

# Video Motion Capture

Christoph Bregler

Jitendra Malik

Computer Science Division  
University of California, Berkeley  
Berkeley, CA 94720-1776  
bregler@cs.berkeley.edu, malik@cs.berkeley.edu

## Abstract

This paper demonstrates a new vision based motion capture technique that is able to recover high degree-of-freedom articulated human body configurations in complex video sequences. It does not require any markers, body suits, or other devices attached to the subject. The only input needed is a video recording of the person whose motion is to be captured. For visual tracking we introduce the use of a novel mathematical technique, the product of exponential maps and twist motions, and its integration into a differential motion estimation. This results in solving simple linear systems, and enables us to recover robustly the kinematic degrees-of-freedom in noise and complex self occluded configurations. We demonstrate this on several image sequences of people doing articulated full body movements, and visualize the results in re-animating an artificial 3D human model. We are also able to recover and re-animate the famous movements of Eadweard Muybridge's motion studies from the last century. To the best of our knowledge, this is the first computer vision based system that is able to process such challenging footage and recover complex motions with such high accuracy.

### CR Categories:

**Keywords:** Computer Vision, Animation, Motion Capture, Visual Tracking, Twist Kinematics, Exponential Maps, Muybridge

## 1 Introduction

In this paper, we offer a new approach to motion capture based just on ordinary video recording of the actor performing naturally. The approach does not require any markers, body suits or any other devices attached to the body of the actor. The actor can move about wearing his or her regular clothes. This implies that one can use historical footage—motion capture Charlie Chaplin's inimitable walk, for instance. Indeed in this paper we shall go even further back historically and show motion capture results from Muybridge sequences—the first examples of photographically recorded motion [15].

Motion capture occupies an important role in the creation of special effects. Its application to CG character animation has been much more controversial; SIGGRAPH 97 featured a lively panel debate[4] between its proponents and opponents. Our goal in this paper is not to address that debate. Rather we take it as a given that motion capture, like any other technology, can be correctly or incorrectly applied and we are merely extending its possibilities.

Our approach, from a user's point of view, is rather straightforward. The user marks limb segments in an initial frame; if multiple video streams are available from synchronized cameras, then the limb segments are marked in the corresponding initial frames in all of them. The computer program does the rest—tracking the multiple degrees of freedom of the human body configuration from frame to frame.

Attempts to track the human body without special markers go back quite a few years – we review past work in Sec. 2. However in spite of many years of work in computer vision on this problem, it is fair to describe it as not yet solved. There are many reasons why human body tracking is very challenging, compared to tracking other objects such as footballs, robots or cars. These include

1. *High Accuracy Requirements.* Especially in the context of motion capture applications, one desires to record all the degrees of freedom of the configuration of arms, legs, torso, head etc accurately from frame to frame. At playback time, any error will be instantly noticed by a human observer.
2. *Frequent inter-part occlusion* During normal motion, from any camera angle some parts of the body are occluded by other parts of the body
3. *Lack of contrast* Distinguishing the edge of a limb from, say the torso underneath, is made difficult by the fact that typically the texture or color of the shirt is usually the same in both regions.

Our contribution to this problem is the introduction of a novel mathematical technique, the product of exponential maps and twist motions, and its integration into a differential motion estimation scheme. This formalism will be explained fully in Section 3. The advantage of this particular formulation is that it results in the equations that need to be solved to update the kinematic chain parameters from frame to frame being *linear*. Also the only parameters that need to be solved for are the true degrees of freedom and pose parameters—there are no intermediate stages which may be unnecessarily hard. For instance recovering the local affine motion parameters of each and every limb segment separately is harder than the final goal of knowing the configuration of all the joints from frame to frame—the fact that the joints are constrained to move together reduces considerably the number of degrees of freedom. This in turn provides robustness to self-occlusions, loss of contrast, large motions etc.

We applied this technique to several video recordings of walking people and to the famous photo plates of Eadweard Muybridge. We achieved accurate tracking results with high degree-of-freedom full body models and could successfully re-animate the data. The accompanying video shows the tracking results and the naturalness of the animated motion capture data.

Section 2 reviews previous video tracking techniques, section 3 introduces the new motion tracking framework and its mathematical formulation, section 4 details our experiments, and we discuss the results and future directions in section 5.

## 2 Review

The earliest computer vision attempt to recognize human movements was reported by O'Rourke and Badler [16] working on syn-

thetic images using a 3D structure of rigid segments, joints, and constraints between them.

Marker-free visual tracking on video recordings of human bodies goes back to work by Hogg and by Rohr [8, 18]. Both systems are specialized to one degree-of-freedom walking models. Edge and line features are extracted from images and matched to a cylindrical 3D human body model. Higher degree-of-freedom articulated hand configurations are tracked by Regh and Kanade [17], full body configurations by Gravrila and Davis [7], and arm configurations by Kakadiaris and Metaxas [11] and by Goncalves and Perona [5]. All these approaches are demonstrated in constrained environments with high contrast edge boundaries. In most cases this is achieved by uniform backgrounds, and skintight clothing of uniform color. Also, in order to estimate 3D configurations, a camera calibration is needed. Alternatively, Weng et. al demonstrated how to track full bodies with color features [20], and Ju et. al showed motion based tracking of leg configurations [10]. No 3D kinematic chain models were used in the last two cases.

To the best of our knowledge, there is no system reported so far, which would be able to successfully track accurate high-degree-of-freedom human body configurations in the challenging footage that we will demonstrate here.

### 3 Articulated Tracking

There exist a wide range of visual tracking techniques in the literature ranging from edge feature based to region based tracking, and brute-force search methods to differential approaches.

Edge feature based tracking techniques usually require clean data with high contrast object boundaries. Unfortunately on human bodies such features are very noisy. Clothes have many folds. Also if the left and right leg have the same color and they overlap, they are separated only by low contrast boundaries.

Region based techniques can track objects with arbitrary texture. Such techniques attempt to match areas between consecutive frames. For example if the area describes a rigid planar object, a 2D affine deformation of this area has to be found. This requires the estimation of 6 free parameters that describe this deformation (x/y translation, x/y scaling, rotation, and shear). Instead of exhaustively searching over these parameters, differential methods link local intensity changes to parameter changes, and allow for Newton-step like optimizations.

In the following we will introduce a new region based differential technique that is tailored to articulated objects modeled by kinematic chains. We will first review a commonly used motion estimation framework [2, 19], and then show how this can be extended for our task, using the twist and product of exponential formulation [14].

#### 3.1 Preliminaries

Assuming that changes in image intensity are only due to translation of local image intensity, a parametric image motion between consecutive time frames  $t$  and  $t + 1$  can be described by the following equation:

$$I(x + \mathbf{u}_x(x, y, \phi), y + \mathbf{u}_y(x, y, \phi), t + 1) = I(x, y, t) \quad (1)$$

$I(x, y, t)$  is the image intensity. The motion model  $\mathbf{u}(x, y, \phi) = [\mathbf{u}_x(x, y, \phi), \mathbf{u}_y(x, y, \phi)]^T$  describes the pixel displacement dependent on location  $(x, y)$  and model parameters  $\phi$ . For example, a 2D affine motion model with parameters  $\phi = [a_1, a_2, a_3, a_4, d_x, d_y]^T$  is defined as

$$\mathbf{u}(x, y, \phi) = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} d_x \\ d_y \end{bmatrix} \quad (2)$$

The first-order Taylor series expansion of (1) leads to the commonly used gradient formulation [12]:

$$I_t(x, y) + [I_x(x, y), I_y(x, y)] \cdot \mathbf{u}(x, y, \phi) = 0 \quad (3)$$

$I_t(x, y)$  is the temporal image gradient and  $[I_x(x, y), I_y(x, y)]$  is the spatial image gradient at location  $(x, y)$ . Assuming a motion model of  $K$  degrees of freedom (in case of the affine model  $K = 6$ ) and a region of  $N > K$  pixels, we can write an over-constrained set of  $N$  equations. For the case that the motion model is linear (as in the affine case), we can write the set of equations in matrix form (see [2] for details):

$$\mathbf{H} \cdot \phi + \vec{z} = \vec{0} \quad (4)$$

where  $\mathbf{H} \in \mathbb{R}^{N \times K}$ , and  $\vec{z} \in \mathbb{R}^N$ . The least squares solution to (3) is:

$$\phi = -(\mathbf{H}^T \cdot \mathbf{H})^{-1} \cdot \mathbf{H}^T \vec{z} \quad (5)$$

Because (4) is the first-order Taylor series linearization of (1), we linearize around the new solution and iterate. This is done by warping the image  $I(t + 1)$  using the motion model parameters  $\phi$  found by (5). Based on the re-warped image we compute the new image gradients (3). Repeating this process is equivalent to a Newton-Raphson style minimization.

A convenient representation of the shape of an image region is a probability mask  $w(x, y) \in [0, 1]$ .  $w(x, y) = 1$  declares that pixel  $(x, y)$  is part of the region. Equation (5) can be modified, such that it weights the contribution of pixel location  $(x, y)$  according to  $w(x, y)$ :

$$\phi = -((\mathbf{W} \cdot \mathbf{H})^T \cdot \mathbf{H})^{-1} \cdot (\mathbf{W} \cdot \mathbf{H})^T \vec{z} \quad (6)$$

$\mathbf{W}$  is an  $N \times N$  diagonal matrix, with  $\mathbf{W}(i, i) = w(x_i, y_i)$ . We assume for now that we know the exact shape of the region. For example, if we want to estimate the motion parameters for a human body part, we supply a weight matrix  $\mathbf{W}$  that defines the image support map of that specific body part, and run this estimation technique for several iterations. Section 3.4 describes how we can estimate the shape of the support maps as well.

Tracking over multiple frames can be achieved by applying this optimization technique successively over the complete image sequence.

#### 3.2 Twists and the Product of Exponential Formula

In the following we develop a motion model  $\mathbf{u}(x, y, \phi)$  for a 3D kinematic chain under scaled orthographic projection and show how these domain constraints can be incorporated into one linear system similar to (6).  $\phi$  will represent the 3D pose and angle configuration of such a kinematic chain and can be tracked in the same fashion as already outlined for simpler motion models.

##### 3.2.1 3D pose

The pose of an object relative to the camera frame can be represented as a rigid body transformation in  $\mathbb{R}^3$  using homogeneous coordinates (we will use the notation from [14]):

$$q_c = \mathbf{G} \cdot q_o \quad \text{with} \quad \mathbf{G} = \begin{bmatrix} r_{1,1} & r_{1,2} & r_{1,3} & d_x \\ r_{2,1} & r_{2,2} & r_{2,3} & d_y \\ r_{3,1} & r_{3,2} & r_{3,3} & d_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (7)$$

$q_o = [x_o, y_o, z_o, 1]^T$  is a point in the object frame and  $q_c = [x_c, y_c, z_c, 1]^T$  is the corresponding point in the camera frame. Using scaled orthographic projection with scale  $s$ , the point  $q_c$  in the camera frame gets projected into the image point  $[x_{im}, y_{im}]^T = s \cdot [x_c, y_c]^T$ .

The 3D translation  $[d_x, d_y, d_z]^T$  can be arbitrary, but the rotation matrix:

$$\mathbf{R} = \begin{bmatrix} r_{1,1} & r_{1,2} & r_{1,3} \\ r_{2,1} & r_{2,2} & r_{2,3} \\ r_{3,1} & r_{3,2} & r_{3,3} \end{bmatrix} \in SO(3) \quad (8)$$

has only 3 degrees of freedom. Therefore the rigid body transformation  $\mathbf{G} \in SE(3)$  has a total of 6 degrees of freedom.

Our goal is to find a model of the image motion that is parameterized by 6 degrees of freedom for the 3D rigid motion and the scale factor  $s$  for scaled orthographic projection. *Euler angles* are commonly used to constrain the rotation matrix to  $SO(3)$ , but they suffer from singularities and don't lead to a simple formulation in the optimization procedure (for example [1] propose a 3D ellipsoidal tracker based on Euler angles). In contrast, the *twist* representation provides a more elegant solution [14] and leads to a very simple linear representation of the motion model. It is based on the observation that every rigid motion can be represented as a rotation around a 3D axis and a translation along this axis. A twist  $\xi$  has two representations: (a) a 6D vector, or (b) a  $4 \times 4$  matrix with the upper  $3 \times 3$  component as a skew-symmetric matrix:

$$\xi = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \omega_x \\ \omega_y \\ \omega_z \end{bmatrix} \quad \text{or} \quad \hat{\xi} = \begin{bmatrix} 0 & -\omega_z & \omega_y & v_1 \\ \omega_z & 0 & -\omega_x & v_2 \\ -\omega_y & \omega_x & 0 & v_3 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (9)$$

$\omega$  is a 3D unit vector that points in the direction of the rotation axis. The amount of rotation is specified with a scalar angle  $\theta$  that is multiplied by the twist:  $\xi\theta$ . The  $v$  component determines the location of the rotation axis and the amount of translation along this axis. See [14] for a detailed geometric interpretation. For simplicity, we drop the constraint that  $\omega$  is unit, and discard the  $\theta$  coefficient. Therefore  $\xi \in \mathbb{R}^6$ .

It can be shown [14] that for any arbitrary  $\mathbf{G} \in SE(3)$  there exists a  $\xi \in \mathbb{R}^6$  twist representation.

A twist can be converted into the  $\mathbf{G}$  representation with following exponential map:

$$\begin{aligned} \mathbf{G} &= \begin{bmatrix} r_{1,1} & r_{1,2} & r_{1,3} & d_x \\ r_{2,1} & r_{2,2} & r_{2,3} & d_y \\ r_{3,1} & r_{3,2} & r_{3,3} & d_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \\ &= e^{\hat{\xi}} = \mathbf{I} + \hat{\xi} + \frac{(\hat{\xi})^2}{2!} + \frac{(\hat{\xi})^3}{3!} + \dots \end{aligned} \quad (10)$$

### 3.2.2 Twist motion model

At this point we would like to track the 3D pose of a rigid object under scaled orthographic projection. We will extend this formulation in the next section to a kinematic chain representation. The pose of an object is defined as  $[s, \xi]^T = [s, v_1, v_2, v_3, \omega_x, \omega_y, \omega_z]^T$ . A point  $q_o$  in the object frame is projected to the image location  $(x_{im}, y_{im})$  with:

$$\begin{bmatrix} x_{im} \\ y_{im} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \cdot s \cdot e^{\hat{\xi}} \cdot q_o \quad (11)$$

The image motion of point  $(x_{im}, y_{im})$  from time  $t$  to time  $t + 1$  is:

$$\begin{aligned} \begin{bmatrix} u_x \\ u_y \end{bmatrix} &= \begin{bmatrix} x_{im}(t+1) - x_{im}(t) \\ y_{im}(t+1) - y_{im}(t) \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \cdot \left( s(t+1) \cdot e^{\hat{\xi}(t+1)} \cdot q_o - s(t) \cdot e^{\hat{\xi}(t)} \cdot q_o \right) \\ &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \cdot \left( (1 + s') \cdot e^{\hat{\xi}'} - \mathbf{I} \right) \cdot s(t) q_c \end{aligned} \quad (12)$$

$$\begin{aligned} \text{with } \hat{\xi}(t+1) &= \hat{\xi}(t) + \hat{\xi}' \\ s(t+1) &= s(t) \cdot (1 + s') \end{aligned}$$

Using the first order Taylor expansion from (10) we can approximate:

$$(1 + s') \cdot e^{\hat{\xi}} \approx (1 + s') \cdot \mathbf{I} + (1 + s') \cdot \hat{\xi} \quad (13)$$

and can rewrite (12) as:

$$\begin{bmatrix} u_x \\ u_y \end{bmatrix} = \begin{bmatrix} s' & -\omega'_z & \omega'_y & v'_1 \\ \omega'_z & s' & -\omega'_x & v'_2 \end{bmatrix} \cdot q_c \quad (14)$$

with

$$\omega(t+1) = \omega(t) + \frac{1}{1 + s'} \cdot \omega'$$

$$v(t+1) = v(t) + \frac{1}{1 + s'} \cdot v'$$

$\phi = [s', v'_1, v'_2, \omega'_x, \omega'_y, \omega'_z]^T$  codes the relative scale and twist motion from time  $t$  to  $t + 1$ . Note that (14) does not include  $v'_3$ . Translation in the  $Z$  direction of the camera frame is not measurable under scaled orthographic projection.

Equation (14) describes the image motion of a point  $(x_i, y_i)$  in terms of the motion parameters  $\phi$  and the corresponding 3D point  $q_c(i)$  in the camera frame. The 3D point  $q_c(i)$  is computed by intersecting the camera ray of the image point  $(x_i, y_i)$  with the 3D model. In this paper we assume that the body segments can be approximated by ellipsoidal 3D blobs. Therefore  $q_c$  is the solution of a quadratic equation. This computation has to be done only once for each new image. It is outside the Newton-Raphson iterations. It could be replaced by more complex models and rendering algorithms.

Inserting (14) into (3) leads to:

$$\begin{aligned} I_t + I_x \cdot [s', -\omega'_z, \omega'_y, v'_1] \cdot q_c + I_y \cdot [\omega'_z, s', -\omega'_x, v'_2] \cdot q_c &= 0 \\ \Leftrightarrow I_t(i) + H_i \cdot [s, v'_1, v'_2, \omega'_x, \omega'_y, \omega'_z]^T &= 0 \end{aligned} \quad (15)$$

$$\text{with } I_t := I_t(x_i, y_i), I_x := I_x(x_i, y_i), I_y := I_y(x_i, y_i)$$

For  $N$  pixel positions we have  $N$  equations of the form (15). This can be written in matrix form:

$$\mathbf{H} \cdot \phi + \mathbf{z} = 0 \quad (16)$$

with

$$\mathbf{H} = \begin{bmatrix} H_1 \\ H_2 \\ \vdots \\ H_N \end{bmatrix} \quad \text{and} \quad \mathbf{z} = \begin{bmatrix} I_t(x_1, y_1) \\ I_t(x_2, y_2) \\ \vdots \\ I_t(x_N, y_N) \end{bmatrix}$$

Finding the least-squares solution (3D twist motion  $\phi$ ) for this equation is done using (6).

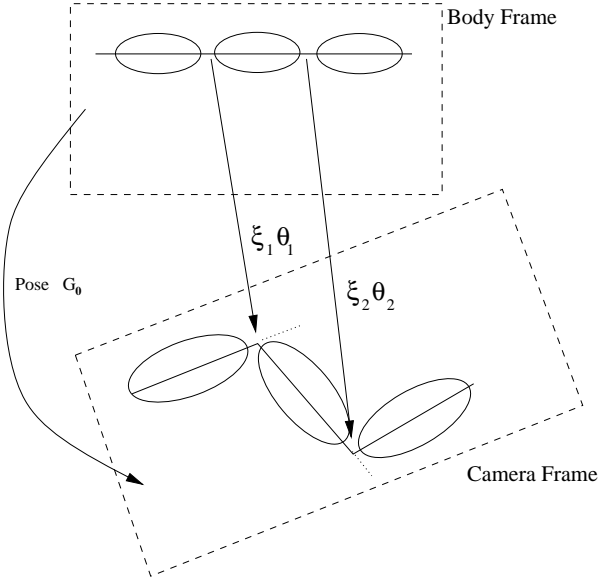


Figure 1: Kinematic chain defined by twists

### 3.2.3 Kinematic chain as a Product of Exponentials

So far we have parameterized the 3D pose and motion of a body segment by the 6 parameters of a twist  $\xi$ . Points on this body segment in a canonical object frame are transformed into a camera frame by the mapping  $\mathbf{G}_0 = e^{\xi}$ . Assume that a second body segment is attached to the first segment with a joint. The joint can be defined by an axis of rotation in the object frame. We define this rotation axis in the object frame by a 3D unit vector  $\omega_1$  along the axis, and a point  $q_1$  on the axis (figure 1). This is a so called revolute joint, and can be modeled by a twist ([14]):

$$\xi_1 = \begin{bmatrix} -\omega_1 \times q_1 \\ \omega_1 \end{bmatrix} \quad (17)$$

A rotation of angle  $\theta_1$  around this axis can be written as:

$$\mathbf{g}_1 = e^{\xi_1 \cdot \theta_1} \quad (18)$$

$$(19)$$

The global mapping from object frame points on the first body segment into the camera frame is described by the following product:

$$\begin{aligned} \mathbf{g}(\theta_1) &= \mathbf{G}_0 \cdot e^{\xi_1 \cdot \theta_1} \\ q_c &= \mathbf{g}(\theta_1) \cdot q_o \end{aligned} \quad (20)$$

If we have a chain of  $K+1$  segments linked with  $K$  joints (kinematic chain) and describe each joint by a twist  $\xi_k$ , a point on segment  $k$  is mapped from the object frame into the camera frame dependent on  $\mathbf{G}_0$  and angles  $\theta_1, \theta_2, \dots, \theta_k$ :

$$\mathbf{g}_k(\theta_1, \theta_2, \dots, \theta_k) = \mathbf{G}_0 \cdot e^{\xi_1 \cdot \theta_1} \cdot e^{\xi_2 \cdot \theta_2} \cdot \dots \cdot e^{\xi_k \cdot \theta_k} \quad (21)$$

This is called the **product of exponential maps** for kinematic chains.

The velocity of a segment  $k$  can be described with a twist  $V_k$  that is a linear combination of twists  $\xi'_1, \xi'_2, \dots, \xi'_k$  and the angular velocities  $\dot{\theta}_1, \dot{\theta}_2, \dots, \dot{\theta}_k$  (see [14] for the derivations):

$$\begin{aligned} V_k &= \xi'_1 \cdot \dot{\theta}_1 + \xi'_2 \cdot \dot{\theta}_2 + \dots + \xi'_k \cdot \dot{\theta}_k \\ \xi'_k &= \mathbf{Ad}_{e^{\xi_1 \theta_1} \dots e^{\xi_{k-1} \theta_{k-1}}} \xi_k \end{aligned} \quad (22)$$

$\mathbf{Ad}_g$  is the adjoint transformation associated with  $g$ .<sup>1</sup>

Given a point  $q_c$  on the  $k$ 'th segment of a kinematic chain, its motion vector in the image is related to the angular velocities by:

$$\begin{bmatrix} \mathbf{u}_x \\ \mathbf{u}_y \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \cdot [\hat{\xi}'_1 \cdot \dot{\theta}_1 + \hat{\xi}'_2 \cdot \dot{\theta}_2 + \dots + \hat{\xi}'_k \cdot \dot{\theta}_k] \cdot q_c \quad (23)$$

Recall (15) relates the image motion of a point  $q_c$  to changes in pose  $\mathbf{G}_0$ . We combine (15) and (23) to relate the image motion to the combined vector of pose change and angular change  $\Phi = [s', v'_1, v'_2, \omega'_x, \omega'_y, \omega'_z, \dot{\phi}_1, \dot{\phi}_2, \dots, \dot{\phi}_K]^T$ :

$$I_t + H_i \cdot [s, v'_1, v'_2, \omega'_x, \omega'_y, \omega'_z]^T + J_i \cdot [\dot{\theta}_1, \dot{\theta}_2, \dots, \dot{\theta}_K]^T = 0 \quad (24)$$

$$[\mathbf{H}, \mathbf{J}] \cdot \Phi + \vec{z} = 0 \quad (25)$$

with

$$\mathbf{J} = \begin{bmatrix} J_1 \\ J_2 \\ \vdots \\ J_N \end{bmatrix} \quad \text{and } \mathbf{H}, \vec{z} \text{ as before}$$

$$J_i = [J_{i1}, J_{i2}, \dots, J_{iK}]$$

$$J_{ik} = \begin{cases} [I_x, I_y] \cdot \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \cdot \hat{\xi}_k \cdot q_c \\ 0 \end{cases} \quad \text{if pixel } i \text{ is on a segment that} \\ \quad \quad \quad \text{is not affected by joint } \xi_k$$

The least squares solution to (25) is:

$$\Phi = -([\mathbf{H}, \mathbf{J}]^T \cdot [\mathbf{H}, \mathbf{J}])^{-1} \cdot [\mathbf{H}, \mathbf{J}]^T \cdot \vec{z} \quad (26)$$

$\Phi$  is the new estimate of the pose and angular change between two consecutive images. As outlined earlier, this solution is based on the assumption that the local image intensity variations can be approximated by the first-order Taylor expansion (3). We linearize around this new solution and iterate. This is done in warping the image  $I(t+1)$  using the solution  $\Phi$ . Based on the re-warped image we compute the new image gradients. Repeating this process of warping and solving (26) is equivalent to a Newton-Raphson style minimization.

### 3.3 Multiple Camera Views

In cases where we have access to multiple synchronized cameras, we can couple the different views in one equation system. Let's assume we have  $C$  different camera views at the same time. View  $c$  corresponds to following equation system (from (25)):

$$[\mathbf{H}_c, \mathbf{J}_c] \cdot \begin{bmatrix} \Omega_c \\ \dot{\phi}_1 \\ \dot{\phi}_2 \\ \vdots \\ \dot{\phi}_K \end{bmatrix} + \vec{z}_c = 0 \quad (27)$$

$\Omega_c = [s'_{c,1}, v'_{c,1}, v'_{c,2}, \omega'_{c,x}, \omega'_{c,y}, \omega'_{c,z}]^T$  describes the pose seen from view  $c$ . All views share the same angular parameters, because

<sup>1</sup>  $\mathbf{Ad}_g = \begin{bmatrix} R & \hat{p} \cdot R \\ 0 & R \end{bmatrix}$ , and  $g = \begin{bmatrix} R & p \\ 000 & 1 \end{bmatrix}$

the cameras are triggered at the same time. We can simply combine all  $C$  equation systems into one large equation system:

$$\begin{bmatrix} \mathbf{H}_1 & \mathbf{0} & \dots & \mathbf{0} & \mathbf{J}_1 \\ \mathbf{0} & \mathbf{H}_2 & \dots & \mathbf{0} & \mathbf{J}_2 \\ \dots & \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{H}_C & \mathbf{J}_C \end{bmatrix} \cdot \begin{bmatrix} \Omega_1 \\ \Omega_2 \\ \dots \\ \Omega_C \\ \phi_1 \\ \phi_2 \\ \dots \\ \phi_K \end{bmatrix} + \begin{bmatrix} \vec{z}_1 \\ \vec{z}_2 \\ \dots \\ \vec{z}_C \end{bmatrix} = \mathbf{0} \quad (28)$$

Operating with multiple views has three main advantages. The estimation of the angular parameters is more robust: (1) the number of measurements and therefore the number of equations increases with the number of views, (2) some angular configurations might be close to a singular pose in one view, whereas they can be estimated in a orthogonal view much better. (3) With more camera views, the chance decreases that one body part is occluded in all views.

### 3.4 Adaptive Support Maps using EM

As in (3), the update can be constrained to estimate the motion only in a weighted support map  $\mathbf{W}_k$  for each segment  $k$  using:

$$\Phi = - \left( (\mathbf{W}_k \cdot [\mathbf{H}, \mathbf{J}])^T \cdot [\mathbf{H}, \mathbf{J}] \right)^{-1} \cdot (\mathbf{W}_k \cdot [\mathbf{H}, \mathbf{J}])^T \vec{z} \quad (29)$$

We approximate the shape of the body segments as ellipsoids, and can compute the support map as the projection of the ellipsoids into the image. Such a support map usually covers a larger region, including pixels from the environment. That distracts the exact motion measurement. Robust statistics would be one solution to this problem [3]. Another solution is an EM-based layered representation [6, 9]. It is beyond the scope of this paper to describe this method in detail, but we would like to outline the method briefly: We start with an initial guess of the support map (ellipsoidal projection in this case). Given the initial  $\mathbf{W}_k$ , we compute the motion estimate  $\Phi$  (M-step). Given such a  $\Phi$  we can compute for each pixel location the probability that it complies with the motion model defined by  $\Phi$ . We do this for each blob and the background (dominant motion) and normalize the sum of all probabilities per pixel location to 1. This results in new  $\mathbf{W}_k$  maps that are better “tuned” to the real shape of the body segment. In this paper we repeat the EM iteration only once.

### 3.5 Tracking Recipe

We summarize the algorithm for tracking the pose and angles of a kinematic chain in an image sequence:

- **Input:**  $I(t), I(t+1), \mathbf{G}_0(t), \theta_1(t), \theta_2(t), \dots, \theta_K(t)$   
(Two images and the pose and angles for the first image).
- **Output:**  $\mathbf{G}_0(t+1), \theta_1(t+1), \theta_2(t+1), \dots, \theta_K(t+1)$ .  
(Pose and angles for second image).

1. Compute for each image location  $(x_i, y_i)$  in  $I(t)$  the 3D point  $q_c(i)$  (using ellipsoids or more complex models and rendering algorithm).
2. Compute for each body segment the support map  $W_k$ .
3. Set  $\mathbf{G}_0(t+1) := \mathbf{G}_0(t), \forall k: \theta_k(t+1) := \theta_k(t)$ .

4. Iterate:

- (a) Compute spatiotemporal image gradients:  $I_t, I_x, I_y$ .
- (b) Estimate  $\Phi$  using (29)
- (c) Update  $G_0(t+1) := G_0(t+1) \cdot (1+s') \cdot e^{\frac{\xi^t}{1+s'}}$
- (d)  $\forall k$  Update  $\theta_k(t+1) := \theta_k(t+1) + \dot{\theta}_k$ .
- (e)  $\forall k$  Warp the region inside  $W_k$  of  $I(t+1)$  by  $\mathbf{G}_0(t+1) \cdot g_k(t+1) \cdot (\mathbf{G}^t \cdot g_k(t))^{-1}$ .

## 3.6 Initialization

The visual tracking is based on an initialized first frame. We have to know the initial pose and the initial angular configuration. If more than one view is available, all views for the first time step have to be known. A user clicks on the 2D joint locations in all views at the first time step. Given that, the 3D pose and the image projection of the matching angular configuration is found in minimizing the sum of squared differences between the projected model joint locations and the user supplied model joint locations. The optimization is done over the poses, angles, and body dimensions. Example body dimensions are “upper-leg-length”, “lower-leg-length”, or “shoulder-width”. The dimensions and angles have to be the same in all views, but the pose can be different. Symmetry constraints, that the left and right body lengths are the same, are enforced as well. Minimizing only over angles, or only over model dimensions results in linear equations similar to what we have shown so far. Unfortunately the global minimization criteria over all parameters is a tri-linear equation system, that cannot be easily solved by simple matrix inversions. There are several possible techniques for minimizing such functions. We achieved good results with a Quasi-Newton method and a mixed quadratic and cubic line search procedure.

## 4 Results

We applied this technique to video recordings in our lab and to photo-plate sequence of Eadweard Muybridge’s motion studies.

### 4.1 Single camera recordings

Our lab video recordings were done with a single camera. Therefore the 3D pose and some parts of the body can not be estimated completely. Figure 2 shows one example sequences of a person walking in a frontoparallel plane. We defined a 6 DOF kinematic structure: One blob for the body trunk, three blobs for the frontal leg and foot, connected with a hip joint, knee joint, and ankle joint, and two blobs for the arm connected with a shoulder and elbow joint. All joints have an axis orientation parallel to the  $Z$ -axis in the camera frame. The head blob was connected with one joint to the body trunk. The first image in figure 2 shows the initial blob support maps.

After the hand-initialization we applied the motion tracker to a sequence of 53 image frames. We could successfully track all body parts in this video sequence (see video). The video shows that the appearance of the upper leg changes significantly due to moving folds on the subject’s jeans. The lower leg appearance does not change to the same extent. The constraints were able to enforce compatible motion vectors for the upper leg, based on more reliable measurements on the lower leg.

We can compare the estimated angular configurations with motion capture data reported in the literature. Murray, Brought, and

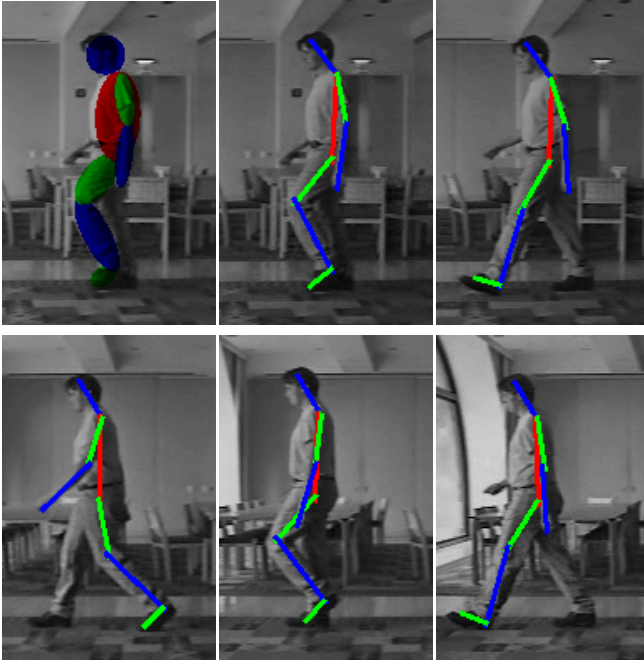


Figure 2: Example configurations of the estimated kinematic structure. First image shows the support maps of the initial configuration. In subsequent images the white lines show blob axes. The joint is the position on the intersection of two axes.

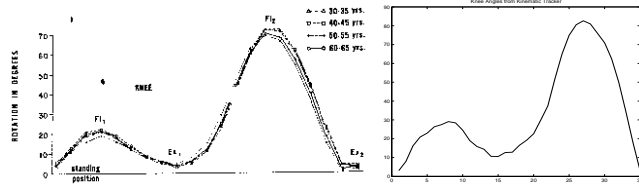


Figure 3: Comparison of a) data from [Murray et al] (left) and b) our motion tracker (right).

Kory published [13] such measurements for the hip, knee, and ankle joints. We compared our motion tracker measurements with the published curves and found good agreement. Figure 4.1a shows the curves for the knee and ankle reported in [13], and figure 4.1b shows our measurements.

We also experimented with a walking sequence of a subject seen from an oblique view with a similar kinematic model. As seen in figure 4, we tracked the angular configurations and the pose successfully over the complete sequence of 45 image frames. Because we use a scaled orthographic projection model, the perspective effects of the person walking closer to the camera had to be compensated by different scales. The tracking algorithm could successfully estimate the scale changes.

## 4.2 Digital Muybridge

The final set of experiments was done on historic footage recorded by Eadward Muybridge in 1884. His methods are of independent interest, as they predate motion pictures. Muybridge had his models walk in an open shed. Parallel to the shed was a fixed battery of 24 cameras. Two portable batteries of 12 cameras each were positioned at both ends of the shed, either at an angle of 90 deg relative to the shed or an angle of 60 deg. Three photographs were take

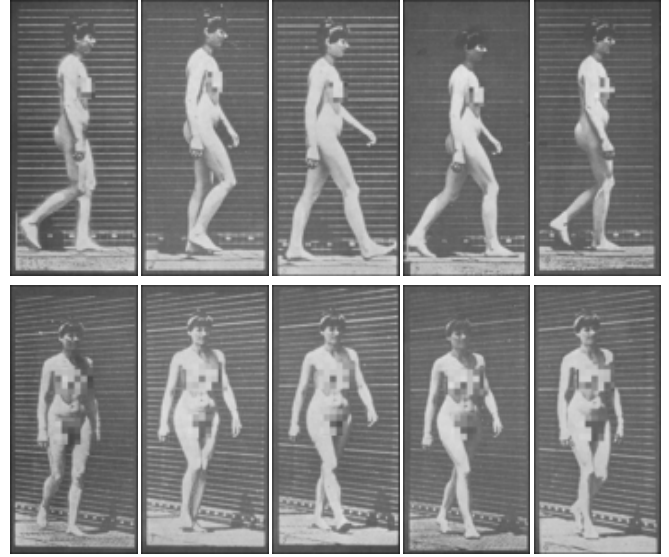


Figure 5: Eadward Muybridge, The Human Figure in Motion, Plate 97: Woman Walking. The first 5 frames show part of a walk cycle from one example view, and the second 5 frames show the same time steps from a different view

simultaneously, one from each battery. The effective ‘framerate’ of his technique is about two times lower than current video frame rates; a fact which makes tracking a harder problem.. It is to our advantage that he took for each time step three pictures from different viewpoints.

Figure 4.2 and figure 4.2 shows example photo plates. We could initialize the 3D pose by labeling all three views of the first frame and running the minimization procedure over the body dimensions and poses. Figure 4.2 shows one example initialization. Every body segment was visible in at least one of the three camera views, therefore we could track the left and the right side of the person. We applied this technique to a walking woman and a walking man. For the walking woman we had 10 time steps available that contained 60 % of a full walk cycle (figure 4.2). For this set of experiments we extended our kinematic model to 19 DOFs. The two hip joints, the two shoulder joints, and the neck joint, were modeled by 3 DOFs. The two knee joints and two elbow joints were modeled just by one rotation axis. Figure 4.2 shows the tracking results with the model overlaid. As you see, we could successfully track the complete sequence. To animate the tracking results we mirrored the left and right side angles to produce the remaining frames of a complete walk cycle. We animated the 3D motion capture data with a stick figure model and a volumetric model (figure 10), and it looks very natural. The video shows some of the tracking and animation sequences from several novel camera views, replicating the walk cycle performed over a century ago on the grounds of University of Pennsylvania.

For the visualization of the walking man sequence, we did not apply the mirroring, because he was carrying a boulder on his shoulder. This made the walk asymmetric. We re-animated the original tracked motion (figure 4.2) capture data for the man, and it also looked very natural.

Given the successful application of our tracking technique to multi-view data, we are planning to record with higher frame-rates our own multi-view video footage. We also plan to record a wider range of gestures.





Figure 4: Example configurations of the estimated kinematic structure of a person seen from an oblique view.

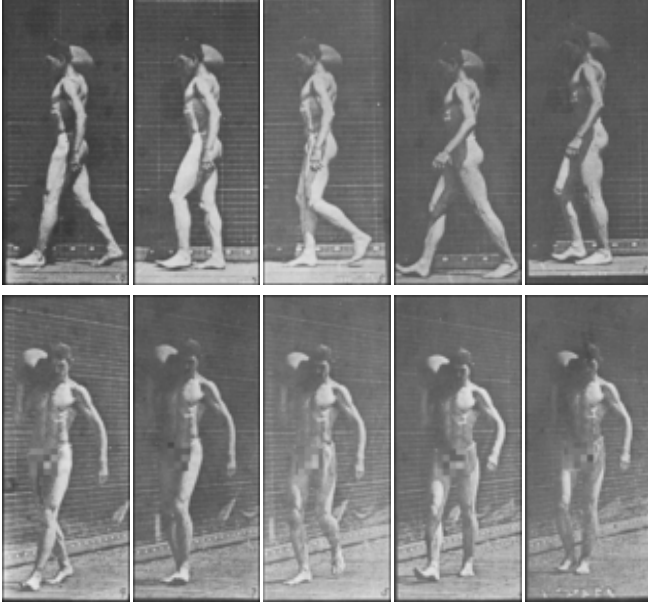


Figure 6: Eadweard Muybridge, *The Human Figure in Motion*, Plate 7: Man walking and carrying 75-LB boulder on shoulder. The first 5 frames show part of a walk cycle from one example view, and the second 5 frames show the same time steps from a different view

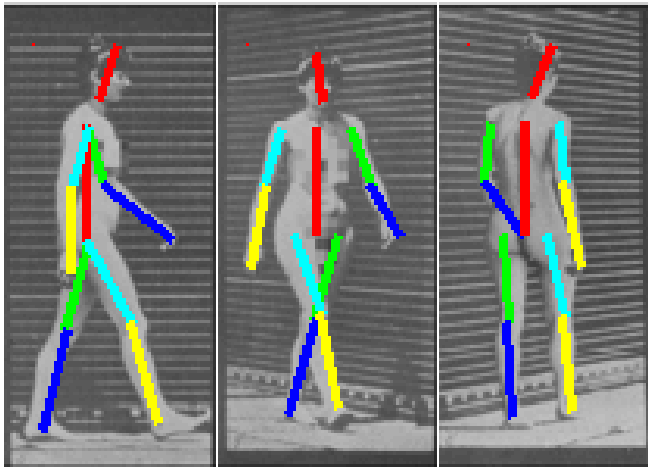


Figure 7: Initialization of Muybridge's Woman Walking: This visualizes the initial angular configuration projected to 3 example views.

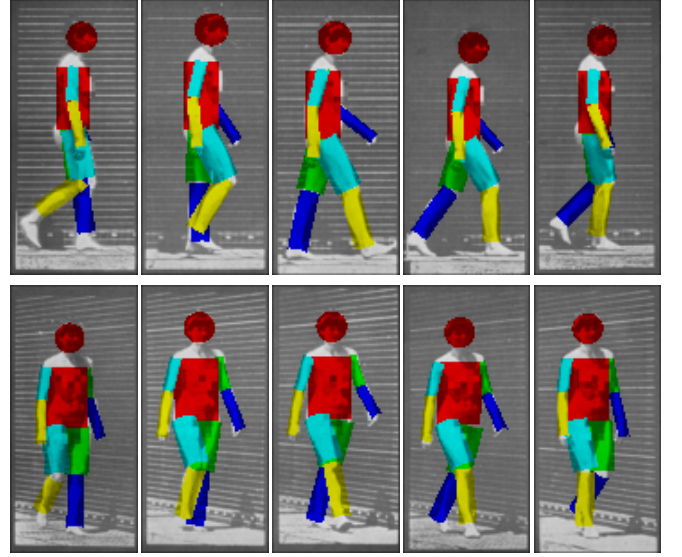


Figure 8: Muybridge's Woman Walking: Motion Capture results. This shows the tracked angular configurations and its volumetric model projected to 2 example views.

## 5 Conclusion

In this paper, we have developed and demonstrated a new technique for video motion capture. The approach does not require any markers, body suits or any other devices attached to the body of the actor. The actor can move about wearing his or her regular clothes. We demonstrated results on video recordings of people walking both in frontoparallel and oblique views, as well as on the classic Muybridge photographic sequences recorded more than a century ago.

Visually tracking human motion at the level of individual joints is a very challenging problem. Our results are due, in large measure, to the introduction of a novel mathematical technique, the product of exponential maps and twist motions, and its integration into a differential motion estimation scheme. The advantage of this particular formulation is that it results in the equations that need to be solved to update the kinematic chain parameters from frame to frame being linear, and that it is not necessary to solve for any redundant or unnecessary variables.

Future work will concentrate on dealing with very large motions, as may happen, for instance, in videotapes of high speed running. The approach developed in this paper is a differential method, and therefore may be expected to fail when the motion from frame-to-frame is very large. We propose to augment the technique by the use of an initial coarse search stage. Given a close enough starting value, the differential method will converge correctly.

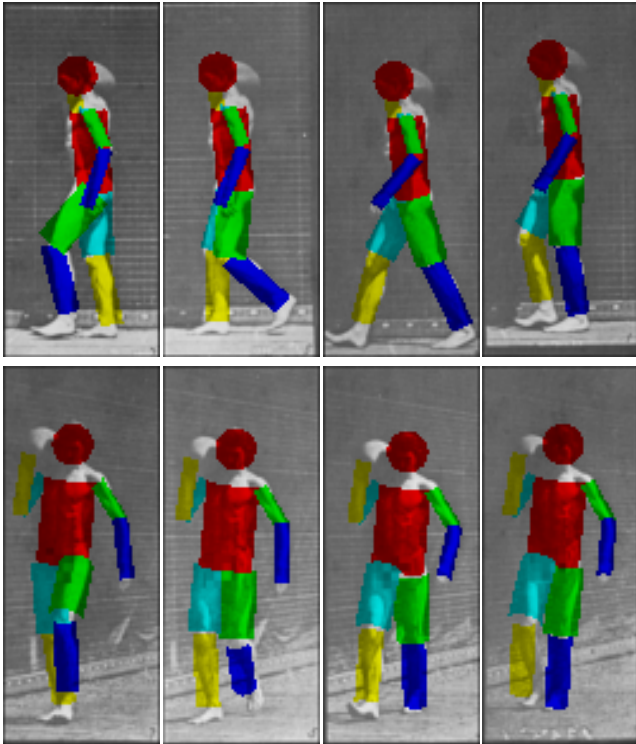


Figure 9: Muybridge's Man Walking: Motion Capture results. This shows the tracked angular configurations and its volumetric model projected to 2 example views.

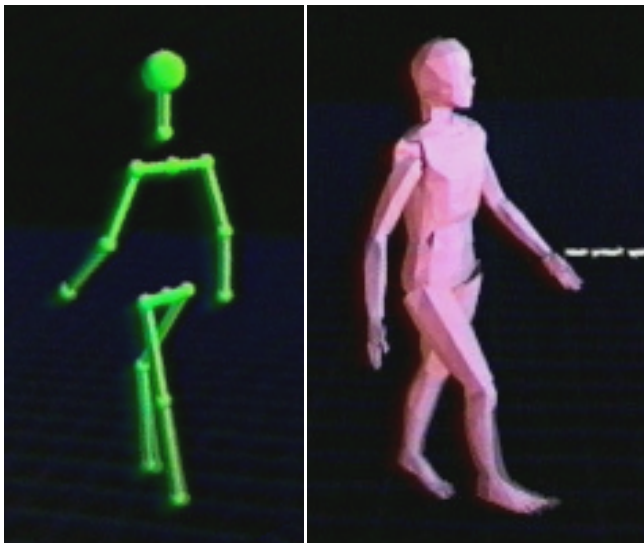


Figure 10: Computer models used for the animation of the Muybridge motion capture. Please check out the video to see the quality of the animation.

## Acknowledgements

We would like to thank Charles Ying for creating the Open-GL animations and video editing, Shankar Sastry, Lara Crawford, Jerry Feldman, John Canny, and Jianbo Shi for fruitful discussions, Chad Carson for helping to write this document, and Interval Research Corp, and the California State MICRO program for supporting this research.

## References

- [1] S. Basu, I.A. Essa, and A.P. Pentland. Motion regularization for model-based head tracking. In *International Conference on Pattern Recognition*, 1996.
- [2] J.R. Bergen, P. Anandan, K.J. Hanna, and R. Hingorani. Hierarchical model-based motion estimation. In *ECCV*, pages 237–252, 1992.
- [3] M.J. Black and P. Anandan. The robust estimation of multiple motions: Parametric and piecewise-smooth flow fields. *Computer Vision and Image Understanding*, 63(1):75–104, Jan 1996.
- [4] G. Cameron, A. Bustanoby, K. Cope, S. Greenberg, C. Hayes, and O. Ozoux. Panel on motion capture and cg character animation. *SIGGRAPH 97*, pages 442–445, 1997.
- [5] L. Concalves, E.D. Bernardo, E. Ursella, and P. Perona. Monocular tracking of the human arm in 3d. In *Proc. Int. Conf. Computer Vision*, 1995.
- [6] A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society B*, 39, 1977.
- [7] D.M. Gavrila and L.S. Davis. Towards 3-d model-based tracking and recognition of human movement: a multi-view approach. In *Proc. of the Int. Workshop on Automatic Face- and Gesture-Recognition, Zurich, 1995*, 1995.
- [8] D. Hogg. A program to see a walking person. *Image Vision Computing*, 5(20), 1983.
- [9] A. Jepson and M.J. Black. Mixture models for optical flow computation. In *Proc. IEEE Conf. Computer Vision Pattern Recognition*, pages 760–761, New York, 1993.
- [10] S.X. Ju, M.J. Black, and Y. Yacoob. Cardboard people: A parameterized model of articulated motion. In *2nd Int. Conf. on Automatic Face- and Gesture-Recognition, Killington, Vermont*, pages 38–44, 1996.
- [11] I.A. Kakadiaris and D. Metaxas. Model-based estimation of 3d human motion with occlusion based on active multi-viewpoint selection. In *CVPR*, 1996.
- [12] B.D. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. *Proc. 7th Int. Joint Conf. on Art. Intell.*, 1981.
- [13] M.P. Murray, A.B. Drought, and R.C. Kory. Walking patterns of normal men. *Journal of Bone and Joint Surgery*, 46-A(2):335–360, March 1964.
- [14] R.M. Murray, Z. Li, and S.S. Sastry. *A Mathematical Introduction to Robotic Manipulation*. CRC Press, 1994.



- [15] Eadweard Muybridge. *The Human Figure In Motion*. Various Publishers, latest edition by Dover Publications, 1901.
- [16] J. O'Rourke and N. I. Badler. Model-based image analysis of human motion using constraint propagation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2(6):522–536, November 1980.
- [17] J.M. Rehg and T. Kanade. Model-based tracking of self-occluding articulated objects. In *Proc. Int. Conf. Computer Vision*, 1995.
- [18] K. Rohr. Incremental recognition of pedestrians from image sequences. In *Proc. IEEE Comput. Soc. Conf. Comput. Vision and Pattern Recogn.*, pages 8–13, New York City, June, 1993.
- [19] J. Shi and C. Tomasi. Good features to track. In *CVPR*, 1994.
- [20] C. Wren, A. Azarbayejani, T. Darrell, and A. Pentland. Pfinder: Real-time tracking of the human body. In *SPIE Conference on Integration Issues in Large Commercial Media Delivery Systems*, volume 2615, 1995.