape and motion parameters as the adjustment is computed. For example, it's insufficient to simply stage the computation, and use the leftovers from the motion parameter computation to feed a computation which determines how the shape parameters could have changed over time to explain the remaining observed motion [14]. This method, while effective in image-coding, simply treats shape parameters as second-class motion parameters. We propose computing an adjustment to the shape parameters that minimizes the residual of the motion estimate. In other words, we determine a new configuration for which the motion parameters would have produced a smaller residual in the first place.

## 1.1 Related work

Residuals from one computation are often used as input to another that follows it. For instance, coarse-to-fine methods typically compute a solution at a particular level of detail, and continue on to finer levels after subtracting away its current guess. The method presented here proceeds somewhat differently. In this case, we start with a model-based optical flow computation which provides estimates of motion parameters. The residuals from this computation are then used to adjust *both* the shape and motion parameters of the model.

Our results should also be distinguished from the large body of work on structure from motion which estimates shape and motion using an optical flow field, reviewed in [2]. There has also been a great deal of work on the structure from motion problem using feature correspondences, which is surveyed in [13]. Another approach aligns locations on a model with image features using displacements or gradient fields by solving a template alignment optimization problem [16, 17, 25, 26].

This paper describes an alternative to these structure from motion methods. Our method is coupled to a model-based optical flow computation using a deformable model. Instead of performing direct surface reconstruction from optical flow, our method indirectly adapts model parameters so that the residual of the model-based optical flow is reduced. Section 4 contains a discussion of precisely how these methods differ from our technique.

## 1.2 Shape and motion separation

The starting point for our method is an intuitive distinction between shape and motion; the process of model design must encode information about a class of objects by allowing for the categorization of model parameters as describing either variation in shape or motion. The shape parameters are a *static* quantity for a particular observed object, and describe its unchanging geometric features. The motion parameters are a *dynamic* quantity, which change when the observed object moves or deforms. Of course, there is no guarantee that the shape and motion of some class of objects is separable; this is a simplifying assumption that we make, and will only apply to a certain degree of accuracy. For example, with human faces, shape parameters describe an individual's appearance, while motion parameters encode the location of their head, as well as their facial displays and expressions. This division is often built into face models [4, 15, 18, 24] to simplify model construction or estimation, and has been used to facilitate learning the variability of motions for a class of objects [22].

The ultimate goal of this separation is to produce a simpler estimation problem. During esti-

mation, the change in the shape parameters should tend to zero as the shape of the observed object is established. Once this occurs, fitting need only continue for the motion parameters. Therefore, during model design, the separation into shape and motion should encode as many of the model deformations with shape parameters as possible. This decision also leads to a more efficient tracking system. This distinction leads us to develop a method where changes in the image are initially attributed entirely to motion, but then the residual from the reconstructed motion is used to correct errors in the shape and motion parameters.

For models with separate parameters for shape and motion, certain cues such as optical flow are appropriately used only for the estimation of motion parameters (and not shape parameters). While in some cases updating all the parameters (shape and motion) based on the flow can result in smaller deviations [14], this is missing the point of separating shape and motion in the first place, and is in conflict with the view of the shape parameters as having static values.

### **1.3 Using residuals**

A significant error in the current model estimate will interfere with the optical flow estimates of the motion, since the model and image will be misaligned. For example, if the estimate of a “nose protrusion” parameter is inaccurate, the pattern of motion that the model would predict during a head turn would be incorrect in the local region of the nose. However, this interference is quite systematic, which enables the adjustment of the current shape and motion estimate. This adjustment aims to correct the error which interfered with the flow computation. Continuing with the above example, we would aim to adjust the nose protrusion parameter in a way that improves the motion model’s accuracy. We will perform this adjustment using the residual from the optical flow constraint equation.

From this residual, each pixel used in the optical flow computation supplies one piece of information which is then used to determine how the parameters can be corrected to minimize the in-

interference. However, some pixels will not supply any useful information. And even worse, many of the pixels will include distracting information resulting from optical flow linearization, optical flow constraint violations (such as lighting changes, shadows, or specularities), motion estimation errors, and noise. As a result, we must be sure to use a sufficient number of pixels, as well as to avoid adjusting the parameters based on distracting information. This second point is addressed by the following empirically determined result: residual contributions which result from small errors in the estimated model parameters are significantly larger than those caused by distracting sources (such as optical flow constraint violations and linearization). We demonstrate this empirical result in the context of face tracking, and it is likely to apply in other domains where model-based optical flow based tracking is successful.

Our method is built on top of the model-based face tracking framework described in [7]. This framework used a model-based optical flow computation as a hard constraint on the motion of a deformable model. Aside from flow, the motion of the deformable model was determined by aligning the model with image features (edges). Using features prevented the accumulation of tracking error, which would have otherwise been a problem using flow alone. In Section 5, we will show how using residuals improves the shape and motion estimates of a face. When used with [7], most of the improvement is in the shape parameters, as the motion parameters are already quite accurate based on the estimates using the optical flow and edges.

## 1.4 Outline

In our method, changes in the image are initially attributed entirely to motion, but then the error in the reconstructed motion is used to more accurately extract both shape and motion parameters of the object being tracked. After a brief review of deformable models in Section 2, we discuss existing approaches to model-based optical flow in Section 3. Section 4 describes our method for adjusting the model parameters using residuals from a model-based optical flow computation. Finally, we

present experiments in Section 5 demonstrating our method, along with some discussion.

## 2 Deformable models

Deformable models [17, 21, 25] are parameterized shapes that deform due to forces according to physical laws. For vision applications, physics provides a useful analogy for treating shape estimation [17], where forces are determined from visual cues such as edges in an image. The deformations that result produce a shape that agrees with the data.

The shape of the deformable model  $\mathbf{x}$  is parameterized by a vector of values  $\mathbf{q}$  and is defined over a domain  $\Omega$  which can be used to identify specific points on the model; a particular point on the model is written as  $\mathbf{x}(\mathbf{q}; \mathbf{u})$  with  $\mathbf{u} \in \Omega$ , although the dependency of  $\mathbf{x}$  on  $\mathbf{q}$  is often omitted. The goal of shape and motion estimation is to recover the value of  $\mathbf{q}$  over time from a sequence of images. For this paper, we will be using the three-dimensional parameterized face model from [7].

As stated earlier, to distinguish the processes of shape estimation and motion tracking, the parameters in  $\mathbf{q}$  are rearranged and separated into  $\mathbf{q}_b$  (the basic shape of the object) and  $\mathbf{q}_m$  (rigid and non-rigid motion), so that  $\mathbf{q} = (\mathbf{q}_b^\top, \mathbf{q}_m^\top)^\top$ . With our face model,  $\mathbf{q}_b$  describes an individual's appearance, while  $\mathbf{q}_m$  encodes the location of their head, as well as their facial displays and expressions. Figure 1 shows examples of the face model undergoing various shape deformations (showing four different individuals), motion deformations (showing brow raising and frowning, smiling, and mouth opening) and finally two examples of when several deformations are applied at once.

The model  $\mathbf{x}$  is formed by applying deformation functions to the underlying shape  $\mathbf{s}$ . For this paper, the underlying face model  $\mathbf{s}$  is a polygon mesh (shown in the center of Figure 1). There are separate deformation functions for shape ( $\mathbf{T}_b$ ) and for motion ( $\mathbf{T}_m$ ). The shape deformation is

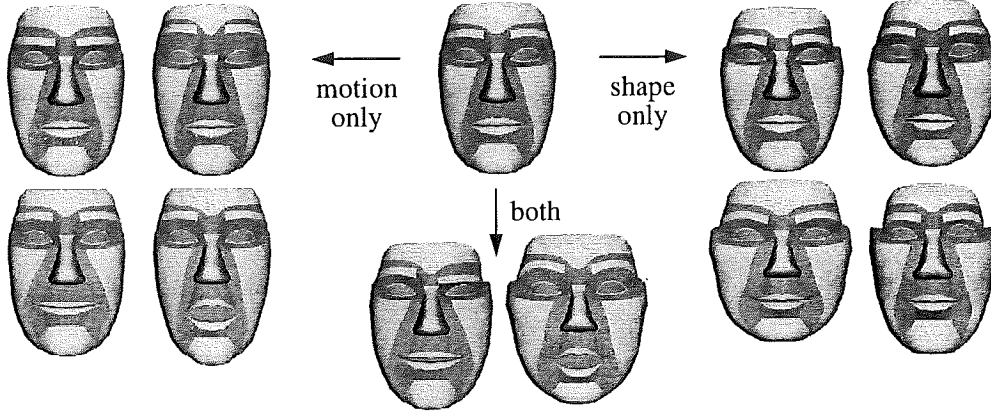


Figure 1: Example parameterized deformations of the face model (with separate parameters for shape and motion)

applied first, so that:

$$\mathbf{x}(\mathbf{q}; \mathbf{u}) = \mathbf{T}_m(\mathbf{q}_m; \mathbf{T}_b(\mathbf{q}_b; \mathbf{s}(\mathbf{u}))) \quad (1)$$

The shape deformation  $\mathbf{T}_b$  uses the parameters  $\mathbf{q}_b$  to deform the underlying shape  $\mathbf{s}$ . On top of this is the motion deformation  $\mathbf{T}_m$  with parameters  $\mathbf{q}_m$ , which includes a rigid translation and rotation (head motion), as well as non-rigid deformations (facial expressions and displays).

When modeling an object viewed in images,  $\mathbf{x}$  must include a camera projection, resulting in a two-dimensional model called  $\mathbf{x}_p$ , which is projected flat from the original three-dimensional model.

## 2.1 Kinematics and Dynamics

The kinematics of the model are determined in terms of the parameter velocities  $\dot{\mathbf{q}}$ . As the shape changes, the velocity at a point  $\mathbf{u}$  on the model is given by:

$$\dot{\mathbf{x}}(\mathbf{u}) = \mathbf{L}(\mathbf{q}; \mathbf{u})\dot{\mathbf{q}} \quad (2)$$

where  $\mathbf{L} = \partial \mathbf{x} / \partial \mathbf{q}$  is the model Jacobian [17]. For reasons of conciseness, the dependency of  $\mathbf{L}$  on  $\mathbf{q}$  is often omitted.

We view  $\mathbf{L}$  as consisting of components that correspond to  $\mathbf{q}_b$  and  $\mathbf{q}_m$ , so that it can be written as  $[\mathbf{L}_b \ \mathbf{L}_m]$ . The Jacobian of the projected model  $\mathbf{x}_p$  is written as  $\mathbf{L}_p$ , and is decomposed into components for  $\mathbf{q}_b$  and  $\mathbf{q}_m$  as  $[\mathbf{L}_{bp} \ \mathbf{L}_{mp}]$ . In Section 3, we use this distinction to attribute observed motion as resulting solely from changes in the motion parameters.

The models defined above are useful for applications such as shape and motion estimation when used in a physics-based framework [17]. These techniques are a form of optimization whereby the deviation between the model and the data is minimized. The optimization is performed by integrating differential equations derived from the Euler-Lagrange equations of motion. These equations are simplified in a standard manner [17], and in this case result in:

$$\dot{\mathbf{q}} = \mathbf{f}_q \tag{3}$$

where the applied forces  $\mathbf{f}_q$  are computed from two-dimensional image forces  $\mathbf{f}_{\text{image}}$  as:

$$\mathbf{f}_q = \sum_j \mathbf{L}_p(\mathbf{u}_j)^\top \mathbf{f}_{\text{image}}(\mathbf{u}_j) \tag{4}$$

The distribution of forces on the model is based in part on forces computed from the edges of an input image [17]. With that, and given an adequate model initialization, these forces will align features on the model with image features, thereby determining appropriate parameter values. The dynamic system in (3) is solved by integrating over time, using standard (explicit) differential equation integration techniques, such as Euler integration:

$$\mathbf{q}(t + \Delta t) = \mathbf{q}(t) + \dot{\mathbf{q}}(t)\Delta t \tag{5}$$

The initialization which specifies the value of  $\mathbf{q}(0)$  is described in Section 5.

### 3 Model-based optical flow

The optical flow is defined as the apparent motion of brightness patterns across an image [11]. Attempting to use this information in applications such as object tracking requires assumptions about the objects being viewed. Most common is the assumption that particular locations on viewed objects do not change in brightness. This brightness constancy assumption leads to the formulation of the well-known optical flow constraint equation for the image  $I$  at a particular pixel (the assumption manifests itself as the zero on the right-hand-side):

$$\nabla I \begin{bmatrix} u \\ v \end{bmatrix} + I_t = 0 \quad (6)$$

where  $\nabla I = [I_x \ I_y]$  are the spatial derivatives and  $I_t$  is the temporal derivative of the image intensity.  $u$  and  $v$  are the components of the image velocities.

The model-based optical flow constraint equation is a reformulation of (6) in terms of a model's motion parameters  $\mathbf{q}_m$ . When viewing a model under projection, there exists a unique model point  $\mathbf{u} \in \Omega$  which corresponds to a particular pixel (except on occluding boundaries and situations involving transparency). In a model-based approach, the image velocities  $u$  and  $v$  are specified by projected velocities of points on the model  $\dot{\mathbf{x}}_p(\mathbf{u})$ :

$$\begin{bmatrix} u \\ v \end{bmatrix} = \dot{\mathbf{x}}_p(\mathbf{u}) = \mathbf{L}_{mp}(\mathbf{u})\dot{\mathbf{q}}_m \quad (7)$$

Note that only the changes resulting from the motion parameters  $\mathbf{q}_m$  are included, as optical flow velocities do not reflect changes in the shape parameters  $\mathbf{q}_b$ , so that only  $\mathbf{L}_{mp}$  is used here (and not

all of  $\mathbf{L}_p$ ). The model-based optical flow constraint equation rewrites (6) using (7):

$$\nabla \mathbf{I} \mathbf{L}_{m,p}(\mathbf{u}) \dot{\mathbf{q}}_m + \mathbf{I}_t = 0 \quad (8)$$

When considered over a set of  $n$  pixels, a stacked set of instances of (8) can be written in matrix form (where each row corresponds to an individual pixel):

$$\mathbf{B} \dot{\mathbf{q}}_m + \mathbf{I}_t = 0 \quad (9)$$

Formulations similar to (9) (although superficially appearing quite different) can be found in [1, 3, 6, 12, 14, 15, 19, 20].

The value of  $\dot{\mathbf{q}}_m$  in (9) can be determined by least-squares minimization:

$$\min_{\dot{\mathbf{q}}_m} [\mathbf{B} \dot{\mathbf{q}}_m + \mathbf{I}_t]^2 \quad (10)$$

where  $[\mathbf{v}]^2 = \mathbf{v}^\top \mathbf{v}$ . An iterative approach to solving this is taken in [3, 12, 14, 19, 20]. Alternatively, it can be solved in a single step using the pseudo-inverse (where  $\mathbf{B}^+$  is the pseudo-inverse of  $\mathbf{B}$  [23]) [6, 20, 15]:

$$\dot{\mathbf{q}}_m = -\mathbf{B}^+ \mathbf{I}_t \quad (11)$$

This is the linear least-squares solution; it linearizes by assuming  $\mathbf{L}_{m,p}$  is constant (ignoring its dependency on  $\mathbf{q}$ ). This is the solution method we employ here.

The derivation of (6) involves the truncation of a Taylor series, and as a result requires relatively small motions between frames. To address problems with estimating larger motions (or to implement coarse-to-fine methods), some iterative approaches transform the model geometry at each iter-

ation using the previous motion estimate [14], while others undo the previous estimate by warping the input images [3].

The most serious difficulty for these techniques is combating tracking drift. Using only velocity information, small estimation errors accumulate over time. The solution to this problem is to include other information (such as features or edges) to prevent errors from building up [7, 14, 15].

### 3.1 Optical Flow Residuals

Unlike image-based optical flow techniques, these model-based methods do not require assumptions about the smoothness of the flow field to determine a solution, as the number of pixels providing useful information is sufficiently greater than the number of motion parameters. Of course, now that the solution is over-determined, there will be a residual from the least-squares solution, given the solution  $\hat{\mathbf{q}}_m$ :

$$R = \mathbf{B}\hat{\mathbf{q}}_m + \mathbf{I}_t \quad (12)$$

The residual  $R$  is a vector having dimension  $n$  (the number of pixels used in the flow computation).

The residual comes from many sources. Aside from measurement noise, most obviously are linearization errors, resulting from ignoring the higher order terms in (6) (which were truncated in the Taylor series). With the linear least squares solution, the linearization of  $\mathbf{L}_{mp}$  is also a culprit. Others include violations of the brightness constancy assumption, such as lighting changes, shadows, and specularities. Finally, shape and motion estimation errors (deviation between  $\mathbf{q}$  and its actual value) will prevent the model from properly aligning with the image, and will cause a sizeable increase in the residual.

In fact, we claim that if significant errors are present in the estimated shape and motion, they will be the primary contributors to the residual. We support this claim empirically in Section 5. This

means that the residual is a valuable piece of information that can be used for estimating shape and motion. The next section describes how to compute small adjustments to the shape and motion parameters which reduce the residual.

## 4 Adjusting parameters using residuals

To provide insight into this method, it's worth comparing the use of a three-dimensional model under projection in a model-based optical flow framework [7, 15], with a two-dimensional model of image motion [4]. For a particular value of  $\mathbf{q}$ , the 3D model, when projected as a 2D model in the image, is actually quite comparable with a 2D image motion model. The key difference is that the form of the 3D model when projected into the image changes with  $\mathbf{q}$ —so that changes in parameters of the 3D model directly affects its projected motion in the image. So we can now ask: how could this 2D motion parameterization have been different (by changing  $\mathbf{q}$ ), which would have resulted in a smaller residual? The remainder of this section describes how we answer this question by adjusting the model's parameters using the residuals.

There are various approaches to using this residual information. One possible approach explains the residual as directly resulting from shape deviation—this is basically a structure from motion approach. In other words, the leftover motion not accounted for by motion parameters is used to update the shape using the same formulation as for determining motion (from Section 3), as in:

$$\mathbf{B}\dot{\mathbf{q}}_m + \mathbf{B}_b\dot{\mathbf{q}}_b + \mathbf{I}_t = \mathbf{0} \quad (13)$$

where the construction of  $\mathbf{B}_b$  is analogous to  $\mathbf{B}$ , but uses  $\mathbf{L}_b$  instead of  $\mathbf{L}_m$ . Instead of solving one large system, the system in (13) is decoupled, and is solved for motion first, and then for shape in

terms of the residual  $R$ :

$$\mathbf{B}_b \dot{\mathbf{q}}_b + R = \mathbf{0} \quad \Rightarrow \quad \dot{\mathbf{q}}_b = -\mathbf{B}_b^+ R \quad (14)$$

This method is closely related to one described by Koch [14]. In that work, shape parameters are actually updated in two steps; first using discrepancies parallel to the line of sight, then those that are perpendicular to the line of sight. This problem involves solving the following minimization:

$$\min_{\dot{\mathbf{q}}_b} [\mathbf{B} \dot{\mathbf{q}}_m + \mathbf{B}_b \dot{\mathbf{q}}_b + \mathbf{I}_t]^2 \quad (\text{given } \dot{\mathbf{q}}_m) \quad (15)$$

This is basically a staged version of the problem in (10), where  $\dot{\mathbf{q}}_m$  is determined first, followed by  $\dot{\mathbf{q}}_b$ .

This is a reasonable approach in the context of image coding, where image fidelity is of much greater importance than the accuracy of the face shape estimate—the face shape is deformed to account for the tracking errors in motion. This produces a face shape that results in a better image, but does not necessarily estimate the actual face shape of the subject. Plus, given our distinction between shape and motion parameters, it does not make sense to adjust the shape parameters  $\mathbf{q}_b$  directly from observed velocities, as in [14], since the true value of  $\mathbf{q}_b$  is a static quantity.

Instead of this, our approach determines the small change in  $\mathbf{q}$  that effects the largest reduction in  $R$ . Let  $\Delta \mathbf{q}$  be the current deviation of  $\mathbf{q}$  from its true value (not including the motion extracted in  $\dot{\mathbf{q}}_m$ )—this includes both the shape error and the accumulated motion error. We assume  $\Delta \mathbf{q}$  is of sufficiently small magnitude so that the first-order approximation to  $\mathbf{L}_m$  using its Taylor-series expansion is sufficiently accurate:

$$\mathbf{L}_{mp}(\mathbf{u}; \mathbf{q} + \Delta\mathbf{q}) \approx \mathbf{L}_{mp}(\mathbf{u}; \mathbf{q}) + \frac{\partial \mathbf{L}_{mp}(\mathbf{u}; \mathbf{q})}{\partial \mathbf{q}} \Delta\mathbf{q} \quad (16)$$

Combining this approximation of  $\mathbf{L}_{mp}$  with the model-based optical flow constraint equation (8) results in:

$$\nabla \mathbf{I} \mathbf{L}_{mp}(\mathbf{u}) \dot{\mathbf{q}}_m + \nabla \mathbf{I} \left( \frac{\partial \mathbf{L}_{mp}(\mathbf{u})}{\partial \mathbf{q}} \Delta\mathbf{q} \right) \dot{\mathbf{q}}_m + \mathbf{I}_t = 0 \quad (17)$$

where  $\partial \mathbf{L}_{mp} / \partial \mathbf{q}$  is part of the model Hessian matrix (a rank 3 tensor). It is written here “curried” with  $\Delta\mathbf{q}$  so that the parenthesized sub-expression here is a matrix.

When (17) is considered over  $n$  pixels from the input image, this results in the system:

$$\mathbf{B} \dot{\mathbf{q}}_m + (\mathbf{G} \dot{\mathbf{q}}_m) \Delta\mathbf{q} + \mathbf{I}_t = \mathbf{0} \quad (18)$$

$$\text{where } \mathbf{G} = \begin{bmatrix} \left( \nabla \mathbf{I}_1 \frac{\partial \mathbf{L}_{mp}(\mathbf{u}_1)}{\partial \mathbf{q}} \right)^\top \\ \vdots \\ \left( \nabla \mathbf{I}_n \frac{\partial \mathbf{L}_{mp}(\mathbf{u}_n)}{\partial \mathbf{q}} \right)^\top \end{bmatrix} \quad (19)$$

The subscripts  $[1 \dots n]$  in the construction of  $\mathbf{G}$  correspond to a particular row in (18). The transpositions performed in the construction of  $\mathbf{G}$  allow it now to be curried with  $\dot{\mathbf{q}}_m$  (this construction transposes the second and third indices of the tensor  $\mathbf{G}$ ). We can now rewrite (18) using the residual

(12) and the value of  $\dot{\mathbf{q}}_m$  supplied by the model-based optical flow solution in (11):

$$-(\mathbf{GB}^+\mathbf{I}_t)\Delta\mathbf{q} + R = \mathbf{0} \quad (20)$$

Then,  $\Delta\mathbf{q}$  is determined using the pseudo-inverse:

$$\Delta\mathbf{q} = (\mathbf{GB}^+\mathbf{I}_t)^+ R \quad (21)$$

This least squares solution determines the best set of small changes in  $\mathbf{q}_b$  and  $\mathbf{q}_m$  that minimize the optical flow residual (12), given the linearization of  $\mathbf{L}_{m,p}$  in (16). This effectively solves the following minimization:

$$\min_{\Delta\mathbf{q}} [\mathbf{B}(\mathbf{q} + \Delta\mathbf{q})\dot{\mathbf{q}}_m + \mathbf{I}_t]^2 \quad (\text{given } \dot{\mathbf{q}}_m) \quad (22)$$

where  $\mathbf{B}(\mathbf{q} + \Delta\mathbf{q})$  is approximated as  $\mathbf{B}(\mathbf{q}) + \partial\mathbf{B}/\partial\mathbf{q} \cdot \Delta\mathbf{q}$ . The intuition for this analysis—that we’re solving a minimization process that is faithful to the distinction between shape and motion parameters—is realized here in the formulation of a minimization problem very different from (15). Note that it is possible that the new value of  $R$  would be smaller if (15) is used over (22), but this goes against the assumption that  $\mathbf{q}_b$  are static parameters, and would result in an inappropriate estimate.

## 4.1 Solution improvement

The value of  $\Delta\mathbf{q}$  from the previous section specifies an absolute update to the state (unrelated to the current time step  $\Delta t$ )—we could simply add  $\Delta\mathbf{q}$  after each iteration.

$$\mathbf{q}(t + \Delta t) = \mathbf{q}(t) + \dot{\mathbf{q}}\Delta t + \Delta\mathbf{q} \quad (23)$$

However, if information is available about the uncertainty in the estimate of  $\dot{\mathbf{q}}$ , we can add in  $\Delta\mathbf{q}$  in a more principled manner.

In [9], a model-based optical flow solution is combined with edge information using an extended Kalman filter, which results in a filtered estimate (using both the flow and edge information) of  $\dot{\mathbf{q}}$  and a covariance estimate  $\Lambda_{\dot{\mathbf{q}}}$ .

If we assume the distracting sources in  $R$  (aside from  $\Delta\mathbf{q}$ ) can be modeled as (zero mean) Gaussian disturbances, then from (20), the covariance of  $\Delta\mathbf{q}$  is:

$$\Lambda_{\Delta\mathbf{q}} = \left( (\mathbf{GB} + \mathbf{I}_t)^\top \Lambda_R^{-1} (\mathbf{GB} + \mathbf{I}_t) \right)^{-1} \quad (24)$$

where  $\Lambda_R$  is the covariance of  $R$  (which we choose to be a diagonal matrix with diagonal entry  $\sigma_R$ ). The value of  $\sigma_R$  represents the contributions to  $R$  from sources other than shape and motion estimation errors, and is determined in Section 5.3 from experiments where  $\mathbf{q}_b$  (the shape) is known in advance. Using this covariance information ( $\Lambda_{\dot{\mathbf{q}}}$  and  $\Lambda_{\Delta\mathbf{q}}$ ), the new value of  $\mathbf{q}$  is found using cue integration techniques [5, 10], which weight  $\dot{\mathbf{q}}\Delta t$  and  $\Delta\mathbf{q}$  together based on their uncertainty.

The method in (23), which does not use the covariance information, is not particularly robust. At first [8], this was attributed to there being a poor linear approximation of  $\mathbf{L}_{mp}$ , which was tested for by determining if the residual actually does decrease with the addition of  $\Delta\mathbf{q}$ . It seems this was actually due to the distracting sources having a significant effect on the value of  $\Delta\mathbf{q}$ .

## 4.2 Implementation

Solving (21) is made more efficient by omitting parameters in the construction of  $\mathbf{G}$  which cannot be affected based on  $\dot{\mathbf{q}}_m$ . For example, if there is no motion extracted in the eyebrow region of the face, then there is no reason to include eyebrow shape parameters in  $\mathbf{G}$ . Another example is when the extent of a parameter is simply not visible in the image. At any point in time, typically about

half of the shape parameters of the face model can be omitted from the computations.

The process of determining  $\Delta\mathbf{q}$  can also be iterated, solving (11) and (21) repeatedly to obtain a greater improvement. For the applications here, the linear approximation in (16) is relatively accurate with this face model, due to the fact that most of the model parameterization is linear scaling. As a result, only the single iteration is performed.

## 5 Experiments

This section describes a number of face tracking experiments on two image sequences. We demonstrate how the adjustment method using residuals is a significant improvement over the framework in [9]. We also justify our assumption that parameter estimation errors are the leading contributor to the residuals. For the remainder of this section, we will be comparing two frameworks. First, is the original framework from [9], which uses optical flow and edges (and is basically an Extended Kalman filter version of [7]). Second is the framework presented here, which is used on top of the framework of [9]. It determines  $\Delta\mathbf{q}$  using (21) which is statistically combined with the original solution from [9] using (24).

### 5.1 Setup

The original image sequences are 8 bit gray images at NTSC resolution (480 vertical lines). In the sequences, the width of the face in the image averages 200 pixels. A single subject is used in both experiments presented here. The shape (determined by  $\mathbf{q}_b$ ) is validated using a Cyberware range scan of the subject, shown in Figure 2.

The entire estimation process is automatic, except for the initialization, which requires the manual specification of several landmark features in the first frame of the sequence (the eyebrow centers, eye corners, nose tip, and mouth corners). The subject must also be at rest and (approximately) fac-

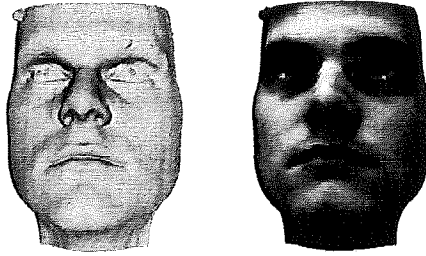


Figure 2: Range scan of the subject (shaded and textured)

ing forward. Experience has shown that the initialization process is robust to small displacements (i.e. several pixels) in the selected landmark points. Details of this initialization process are provided in [7].

For each of the tracking examples, several frames from the image sequence are displayed, cropped appropriately. Below each, the same sequence is shown with the estimated face superimposed. In each case, a model initialization is performed as described below. The initialization process usually takes about 2 minutes of computation. Afterwards, processing each frame using the method in [9] takes approximately 1.4 seconds each. The error residual computation adds an additional 8 seconds per frame (all computation times are measured on a 175 MHz R10000 SGI O2). For both methods, 120 pixels are used in the optical flow computation (the value of  $n$ ).

## 5.2 Estimation experiments

The shape estimation validation experiment in Figure 3 shows the subject making a series of non-rigid face motions: opening his mouth in (b) and (c), smiling in (d) through (e), and finally raising his eyebrows in (f). At each frame, Figure 4 shows the extracted shape results as compared against the range scan of the subject, for both techniques. Note that for this comparison, all motion parameters are ignored, so that only the shape is compared (ground truth for motion is not available). The RMS error is computed using the nodes of the model, and also includes a uniform scaling of the model so that the two faces are the same scale (this eliminates the depth ambiguity—in this case,

the estimated model was compared at 96% scale).

For the estimation system in [9] (the dotted line), the RMS error starts at around 1.7 cm after initialization, and shows a steady reduction over the course of the experiment, ending around 1.3 cm. For the system using error residuals (the solid line), the RMS error again starts at around 1.7 cm, but ends with less error (0.85 cm) compared to the edge-based technique.

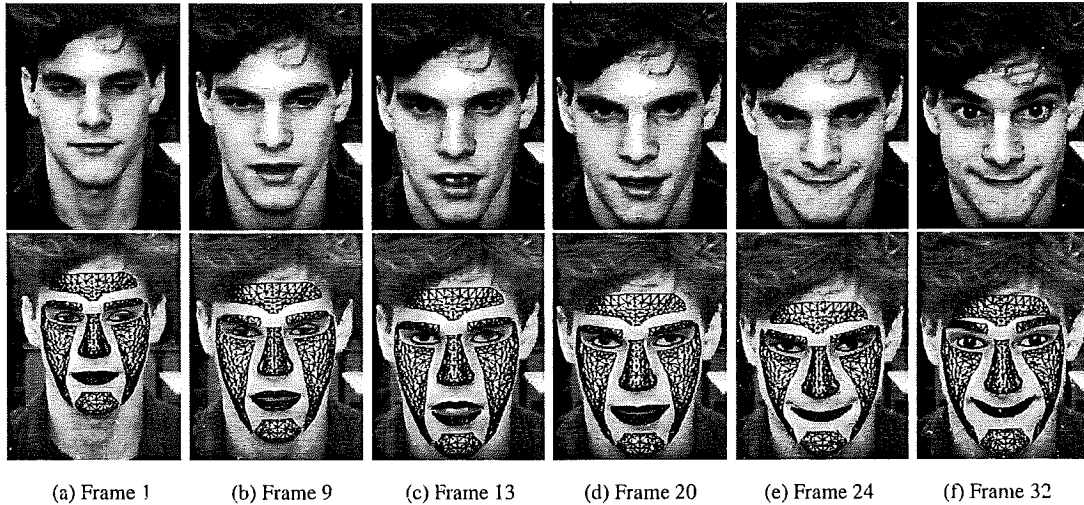


Figure 3: Shape estimation experiment 1

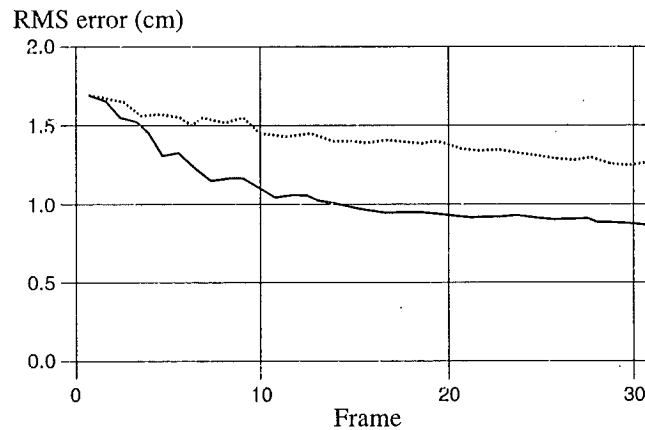


Figure 4: Results of experiment 1

The experiment in Figure 5 shows the subject performing small head motions in (a) through (f) while smiling in (c) and (d), and finishing with a significant head rotation in (g). Using the method from [9] (again, the dotted line), the RMS error starts at around 1.9 cm after initialization, and shows

a gradual reduction over the course of the experiment, ending just under 1 cm, with the large reduction in error around frame 50 corresponding to when the subject turned his head significantly to the side in Figure 5(f) and (g), where the profile view contained good edge information to fit the face shape. For the system using error residuals (the solid line), the RMS error again starts at around 1.9 cm, but this time finishes with just under half of the RMS error as the edge-based technique: around 0.4 cm. In addition, this lower level was reached fairly quickly, showing the advantage of using the error residual technique.

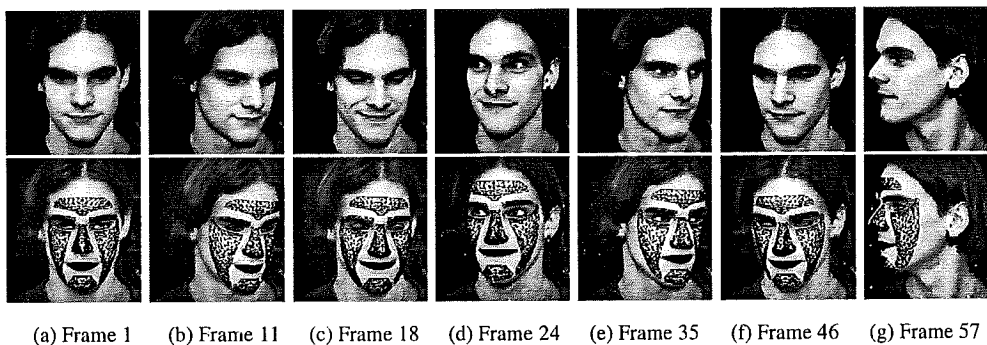


Figure 5: Shape estimation experiment 2

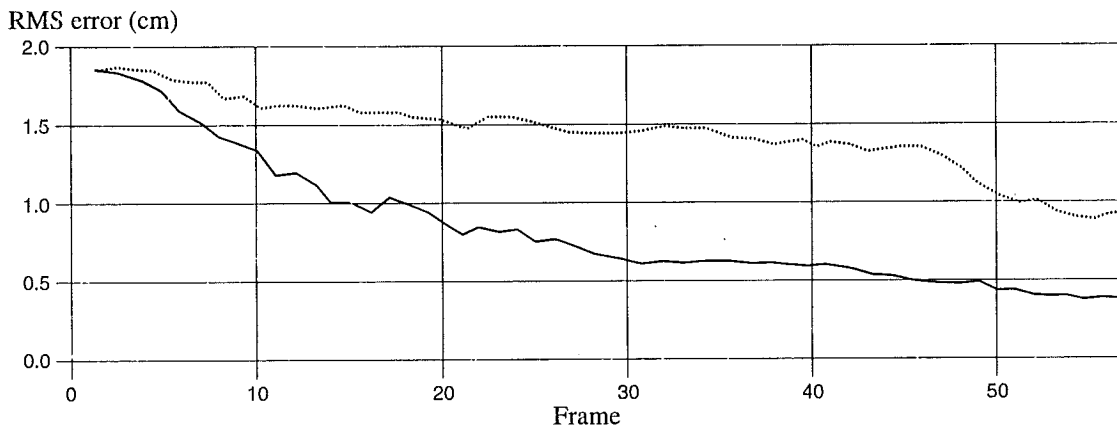


Figure 6: Results of experiment 2

In both experiments, the motion parameter values change appropriately, and at the correct times. Both techniques extracted virtually the same motion parameter values. This is not particularly surprising, as the edge information will maintain fairly accurate estimates of the motion parameters. Thus, this method is mainly a benefit for shape estimation. However, this does not mean that the

motion parameters can be omitted in the formation of  $\mathbf{G}$ , as this would mean that deviations that could be best explained with motion parameters will be incorrectly explained using shape parameters.

### 5.3 Analysis of residuals

The derivation of the method using the residuals in Section 4 assumes that shape error is the leading contributor to the residuals from the motion computation. We now describe our analysis of the residual magnitudes to justify this claim. First, we define the average residual magnitude (in pixel intensity units in the range  $[0, 1]$ ) as:

$$\|R\|_{\text{average}} = \frac{1}{n} \sum_i |R_i| \quad (25)$$

where  $n$  is the number of pixels used (the dimension of  $R$ ).

During both tracking experiments (using the residual-based method), the average residual magnitudes started fairly high: initially 0.18 for the first experiment, and 0.24 for the second. At the end, the values were 0.049 for the first experiment and 0.051 for the second. To isolate the portion of the residuals caused by shape error, both experiments were run again; this time, the initial model shape was taken from the range scan of the subject. As a result, there should be no (or minimal) shape error during the tracking sequence. The residuals that resulted from these experiments had a fairly small magnitude (and were basically the same across both experiments), averaging around 0.035 (ranging from 0.013 to 0.057).

This enforces the validity of our assumption that parameter estimation error is responsible for the bulk of the residual. Aside from this, the data from this analysis provided a residual variance value ( $\sigma_R$ ) of 0.044, used in (24).

## 5.4 Discussion

Besides having improved accuracy over a framework using only optical flow and edges, our framework extracts the shape of the face without needing data from such extreme head poses (such as a profile view). Instead, fewer observations are needed to extract the shape, so that the static part of the estimation problem converges sooner.

In addition, once the static estimation problem is complete, there is no need to perform adjustments in this application, as there was little improvement in the motion estimates. We can also skip the adjustment computation in situations where the residual is small (compared to  $\sigma_R$ ), as they would tend to provide little useful information.

We find that performing adjustments using this method is quite robust, at least to the same level as the underlying flow computation. Of course, in situations where there is large optical flow violation (such as a major lighting change), both the adjustment method and the model-based optical flow computation will produce poor results. Any small deviations are either ignored (due to the cue integration) or are corrected by the edge information. This method is quite complementary to the original framework presented in [7], especially when a Kalman filter is added [9].

## 6 Conclusions and future work

We have presented a novel deformable model technique which uses residuals from a model-based optical flow solution to refine the shape and motion of the model. By using the relationship between the shape and motion parameterizations, small improvements to the parameters are made by minimizing the model-based optical flow residuals. It was the separation of the parameterization which made this computation possible, since additional parameters that did not apply in the model-based flow computations could still be adjusted. While this method is presented in the context of face tracking, it could certainly be applied in other model-based domains with separable parameteriza-

tions.

The adjustment computation, along with cue integration methods seem to be robust to optical flow constraint equation violations or approximations (such as small lighting changes or ignoring higher order image derivative terms). Besides having greater accuracy than a framework using only optical flow and edges, our framework extracts the shape of the face without needing data from extreme head poses (such as a profile view). Instead, much smaller subject motions are required to extract the shape information.

## Acknowledgments

This research is supported by ARO grant DAAH-04-96-1-007; NSF Career Award (IRI-96-24604); AFOSR and ONR-YIP N00014-97-1-0817; and a RuCCS postdoctoral fellowship.

## References

- [1] G. Adiv. Determining 3-D motion and structure from optical flow generated by several moving objects. *IEEE Pattern Analysis and Machine Intelligence*, 7(4):384–401, July 1985.
- [2] John Barron. Computing motion and structure from noisy time-varying image velocity information. Technical Report RBCV-TR-88-24, Department of Computer Science, University of Toronto, 1988.
- [3] J. Bergen, P. Anandan, K. Hanna, and R. Hingorani. Hierarchical model-based motion estimation. In *Proceedings ECCV '92*, pages 237–252, 1992.
- [4] M. Black and Y. Yacoob. Tracking and recognizing rigid and non-rigid facial motions using local parametric models of image motion. In *Proceedings ICCV '95*, pages 374–381, 1995.
- [5] A. Bryson and Y. Ho. *Applied Optimal Control*. Halsted Press, 1975.
- [6] C. Choi, K. Aizawa, H. Harashima, and T. Takebe. Analysis and synthesis of facial image sequences in model-based image coding. *IEEE Circuits and Systems for Video Technology*, 4(3):257–275, 1994.
- [7] D. DeCarlo and D. Metaxas. The integration of optical flow and deformable models with applications to human face shape and motion estimation. In *Proceedings CVPR '96*, pages 231–238, 1996.

- [8] D. DeCarlo and D. Metaxas. Deformable model-based shape and motion analysis from images using motion residual error. In *Proceedings ICCV '98*, pages 113–119, 1998.
- [9] D. DeCarlo and D. Metaxas. Combining information using hard constraints. In *Proceedings CVPR '99*, pages 132–138, 1999.
- [10] H.F. Durrant-Whyte. Consistent integration and propagation of disparate sensor observations. *International Journal of Robotics Research*, 6(3):3–24, 1987.
- [11] B. Horn. *Robot Vision*. McGraw-Hill, 1986.
- [12] B. Horn and E. Weldon. Direct methods for recovering motion. *International Journal of Computer Vision*, 2(1):51–76, June 1988.
- [13] T.S. Huang and A.N. Netravali. Motion and structure from feature correspondences: A review. *PIEEE*, 82(2):252–268, February 1994.
- [14] R. Koch. Dynamic 3-D scene analysis through synthesis feedback control. *IEEE Pattern Analysis and Machine Intelligence*, 15(6):556–568, June 1993.
- [15] H. Li, P. Roivainen, and R. Forchheimer. 3-D motion estimation in model-based facial image coding. *IEEE Pattern Analysis and Machine Intelligence*, 15(6):545–555, June 1993.
- [16] D.G. Lowe. Fitting parameterized three-dimensional models to images. *IEEE Pattern Analysis and Machine Intelligence*, 13(5):441–450, May 1991.
- [17] D. Metaxas and D. Terzopoulos. Shape and nonrigid motion estimation through physics-based synthesis. *IEEE Pattern Analysis and Machine Intelligence*, 15(6):580–591, June 1993.
- [18] Y. Moses, D. Reynard, and A. Blake. Robust real time tracking and classification of facial expressions. In *Proceedings ICCV '95*, pages 296–301, 1995.
- [19] S. Negahdaripour and B. Horn. Direct passive navigation. *IEEE Pattern Analysis and Machine Intelligence*, 9(1):168–176, January 1987.
- [20] A. Netravali and J. Salz. Algorithms for estimation of three-dimensional motion. *AT&T Technical Journal*, 64:335–346, 1985.
- [21] A. Pentland and S. Sclaroff. Closed-form solutions for physically based shape modeling and recognition. *IEEE Pattern Analysis and Machine Intelligence*, 13(7):715–729, 1991.
- [22] D. Reynard, A. Wildenberg, A. Blake, and J. Marchant. Learning dynamics of complex motions from image sequences. In *Proceedings ECCV '96*, pages I:357–368, 1996.
- [23] G. Strang. *Linear algebra and its applications*. Harcourt, Brace, Jovanovich, 1988.
- [24] D. Terzopoulos and K. Waters. Analysis and synthesis of facial image sequences using physical and anatomical models. *IEEE Pattern Analysis and Machine Intelligence*, 15(6):569–579, 1993.

- [25] D. Terzopoulos, A. Witkin, and M. Kass. Constraints on deformable models: Recovering 3D shape and nonrigid motion. *Artificial Intelligence*, 36(1):91–123, 1988.
- [26] A.L. Yuille, D.S. Cohen, and P. Halliman. Feature extraction from faces using deformable templates. *International Journal of Computer Vision*, 8:104–109, 1992.

# RuCCS Technical Report Series

Rutgers Center for Cognitive Science (RuCCS) hosts two electronic archives accessible through our Web Site or anonymous FTP. One contains papers from the Technical Report series (listed below). The prices indicated are only if you wish to order the hard copy of the report (FTP is free). The other series is the Rutgers Optimality Archive (ROA) which is open to all interested in Optimality Theory (of which TR-2 is the founding document). The Archive contains numerous papers authored by a broad spectrum of researchers and a bibliography of hundreds of papers, books and theses in Optimality Theory.

The URL for RuCCS is  
<http://ruccs.rutgers.edu>

Anonymous FTP paths:

RuCCS-TR Archive

/pub/papers/

Rutgers Optimality Archive

/pub/OT/Texts/

Or, to FTP directly, you would:

**ftp ruccs.rutgers.edu**

If you want to explore these facilities from within RuCCS, the entire path is  
 /home/ruccs/ftp/pub/...

<b>TR-1</b> Alan Prince \$7.40 <b>In Defense of the Number 1: Anatomy of a Linear Dynamical Model of Linguistic Generalizations</b>	<b>TR-2</b> Alan Prince Paul Smolensky 12.25 <b>Optimality Theory: Constraint Interaction in Generative Grammar</b>	<b>TR-3</b> John McCarthy Alan Prince 10.50 <b>Prosodic Morphology I: Constraint Interaction and Satisfaction</b>
<b>TR-4</b> Jane Grimshaw \$3.15 <b>Minimal Projection, Heads, and Optimality</b>	<b>TR-5</b> Stephen Stich Ian Ravenscroft \$1.85 <b>What is Folk Psychology?</b>	<b>TR-6</b> Jacob Feldman 10.55 <b>Perceptual Categories and World Regularities</b>
<b>TR-7</b> John McCarthy Alan Prince \$4.30 <b>Generalized Alignment</b>	<b>TR-8</b> Zenon Pylyshyn \$1.85 <b>Some Primitive Mechanisms Underlying Spatial Attention</b>	<b>TR-9</b> Stephen Stich Stephen Laurence \$1.85 <b>Intentionality &amp; Naturalism</b>
<b>TR-10</b> Alan Leslie \$2.30 <b>Pretending and Believing: Issues in the theory of ToMM</b>	<b>TR-11</b> Stephen Stich Shaun Nichols \$2.30 <b>Second Thoughts on Simulation</b>	<b>TR-12</b> Alan Leslie \$2.30 <b>A Theory of Agency</b>
<b>TR-13</b> Jerry Fodor \$1.65 <b>Concepts: A Tutorial Introduction</b>	<b>TR-14</b> Sven Dickinson Dimitri Metaxas \$2.40 <b>Integrating Qualitative and Quantitative Shape Recovery</b>	<b>TR-15</b> David Wilkes Sven Dickinson John Tsotsos \$1.75 <b>A Computational Model of View Degeneracy and its Application to Active Focal Length Control</b>
<b>TR-16</b> Alan Leslie Tim German \$1.95 <b>Knowledge and ability in "theory of mind": One-eyed overview</b>	<b>TR-17</b> Jerry Fodor Ernie Lepore \$1.80 <b>The Red Herring and the Pet Fish; Why Concepts Still Can't be Prototypes</b>	<b>TR-18</b> Suzanne Stevenson \$9.95 <b>A Competitive Attachment Model for Resolving Syntactic</b>
<b>TR-19</b> Jerry Fodor Ernie Lepore \$1.80 <b>What Can't be Evaluated, Can't Be Evaluated; and It Can't Be Supervalued Either</b>	<b>TR-20</b> Ehud Rivlin Sven Dickinson Azriel Rosenfeld \$1.95 <b>Recognition by Functional Parts</b>	<b>TR-21</b> Jacob Feldman \$3.10 <b>Regularity-based Perceptual Grouping</b>

<b>TR-22</b> Zenon Pylyshyn Jacquelyn Burkell \$4.40 <b>Searching through selected subsets of visual displays: A test of the FINST Indexing Hypothesis</b>	<b>TR-23</b> Ilona Kovács \$1.40 <b>Gestalten of today: Early processing of visual contours and surfaces</b>	<b>TR-24</b> Suzanne Stevenson Paola Merlo \$3.05 <b>Lexical Structure and Processing Complexity</b>
<b>TR-25</b> Paola Merlo Suzanne Stevenson \$1.50 <b>Integrating Statistical and Structural Information in a Distributed Architecture for Syntactic Disambiguation</b>	<b>TR-26</b> Sven Dickinson Henrik Christense John Tsotsos \$2.40 Göran Olofsson <b>Active Object Recognition Integrating Attention and Viewpoint Control</b>	<b>TR-27</b> Jerry Fodor Ernie Lepore \$1.60 <b>The Emptiness of the Lexicon: Critical Reflections on J. Pustejovsky's The Generative Lexicon</b>
<b>TR-28</b> Sven Dickinson Dimitri Metaxas Alex Pentland \$2.35 <b>The Role of Model-Based Segmentation in the Recovery of Volumetric Parts from Range Data</b>	<b>TR-29</b> Sven Dickinson Anil K. Jain, Robert Bergevin Roger Munck-Fai Irving Biederman Alex Pentland \$2.20 Jan-Olof Eklundh, <b>Panel Report: The Potential of Geons for Generic 3-D Object Recognition</b>	<b>TR-30</b> Sven Dickinson Dimitri Metaxas \$2.25 <b>Integrating Qualitative and Quantitative Object Representations in the Recovery and Tracking of 3-D Shape</b>
<b>TR-31</b> Alan Prince \$1.55 <b>Gradient Ascent in a Linear Inhibitory Network</b>	<b>TR-32</b> John McCarthy Alan Prince \$8.00 <b>Prosodic Morphology 1986</b>	<b>TR-33</b> Ilona Kovács Ákos Fehér Bela Julesz \$1.00 <b>Medial-point Description of Shape: A Representation for Action Coding and its Psychophysical Correlates</b>
<b>TR-34</b> Jerry Fodor Ernie Lepore \$1.60 <b>Morphemes Matter; The Continuing Case Against Lexical Decomposition</b>	<b>TR-35</b> Thomas V. Papatho Andrei Gorea Akos Feher \$2.70 Tiffany Conway <b>Attention Based Texture Segregation</b>	<b>TR-36</b> Ali Shokoufandeh Ivan Marsic Sven J. Dickinson \$2.70 <b>View Based Object Recognition Using Saliency Maps</b>
<b>TR-37</b> Sven J. Dickinson Dimitri Metaxas \$2.85 <b>Using Aspect Graphs to Control the Recovery and Tracking of Deformable Models</b>	<b>TR-38</b> Zenon Pylyshyn \$3.10 <b>Is Vision Continuous with Cognition? -- The Case for Cognitive Impenetrability of Visual Perception</b>	<b>TR-39</b> Kaleem Siddiqi Ali Shokoufandeh Steven Zucker \$2.40 Sven Dickinson <b>Shock Graphs and Shape Matching</b>
<b>TR-40</b> Brian Scholl Zenon Pylyshyn \$3.50 <b>Tracking Multiple Items Through Occlusion: Clues to Visual Objecthood</b>	<b>TR-41</b> Patrice Tremoulet \$2.50 <b>Individuation and Identification of Physical Objects: Evidence From Human Infants</b>	<b>TR-42</b> Jacob Feldman \$2.50 <b>The Role of Objects in Perceptual Grouping</b>
<b>TR-43</b> Ernest Lepore Kirk Ludwig \$2.50 <b>The Semantics and Pragmatics of Complex Demonstratives</b>	<b>TR-44</b> Matthew Stone 12.20 <b>Modality in Dialogue: Planning, Pragmatics, and Computation</b>	<b>TR-45</b> Shaun Nichols Stephen Stich \$4.85 <b>A Cognitive Theory of Pretense</b>
<b>TR-46</b> Jerry Fodor Ernest Lepore \$1.65 <b>All At Sea in Semantic Space: Churchland on Meaning Similarity</b>	<b>TR-47</b> Jerry Fodor Ernest Lepore \$1.65 <b>Why Compositionality Won't Go Away: Reflections on Horwich's 'Deflationary' Theory</b>	<b>TR-48</b> Doug Decarlo 10.50 <b>Generation, Estimation and Tracking of Faces</b>

<b>TR-49</b> Matthew Stone \$2.50 <b>Reference to Possible Worlds</b>	<b>TR-50</b> Shaun Nichols Stephen Stich \$9.00 <b>Reading One's Own Mind: A Cognitive Theory of Self Awareness</b>	<b>TR-51</b> Doug DeCarlo Dimitris Metaxas \$2.45 <b>Adjusting Model Parameters using Model-Based Optical Flow Residuals</b>
All checks should be made payable to: <b>Rutgers University</b> We cannot accept credit card payments or Eurochecks; American Express Money Orders (in \$US) are acceptable. Prices include 1 <sup>st</sup> class postage in the U.S. For postage to other locations, inquire at: <b>admin@ruccs.rutgers.edu or (732)445-0635.</b>		All Orders for TR's should be addressed to: Susan Cosentino Assistant Director Center for Cognitive Science Psych Bldg Addition, Busch Campus Rutgers University - New Brunswick 152 Frelinghuysen Road Piscataway, NJ 08854-8020