

The Computational Perception of Scene Dynamics

RICHARD MANN,* ALLAN JEPSON,*^{†,1} AND JEFFREY MARK SISKIND[‡]

**Department of Computer Science, University of Toronto, 6 Kings College Road, Toronto, Ontario, M5S 3H5 Canada; †Canadian Institute for Advanced Research; ‡Department of Electrical Engineering and Computer Science, University of Vermont, Burlington, Vermont 05405*

Received November 17, 1995; accepted November 20, 1996

Understanding observations of interacting objects requires one to reason about qualitative scene dynamics. For example, on observing a hand lifting a can, we may infer that an “active” hand is applying an upwards force (by grasping) to lift a “passive” can. We present an implemented computational theory that derives such dynamic descriptions directly from camera input. Our approach is based on an analysis of the Newtonian mechanics of a simplified scene model. Interpretations are expressed in terms of assertions about the kinematic and dynamic properties of the scene. The feasibility of interpretations relative to Newtonian mechanics is determined by a reduction to linear programming. Finally, to select plausible interpretations, multiple feasible solutions are compared using a preference hierarchy. We provide computational examples to demonstrate that our model is sufficiently rich to describe a wide variety of image sequences. © 1997 Academic Press

1. INTRODUCTION

Both AI and psychology researchers have argued for the need to represent “causal” information about the world in order to make inferences. In particular, understanding motion sequences requires the observer to postulate forces on objects and force transfer between interacting objects. In this paper we make these ideas precise with an implemented system that can make causal inferences directly from video sequences.

The use of domain knowledge by a vision system has been studied extensively for both static and motion domains. Our work addresses a number of important limitations in prior work:

- Our system models kinematic and dynamic information about the world, in addition to static information available from a single frame.
- Our system embodies a sound inference procedure based on an explicit physical theory.

¹ Correspondence should be sent to A. Jepson, Department of Computer Science, University of Toronto, 6 Kings College Road, Toronto, Ontario M5S 3H5. Tel: 416/978-6488; Fax: 416/978-1455; E-mail: jepson@cs.toronto.edu.

- Our system can represent uncertainty to make inferences from ambiguous input.
- Finally, our system generates such inferences directly from camera input.

We postpone a more detailed discussion of related work until Section 8.

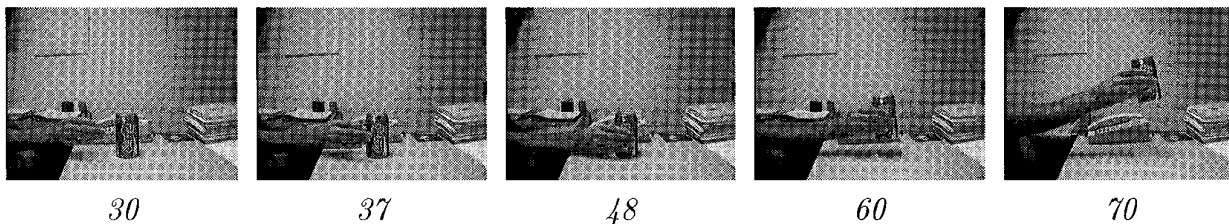
This paper makes three central contributions. First, we provide an ontology suitable for producing interpretations of image sequences in terms of the kinematic and dynamic properties of observed objects. Second, we provide a computational procedure to test the feasibility of such interpretations by reducing the problem to a feasibility test in linear programming. Finally, we provide a theory of preference ordering between multiple interpretations along with an efficient computational procedure to determine such orderings.

2. OVERVIEW

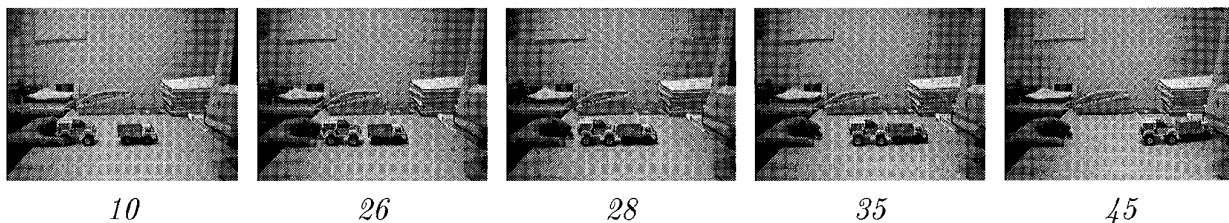
It is useful to first present a brief overview of the main components of our system. Here we emphasize the basic intuition behind each of the components and the prerequisites for their use. In later sections we describe each of the components in detail and present experimental results.

We provide this overview in terms of a single illustrative example, the **coke** sequence, shown on the top row of Fig. 1. (The sequences displayed there are used for the computational examples in Section 7.) In particular, a hand is reaching for, grasping, and then lifting a coke can off a desk top. As mentioned in the Introduction, our eventual goal is to have a machine vision system that, when given image sequences such as this one, can understand the basic force generation and force transfer relationships of the various objects in the scene. For this example sequence the system should understand that the can is initially supported by the table and that the hand (and arm) is possibly an “active object.” Roughly speaking, an active object is something that can generate forces other than those due to gravity, friction, and acceleration. Objects that are not active, such as the table and the coke can, are said to be “passive objects.” Later, during the lifting phase of this

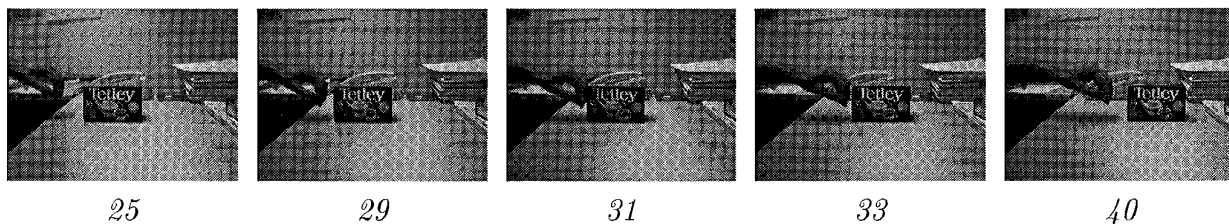
coke



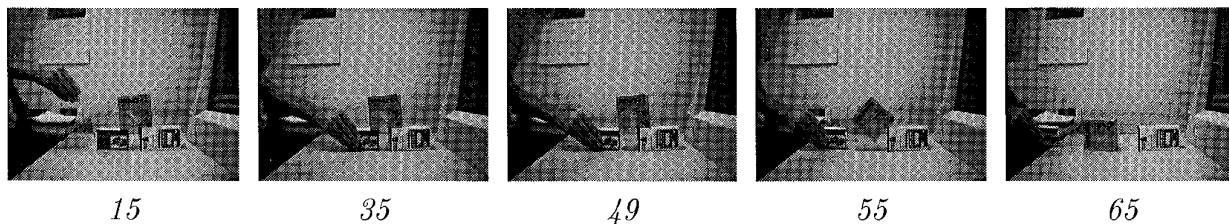
cars



hit



arch



tip

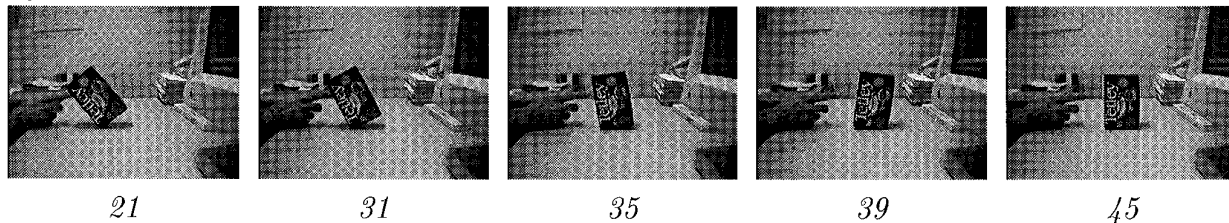


FIG. 1. The example sequences: **coke**, **cars**, **hit**, **arch**, and **tip**. The frame numbers are given below each image.

sequence, if the system remembers that it has determined that the hand is an active object, then it should conclude that the hand is attached to the coke can (i.e., grasping it) and applying an upward force on it.

This goal of understanding forces and dynamics given image sequences is an ambitious one. Obviously, we do not achieve it completely here. However, we do present an implementation that takes us much of the way to our

goal, albeit in a simple domain. The critical simplifications are: (1) objects are represented as rigid 2D polygons in a layered 2D representation; (2) the analysis is based on only the instantaneous motion of the objects, not on the behavior of the objects throughout the sequence. Such a representation provides a crude but useful approximation for a variety of situations, including all the sequences depicted in Fig. 1. Moreover, the simplicity of the domain serves to highlight the basic principles behind our approach. It is our belief that these principles can be applied to general 3D problems and to problems involving integration of inferences through time.

2.1. An Outline of the Approach

To reason about dynamic properties of the objects in the scene, we use the equations of Newtonian mechanics. Given a concrete hypothesis about the scene content, we simply express the corresponding equations of Newtonian mechanics in a suitable form and check to see whether they have a feasible solution.

This physics-based approach places stringent requirements on the richness and the level of detail of the underlying scene model. Indeed, we have the basic requirement that the geometry of the scene must be completely specified in order to express these equations. That is, we need to know object shape and placement, along with the surfaces of contact between objects. In addition, we need estimates for the velocities and accelerations of the observed objects. We also need to assume something about the distribution of mass within each object. For our purposes here, we simply take an object’s center of mass to be its centroid, as if it had a uniform mass distribution. The total mass and the inertial tensor are treated as constrained unknowns when the system checks for a feasible solution. In summary, a strong prerequisite on the hypothesis-generation mechanism is that a rather detailed description of the scene is needed before the physics-based modeling can be applied.

It is convenient to define a *configuration* to be the set of scene properties that are necessarily present, given the image data and any restrictions inherent in the ontology. For example, in the current system, the positions, velocities, and accelerations of the objects are provided by the image observations, and the positions of the centers of mass are fixed, by our ontology, to be at the object centroids.

For a concrete example, consider the lifting phase in the **coke** sequence. Through the use of a tracking algorithm, we obtain the 2D positions, velocities, and accelerations of polygons that roughly describe the shapes of the scene objects. This information, along with a line denoting the table surface, constitutes the configuration. The hand and can polygons provided by the tracker are displayed in Fig. 2, along with other symbols, discussed below, which denote

additional properties of the various interpretations. Note that in this example, the system is not given any further information about the objects, so, in particular, the objects are all considered equally likely to be active.

To supply the scene information not included in the configuration, we consider *assertions* taken from a limited set of possibilities. These assertions correspond to our hypotheses about the basic force generation and force transfer relationships between objects. For brevity, we constrain ourselves here to only those assertions that play a significant role in the **coke** sequence. The full set is provided in Section 3. In particular, here we need the following three types of assertions:

1. CONTACT(o_1, o_2, c)—objects o_1 and o_2 contact in the scene with the region of contact c ;
2. ATTACH(o_1, o_2, p)—objects o_1 and o_2 are attached at some set p of points in the contact region;
3. BODYMOTOR(o)—object o has a “body motor.”

The intuitive meanings of these assertions are: (1) objects o_1 and o_2 are contacting, either in depth or abutting; (2) objects o_1 and o_2 are attached on some set p of points in the contact region; these attachment points are functionally equivalent to rivets, fastening the objects together; (3) object o can generate an arbitrary force and torque on itself, as if it had several thrusters. Note that the attachment assertion is properly understood to be a characterization of the types of forces supported at the attachment points. Attached objects can be pulled, pushed, and sheared without coming apart. Without attachment, contacting objects may separate or slide on each other, depending on the applied forces and on the coefficient of friction. The precise choice of these kinematic and dynamic assertions is not critical for the purpose of this paper. What is important is that they are sufficiently rich to describe a wide variety of phenomena.

Note that, given the hand and can polygons in Fig. 2, there is no evidence that any of these assertions are individually true or false. It is for this reason that neither they, nor their negations, belong to the configuration. Furthermore, notice that some sets of assertions are not *admissible* in that they violate basic constraints for their application. For example, two objects cannot be both attached and not contacting. Moreover, some sets of assertions are not *complete* in that they leave some properties unspecified. We define an *interpretation* to be a combination of a configuration along with an admissible and complete set of assertions that fully specify the scene.

Given any interpretation, we can check the feasibility according to whether or not a “force balance” exists in our physics-based model. This feasibility test can, to a good approximation, be reduced to a linear programming problem (see Section 5).

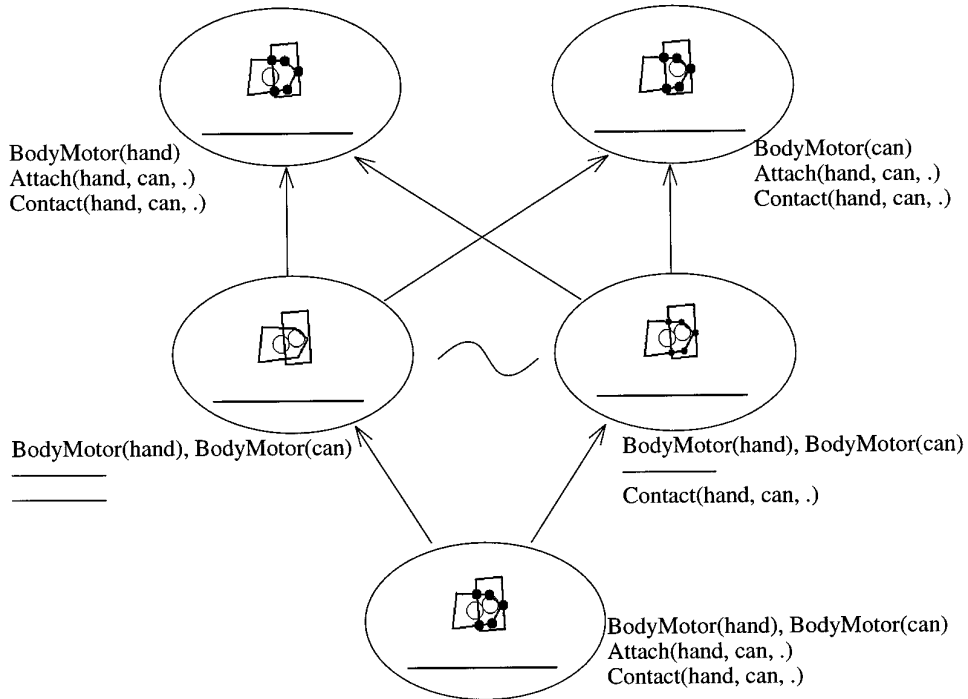


FIG. 2. The preference ordering for the five feasible interpretations of a single frame during the lifting phase of the **coke** sequence. The *arrows* indicate a preference relation between two interpretations while the *tilde* indicates that two interpretations are equally preferred. A large open circle at the object center denotes a BODYMOTOR; the large disks at the vertices of the polygons denote ATTACHED objects; the small disks at the vertices of the polygons denote CONTACTING objects. A textual form of the assertions appears adjacent to each interpretation. The three levels of priority are represented by each line of text. Note that the absence of an assertion denotes its negation.

For the lifting phase of the **coke** sequence, there are five feasible interpretations within our ontology, as displayed in Fig. 2. Note that each interpretation has a feasible force balance, given that we can adjust both the masses and the forces generated by the body motors. In contrast, any interpretation without a body motor is infeasible since the vertical acceleration cannot be “explained” (i.e., balanced with a corresponding force). Similarly, the active hand in the top left interpretation in Fig. 2 must be attached to the can since, in our ontology, we do not model forces in depth. This means, in turn, that the physics-based model cannot generate a frictional force between the hand and the can. So, without attachment, the observed vertical acceleration of the can cannot be balanced with a vertical force. While this sort of “reasoning” seems sophisticated, it simply relies on the feasibility of linear programming problems that the system can pose for itself.

A naive system would generate all feasible interpretations thereby producing all (and only) the interpretations depicted in Fig. 2. Our algorithm, however, does not generate all such feasible interpretations since, as we discuss next, many of these interpretations are not of interest. It turns out that the algorithm need not evaluate all interpretations in order to find the interesting ones.

2.2. Preferred Explanations

Up to this point we have discussed the nature of an interpretation in terms of a scene configuration along with assertions that specify further dynamic and kinematic properties. We pointed out that the feasibility of such an interpretation can be checked using a physics-based model for the forces involved along with the admissibility of the assertions. However, given a fairly rich ontology, it is common for there to be multiple feasible interpretations for a given scene configuration (see Fig. 2).

Indeed, given that a body motor is capable of generating an arbitrary force and torque on the object possessing it, it follows that for any scene configuration there is always at least one trivial interpretation in which every object has a body motor. The bottom three interpretations in Fig. 2 are examples of such trivial solutions. They are guaranteed to pass the force-balance feasibility criterion and, as such, are not very informative interpretations on their own. Rather, we seek interpretations that require, in some specified sense, the *weakest* properties of the various objects.

Model preference relations, as discussed by Richards, Jepson, and Feldman [27], can be used to express suitable preference orderings. The basic idea is simple, namely to

compare two different interpretations in terms of a prioritized set of elementary preference relations. Our current ontology includes the elementary preference for the absence of the assertion $\text{BODYMOTOR}(o)$ for each object o . In other words, we prefer to see a particular object as a passive object, if feasible, given that other elementary preference relations do not contradict this preference. In addition, we prefer the absence of a kinematic constraint of attachment between two objects over the presence of such a constraint. Finally, we are indifferent as to whether we assert that two objects are contacting or not (always assuming admissibility and feasibility).

These three elementary preference relations are taken to be stratified in terms of priority. In particular, the preference for the absence of a body motor is taken to have a higher priority than the preference for no attachment. Finally, we have no elementary preference relation for contact over no contact—these two situations are considered to be equally preferable. This priority ordering is indicated in Fig. 2 by listing the assertions adjacent to each interpretation. Assertions are placed on different lines to indicate their priority. It is convenient to place contact assertions on the bottom row even though there is no preference between different contact assertions.

Given these three elementary preference relations along with the priority ordering, the induced ordering on our five feasible interpretations for the **coke** example is provided simply by a prioritized subset ordering. In particular, suppose one is given two feasible interpretations along with their corresponding assertion sets. If, at the highest level of priority, one interpretation’s assertion set is a strict subset of the other’s, then that interpretation is preferred. For example, the top left interpretation in Fig. 2 is preferred over the three lower interpretations because it does not assert a body motor on the can. In a sense, not needing a body motor is considered to be major simplification in what the various objects are asserted to be capable of doing. Thus the preference against body motors is placed at a high priority.

Alternatively, if the assertion sets at the highest priority do not satisfy a subset relation, the interpretations are considered to be unordered. For example, the top two interpretations in Fig. 2 are considered unordered since the body motor is on different objects. If the sets at the highest priority are the same, we check the assertions at the next lower priority. For example, the three lower interpretations in Fig. 2 all contain the same motor assertions, but they can be further ordered, based on minimizing the attachment. Finally, if interpretations agree on both motor and attachment assertions, we are indifferent as to the preference relation. This occurs in Fig. 2 for the middle two interpretations, since these differ only in terms of their contact relations. (Indifference between the two interpretations is denoted by the *tilde* in Fig. 2.)

This approach to preference ordering will not always yield a unique preferred interpretation. Therefore, multiple percepts are possible. For example, in the **coke** sequence there are two maximally preferred interpretations: either the hand is lifting the can, or vice versa. Given that our ontology includes nothing about hands or cans, that is, we just have moving polygons and notions of mass and force, arriving at these two interpretations for this single frame is intuitively the right thing to do. Indeed, Jepson and Richards [16] and Richards, Jepson, and Feldman [27] propose that such maximal interpretations provide a computational model for a “percept.” Moreover, Richards, Jepson, and Feldman [27] explore the relationship between such preference orderings and qualitative probabilistic models (see also Jepson, Richards, and Knill [17]).

2.3. Implementation and Limitations

Our current implementation has a number of limitations. One limitation is that our system uses a 2D layered representation. Given this limited representation, we can process only fronto-parallel scenes and cannot reason about occlusion or motion in depth. This is a reasonable approximation if we assume that objects move in planes roughly fronto-parallel to the camera. Furthermore, we believe that this limitation is not fundamental to our general approach. In particular, our current system is able to reason about 3D scenes given suitable 3D input. We are currently investigating approaches for tracking 3D deforming objects such as the methods described by Black and Jepson [3].

Another limitation is that our system does not currently integrate information over time or reason about object properties. For example, consider the **coke** sequence once again. During the reaching phase of this sequence, the hand is seen to be above the table and moving horizontally. Our system concludes that the hand must be an active object while the coke can could be passive and at rest on the table. This is the unique maximally preferred solution for this frame. As we have just shown, however, there are multiple interpretations of the scene during the subsequent lifting phase; either the hand or the coke can could be active, but not both. Since our system lacks prior knowledge of object capabilities (such as the fact that hands are typically active, and coke cans are typically passive) and does not integrate conclusions over time (such as noticing that it previously concluded the hand must be active), the system is left with multiple interpretations for the lifting phase.

Finally, our system is limited in that object velocities and accelerations are taken to be continuous functions of time. In particular, we do not consider force impulses that give rise to step discontinuities in the velocities.

Despite these limitations, we are left with a surprisingly rich domain, as is indicated by the variety of computational

examples displayed in Fig. 1. We return to these limitations and ways to overcome them in Section 9.

3. ONTOLOGY

In this section, we discuss the details of the representation used in the implementation and motivate some of the choices made. As discussed above, our current implementation uses a layered 2D representation of the scene. In addition, we assume that we are given estimates for the object velocities and accelerations at some point in time where the motion is continuous.

Given this simplified domain, we now describe the ontology necessary to represent interpretations. We begin with a description of the kinematic and dynamic properties of the configuration, followed by a description of the assertions.

3.1. Kinematic Model

The basic primitive for an object part is a rigid two-dimensional convex polygon. A single *object* is a rigid union of convex polygons. To represent the spatial relationships between objects in the scene we use a *layered* scene model. In our layered model there is no depth ordering. Instead, we represent only whether two objects are in the same layer, in adjacent layers, or in layers separated in depth. Objects can contact either within the same layer or between adjacent layers. The first type of contact, called *abutting contact*, occurs when two objects in the same layer contact at a point or at an edge along their boundary. The second type of contact, called *overlapping contact*, occurs when two objects in adjacent depth layers contact over part of their surfaces and the region of overlap has a nonzero area. We denote both types of contact by $\text{CONTACT}(o_1, o_2, c)$, where o_1 and o_2 are the objects and c is a *contact region* between the two objects.

In addition to position information, the relative motion of objects constrains their allowable contact and layer relations. In the case of abutting contacts, contact is admissible only when the relative velocity of the contacting objects is tangential to the contact region (i.e., objects can slide along their contact region, but they cannot penetrate or separate). In the case of overlapping contacts, the objects must be in different layers. Finally, the relative depth of objects that do not overlap or abut is left unspecified.

3.2. Dynamic Model

In order to perform force balancing we need each object’s center of mass, total mass, and moment of inertia. To determine these properties from image data, we assume that the center of mass of each object is taken to be at that object’s geometric center. Object masses, however,

are left as free parameters (constrained to be positive) that can be adjusted by the physics-based model in order to find a feasible solution.

We also need estimates of the inertial tensors for objects. For the case of two-dimensional motion considered in this paper, the inertial tensor I is a scalar. In order to reflect the uncertainty of the actual mass distribution, we allow a range for I . An upper bound for I is provided by considering an extreme case, where all of the mass is placed at the furthest point from the center. A lower bound is provided by considering an alternate case where all of the mass is distributed uniformly inside a disk inscribed in the object. Together, these provide the constraint

$$\frac{1}{2}Mr_{\min}^2 \leq I \leq Mr_{\max}^2, \quad (1)$$

where M is the object mass and r_{\min} and r_{\max} are the minimum and maximum radii of the object, respectively.

An object is subject to gravitational and inertial forces and to forces and torques resulting from contact with other objects. The dynamics of the object under these forces is obtained from the physics-based model described in detail in Section 5.

Finally, particular objects may be denoted as *ground*. We typically use this for the table top. Forces need not be balanced for objects designated as ground.

3.3. Assertions and Interpretations

We denote an interpretation constructed from a configuration C and an assertion set A as $i = (C, A)$. For our system, C consists of the object positions, velocities, accelerations, polygonal shapes, and centroids. The assertions A describe additional kinematic and dynamic constraints on the objects.

Currently, our implementation uses the two types of kinematic assertions defined in Section 2, $\text{CONTACT}(o_1, o_2, c)$ and $\text{ATTACH}(o_1, o_2, p)$, along with their negations. The admissibility constraints for contact are discussed in Section 3.1 above. The admissibility constraints for attachment require that, in addition to contact, two objects do not exhibit relative motion and that the set p of attachment points is contained within the intersection of the object polygons.

In addition to kinematic assertions, we have the following three types of dynamic assertions, along with their negations:

- $\text{BODYMOTOR}(o)$ —object o has a “body motor” that can generate an arbitrary force and torque on itself;
- $\text{LINEARMOTOR}(o_1, o_2, c)$ —a linear motor exists between the abutting objects o_1 and o_2 . This motor can generate an arbitrary tangential shear force across the motor region c . This region must be contained within the contact region between the objects;

- **ANGULARMOTOR**(o_1, o_2, p)—an angular motor exists at a single point p that can generate an arbitrary torque about that point. p must be within the contact region between the objects.

The notion of a body motor was introduced in Section 2. Linear motors are used to generate a shear force across an abutment (providing an abstraction for the tread on a bulldozer). Angular motors are used to generate torques at joints.

We apply the following admissibility constraints to sets of motor assertions. Body motors are admissible on all objects. Linear motors are admissible only at contacts where the direction for application of tangential force can be defined. Thus linear motors are admissible only at point-to-edge and edge-to-edge abutments but not at point-to-point abutments or overlapping contacts. Angular motors are admissible only at a single point within the contact region between two objects and the objects must be attached at this point.

4. GENERATING HYPOTHESES

In order to demonstrate the applicability of our approach to camera input we have developed a complete implementation for the simplified domain described in Section 3. We describe the various components of our implementation below.

4.1. Configuration

To acquire the position and orientation of the object polygons for each frame we use a *view-based* tracking algorithm similar to the optical flow and stereo disparity algorithms described in [15, 14]. (The full details of the tracking algorithm are described in [22].) In particular, a template image is provided for each object, along with information about where the object is located within the template. Given an initial guess for the positions of the objects in the first frame, the tracking algorithm then estimates the two-dimensional position and orientation of these initial templates throughout the image sequence by matching the templates to each successive frame. The position of the object polygons is then obtained by mapping the original outlines according to these estimated positions.

The process of matching the template to each frame is *robust*. That is, poorly matching pixels, called outliers, are rejected so they do not affect the estimation. For this reason, the shape of the template is not critical so long as it contains a significant portion of the object to be tracked. Note that robustness also allows the tracker to deal with objects that have changing occlusion relationships. (We show such examples in Section 7.)

Given the pose (that is, $x(t)$, $y(t)$, and $\theta(t)$) of the objects in each frame, we estimate the motion by a robust interpola-

tion procedure (see [22] for details). Specifically, at each point in time, we robustly fit a cubic polynomial to the data over a sliding temporal window that is seven frames wide. We then differentiate these polynomials to obtain estimates for the velocity, acceleration, angular velocity, and angular acceleration of each object. Often, a single polynomial will not fit the pose data. For example, if there is a step change in velocity or acceleration (e.g., due to a collision or a change in contact between objects), a single polynomial will provide inaccurate estimates near the discontinuity. Our system avoids this problem by fitting two or more polynomials to the data within each temporal window. In addition, we allow outliers to be excluded from the fitting process. Note that while this approach will eliminate biased estimates of velocity and acceleration near a discontinuity, the velocities and accelerations at the discontinuities themselves are not defined. As described in Section 2, since our system assumes continuous motion, the interpretations formed at motion discontinuities will often be anomalous. We return to this issue when we show experimental results in Section 7.

As discussed in the previous section, this position, velocity, acceleration, and shape data constitute the configuration component C of any interpretation. In order to hypothesize an interpretation, we therefore need to select an admissible and complete set of assertions A .

4.2. Assertions

In order to explore the space of possible interpretations, we must first construct a set of admissible assertions. Given that the allowable forces between objects depend on the contact geometry and the relative motion of the objects [8, 2], an analysis of the scene kinematics is necessary. Since we do not have exact shape or motion information, however, we need a way to determine which contact relations are possible. In general, to determine the possible contact relations we must consider interactions among all the objects. For the purposes of this paper, however, we implement a partial test in which we consider only pairwise constraints between objects. (We will see some examples of the limitations of this test in the experiments in Section 7.)

According to our ontology, contacting objects must either abut or overlap. We can classify the contact type by examining the region of intersection between objects. If the width of the intersection region is greater than a specified tolerance the contact type must be overlap.² If the width of the intersection region is less than that tolerance, we assume the contact type to be abutment. Overlapping contact is always admissible. Abutting contact, however, must satisfy the additional constraint that the relative motion

² For each of the experiments reported in Section 7 we selected a single intersection tolerance. The value used was between 4 and 8 pixels.

of the two objects be tangential to the boundary between those objects.³ Abutting contacts that do not satisfy this constraint are inadmissible. Furthermore, abutment is not admissible for objects in point-to-point contact, except when the objects are attached at the contact point.

An attachment relation between two objects is admissible only when those two objects are in contact and their relative velocities and accelerations at the attachment points are less than a specified tolerance. Given contact between two objects, our ontology allows attachment at any set of points in the contact region. To reduce the number of hypotheses that we need to consider, however, we restrict attachment to the vertices of the perimeter of the contact region. Furthermore, in the system described here we consider only those hypotheses where all such vertices are attached.

Given the contact and attachment relations, the admissibility of the dynamic assertions is determined as follows. `BODYMOTOR`(o) is always admissible for any object o . `LINEARMOTOR`(o_1, o_2, c) is admissible for any pair of abutting objects o_1 and o_2 and c is the contact region. Finally, `ANGULARMOTOR`(o_1, o_2, p) is admissible at any pair of attached objects o_1 and o_2 , where p is the attachment point. Our current implementation allows angular motors only when the contact region is a single point, that is, when the contact type is either point-to-point or point-to-edge. In this case, the point p is uniquely determined.

Given this specification, it is now possible to generate all admissible interpretations for a given frame of the image sequence. These are obtained from the configuration provided by the tracker combined with every admissible and complete set of assertions.

Among all possible interpretations we must select those that are feasible, namely those that satisfy our physical theory. Nominally we would have to consider preference relations between all feasible interpretations. In practice, however, we avoid generating this full set of feasible interpretations by exploiting the structure of our domain to reduce the search space. Before we discuss the search algorithm, however, we first describe the mechanism for checking the feasibility of an interpretation.

5. FEASIBLE INTERPRETATIONS

Given a configuration of the scene objects along with a set of assertions about the kinematic and dynamic properties of the scene, we can use a theory of dynamics to determine whether the interpretation has a feasible force

³ To enforce tangential motion, the magnitude of the normal component of the velocity and acceleration must be below a specified tolerance. For each of the experiments reported in Section 7 a single tolerance value was chosen. The velocity tolerance was between 1.0 and 2.0 pixels/frame. The acceleration tolerance was between 0.5 and 1.0 pixels/frame.²

balance. In particular, we show how the test for consistency within the physical theory can be expressed as a set of algebraic constraints that, when provided with an admissible interpretation, can be tested with linear programming. We use a “force balancing” approach similar to that proposed by Blum *et al.* [4]. Our approach, however, models dynamics as well as static force balancing.

In this section, we present a theory for the general three-dimensional case. The experimental results we describe later were produced using a two-dimensional variant of this theory.

For rigid bodies under continuous motion, the dynamics are described by the Newton–Euler equations of motion [12]. For rigid bodies of nonvarying mass, the appropriate equations are

$$\begin{aligned} \mathbf{F} &= \dot{\mathbf{p}} \\ &= M\dot{\mathbf{v}} \end{aligned} \quad (2)$$

$$\begin{aligned} \mathbf{N} &= \dot{\mathbf{L}} \\ &= \mathbf{I}\dot{\boldsymbol{\omega}} + \boldsymbol{\omega} \times \mathbf{I}\boldsymbol{\omega} \end{aligned} \quad (3)$$

The first equation relates the total applied force \mathbf{F} to the rate of change of linear momentum \mathbf{p} . For bodies with nonvarying mass, this reduces to $M\dot{\mathbf{v}}$, where M is the object mass and $\dot{\mathbf{v}}$ is the acceleration. The second equation relates the total applied torque \mathbf{N} to the rate of change of angular momentum \mathbf{L} . For rigid bodies, the rate of change of angular momentum is given by $\mathbf{I}\dot{\boldsymbol{\omega}} + \boldsymbol{\omega} \times \mathbf{I}\boldsymbol{\omega}$, where \mathbf{I} is the inertial tensor, $\boldsymbol{\omega}$ is the angular velocity, and $\dot{\boldsymbol{\omega}}$ is the angular acceleration.

Given a scene with convex polygonal objects, we can represent the forces between contacting objects by a set of forces acting on the vertices of the convex hull of their contact region [8]. Under this simplification, the equations of motion for each object can be written as

$$M\dot{\mathbf{v}} = M\mathbf{g} + \mathbf{F}_b + \sum_{c \in \Gamma} \delta_c \sum_{p \in c} \mathbf{F}_p \quad (4)$$

$$\mathbf{I}\dot{\boldsymbol{\omega}} + \boldsymbol{\omega} \times \mathbf{I}\boldsymbol{\omega} = \mathbf{N}_b + \sum_{c \in \Gamma} \delta_c \left(\sum_{p \in c} \mathbf{N}_p + (\mathbf{r}_p - \mathbf{r}) \times \mathbf{F}_p \right), \quad (5)$$

where \mathbf{g} is the acceleration due to gravity, \mathbf{r} is the center of gravity of the object, \mathbf{r}_p is the position of point p , and \mathbf{F}_b and \mathbf{N}_b are unknown body forces and torques that may act upon the object. We use the body forces and torques to implement the `BODYMOTOR` assertion in our ontology. Γ denotes the set of contact regions involving the object. For each contact region $c \in \Gamma$ we add the forces \mathbf{F}_p at each contact point $p \in c$. In addition, we allow terms \mathbf{N}_p at each contact point. These are used to implement the `ANGULARMOTOR` assertions in our ontology. Finally, the torques $\delta_c \in \{-1, 1\}$ encode the direction of the contact forces.

The signs of δ_c are arbitrary as long as they are consistent between contacting objects.

5.1. Contact Conditions

Assuming that there are no degenerate contacts we can represent each contact region c by a set of one or more contact points and a vector \mathbf{n}_c normal to the contact region. Contact points that are not asserted to be ATTACHED must obey the normal force constraint

$$\mathbf{F}_p \cdot \mathbf{n}_c \geq 0, \quad (6)$$

where \mathbf{F}_p is the contact force at each point $p \in c$.

In addition, contact points that are not part of a LINEARMOTOR will have tangential forces limited by friction. When the relative velocity at the contact point is zero we constrain the tangential forces by a *Coloumbic friction* model

$$\|\mathbf{F}_p^t\| \leq \mu_r \mathbf{F}_p^n, \quad (7)$$

where \mathbf{F}_p^t is the tangential component and \mathbf{F}_p^n is the normal component of the contact force. The coefficient μ_r is the *coefficient of resting friction*. In the case of objects with relative motion, we limit the tangential forces as above, but use a different coefficient μ_s , the *coefficient of sliding friction*. In general we will have $\mu_s < \mu_r$. In addition to the magnitude constraint, we also need to limit the direction of sliding friction to be *opposing* the direction of motion. This is achieved by the additional constraint

$$\mathbf{F}_p \cdot \mathbf{v}_p \leq 0, \quad (8)$$

where \mathbf{v}_p is the relative velocity of the contact point, p , within the contact plane.

5.2. Testing Feasibility

An interpretation is *feasible* if the motion equations can be satisfied subject to the contact conditions and the bounds on the mass and inertia described in Section 3.2.

We can approximate these constraints by a set of linear equalities and inequalities as follows. We write the normal force and sliding friction constraints directly. The motion equations can be written as a set of linear equalities by taking the x , y , and z components separately. We can approximate the bounds on the magnitude of the tangential component of the frictional forces by restricting the magnitude in a set of two or more component directions.

Finally, we add conditions on the body forces and torques based on the assertions in A . For passive objects, we constrain each of the component directions of the body forces and torques to be zero, within tolerances F_0 and N_0 ,

respectively. Note that the use of tolerances is necessary since the observed motion will never be exactly zero.

Since all of the above equations and inequalities are linear, dynamic feasibility can be reduced to a feasibility test using linear programming [21].

6. PREFERENCES

As described in Section 3, we have a fixed set of elementary preference relations, namely

- $P_{\text{bodymotor}}(o): \neg \text{BODYMOTOR}(o) > \text{BODYMOTOR}(o)$;
- $P_{\text{linearmotor}}(c): \neg \text{LINEARMOTOR}(o_1, o_2, c) > \text{LINEARMOTOR}(o_1, o_2, c)$;
- $P_{\text{angularmotor}}(c): \neg \text{ANGULARMOTOR}(o_1, o_2, p) > \text{ANGULARMOTOR}(o_1, o_2, p)$.

Here \neg denotes the negation of the predicate that follows and $>$ denotes the preference relation. These elementary preference relations all encode the specification that it is preferable not to resort to the use of a motor, all else being equal. These elementary preference relations appear at the highest priority (recall Fig. 2).

At the next level of priority we have

- $P_{\text{attach}}(o_1, o_2, p): \neg \text{ATTACH}(o_1, o_2, p) > \text{ATTACH}(o_1, o_2, p)$,

so the absence of an attachment assertion is also preferred.

Finally, at the lowest level of priority, we have the indifference relation

- $P_{\text{contact}}(o_1, o_2, c): \neg \text{CONTACT}(o_1, o_2, c) \sim \text{CONTACT}(o_1, o_2, c)$,

so the system is indifferent to the presence or absence of contact, all else being equal.

We wish to find interpretations that are maximally-preferred subject to these preference relations. In this paper we consider a special case of the more general orderings described by Richards, Jepson, and Feldman [27], where the elementary preference relations can be of the form $P(x) > Q(x)$ for predicates P and Q . In addition, we adopt the convention that the absence of an assertion indicates its negation. Thus all of the above preferences, except for the indifference to contact, have a particularly simple form: a preference for the negation of an assertion over the assertion itself.

When the elementary preferences can be written in this simple form, the induced preference relation on interpretations is given by prioritized subset ordering on the set of assertions made in the various feasible interpretations. As described in Section 2, we can determine the preference order for any two interpretations by first comparing the assertions made at the highest priority. If the highest priority assertions in one interpretation are a subset of the highest priority assertions in a second interpretation, the

first interpretation is preferred. Otherwise, if the two sets of assertions at this priority are not ordered by the subset relation, that is neither set contains the other, then the two interpretations are considered to be unordered. Finally, in the case that the assertions at the highest priority are the same in both interpretations, then we check the assertions at the next lower priority, and so on. This approach, based upon prioritized ordering of elementary preference relations, is similar to prioritized circumscription [23].

To find maximally preferred interpretations, we search the space of possible interpretations. We perform a breadth-first search, starting with the empty set of assertions, incrementally adding new assertions to this set. Each branch of the search terminates upon finding a minimal set of assertions required for feasible force balancing. Note that because we are indifferent to contacts, we explore every admissible set of contact assertions at each stage of the search. In theory, this search could require the testing of every possible interpretation. In practice, however, we often examine only a fraction of the interpretations since the search terminates upon finding minimal assertion sets.

Furthermore, when the assertions are stratified by a set of priorities, we can achieve significant computational savings by performing the search over each priority level separately. For example, under our preference ordering, we can search for minimal sets of motors using only interpretations that contain all admissible attachments. It is critical to note that this algorithm is correct only because of the special structure of the assertions and the domain. The critical property is that if there is a feasible interpretation $i = (C, A)$ and if A' is the set obtained by adding all of the admissible attachments to A , then the interpretation $i = (C, A')$ is also feasible. This property justifies the algorithm above, where we set all of the lower priority assertions to the most permissive settings during each stage of the minimization. In general we refer to this property as *monotonicity*. (Details of the search algorithm and a proof of correctness are given in [22].)

7. EXPERIMENTAL RESULTS

We have applied our system to several image sequences taken from a desktop environment (see Fig. 1). The sequences were taken from a video camera attached to a SunVideo imaging system. MPEG image sequences were acquired at a rate of 30 frames per second and a resolution of 320×240 pixels. The 24-bit color image sequences were converted to 8-bit grey-scale images used by the tracker.

The input to our system consists of an image sequence, a set of object template images, a polygonal outline of each object, and an estimate for the positions of the objects within the first frame. In addition, we provide an estimate for the *ground plane*, which is designated as a ground object in our ontology. Note that the exact shape of the ground

plane is not critical, so long as it can be provided to the system in the first frame. Given this input, the tracker provides estimates for the object poses and motions in each frame of the sequence. These estimates, along with the polygonal shapes, are used by the interpretation-construction module.

Figures 3 and 4 show the output of the tracker for the **coke** and **tip** sequences, respectively. (Tracking data for all sequences is given in [22].) In each figure, the upper left graph shows the estimates for the x and y component velocities of the object polygons, while the upper right graph shows the estimates for the angular velocities. The lower graphs show the corresponding estimates for the linear and angular accelerations. Note that while the estimates are somewhat noisy, we can clearly interpret the event structure from the graphs. In Fig. 3, for example, we can distinguish the two distinct phases (reaching and lifting) by examining the velocity and acceleration of the hand and the can. Figure 4 shows an example that demonstrates the necessity for the robust interpolation algorithm described in Section 4.1. Even though there are sharp changes in the velocities and accelerations, our interpolation algorithm was able to obtain reasonable estimates of the motion.

Figure 6 shows some of the preferred interpretations found for selected frames from each sequence. (Note that the selected frames do not necessarily match those shown in Fig. 1.) For each sequence we show frames ordered from left to right. A legend of symbols used to indicate assertions is shown in Fig. 5. To simplify the presentation, except for the **cars** sequence, we show only the interpretations involving body motors. For the **cars** sequence we show only those interpretations using linear motors. While the preferred interpretations are often unique, at times there are multiple interpretations, particularly when objects interact. We highlight frames with multiple preferred interpretations by grey shading.

Our machine interpretations are surprisingly intuitive. For example, the difference between interpretations 1 and 2 in frame 63 of the **coke** sequence can be interpreted as the hand “lifting” the can versus the can “lifting” the hand. Similarly, the difference between interpretations 1 and 2 in frame 34 of the **cars** sequence can be interpreted as the rear car “pushing” the front car versus the front car “pulling” the rear car. Note that the system correctly hypothesizes an attachment between the front and rear cars in the “pulling” interpretation, but it does not do so in the “pushing” interpretation. Note that without interframe analysis or prior information about the objects, all of these interpretations are reasonable.

In addition to making inferences about which objects have motors and which objects are attached, our system can also make inferences about the contact relations between objects in the scene. For example, during the reaching

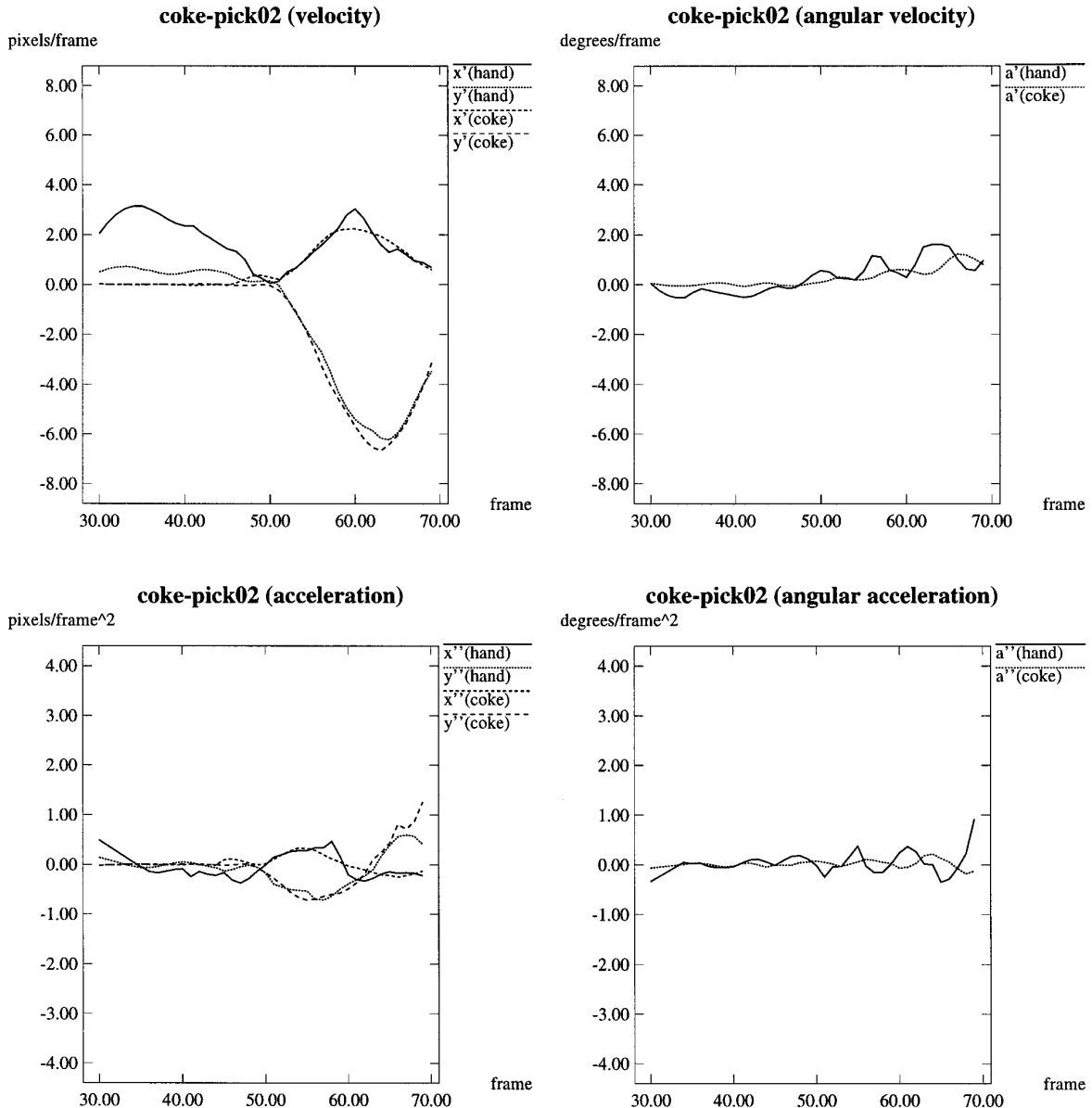


FIG. 3. Tracking results for the **coke** sequence using the view-based tracker followed by the robust interpolation algorithm. The top row shows the estimates for velocity and angular velocity. The bottom row shows the estimates for acceleration and angular acceleration.

phase of the **coke** sequence (frame 32), the system can infer that the can must be contacting (i.e., sitting on) the table. In the remaining examples, however, we do not explicitly consider variations in the contact geometry. Instead we use only the *maximal* set of contacts in each frame.⁴ The reason for this is twofold. First, using the maximal contact set significantly reduces the search space. (At each stage of the search we only need to consider a single contact

⁴Note that because the individual contact assertions are *independent* (they depend only on the configuration C), there will always exist a unique maximal admissible set of contacts.

set.) Second, finding the feasible solutions with the maximal contact set is sufficient to determine the minimal sets of motors and attachments. The reason for this is that for any feasible interpretation, there will always exist a corresponding interpretation with the same set of motors and attachments, but with the maximal contact set.

A physics-based ontology that includes dynamics allows a richer set of descriptions than one based purely on static scenes. This is illustrated by the **arch** example in the fourth row of Fig. 6. In the sequence the top block of the arch changes from being “supported” in frame 45 to “tipping” in frame 52 when the supporting block is removed. Even

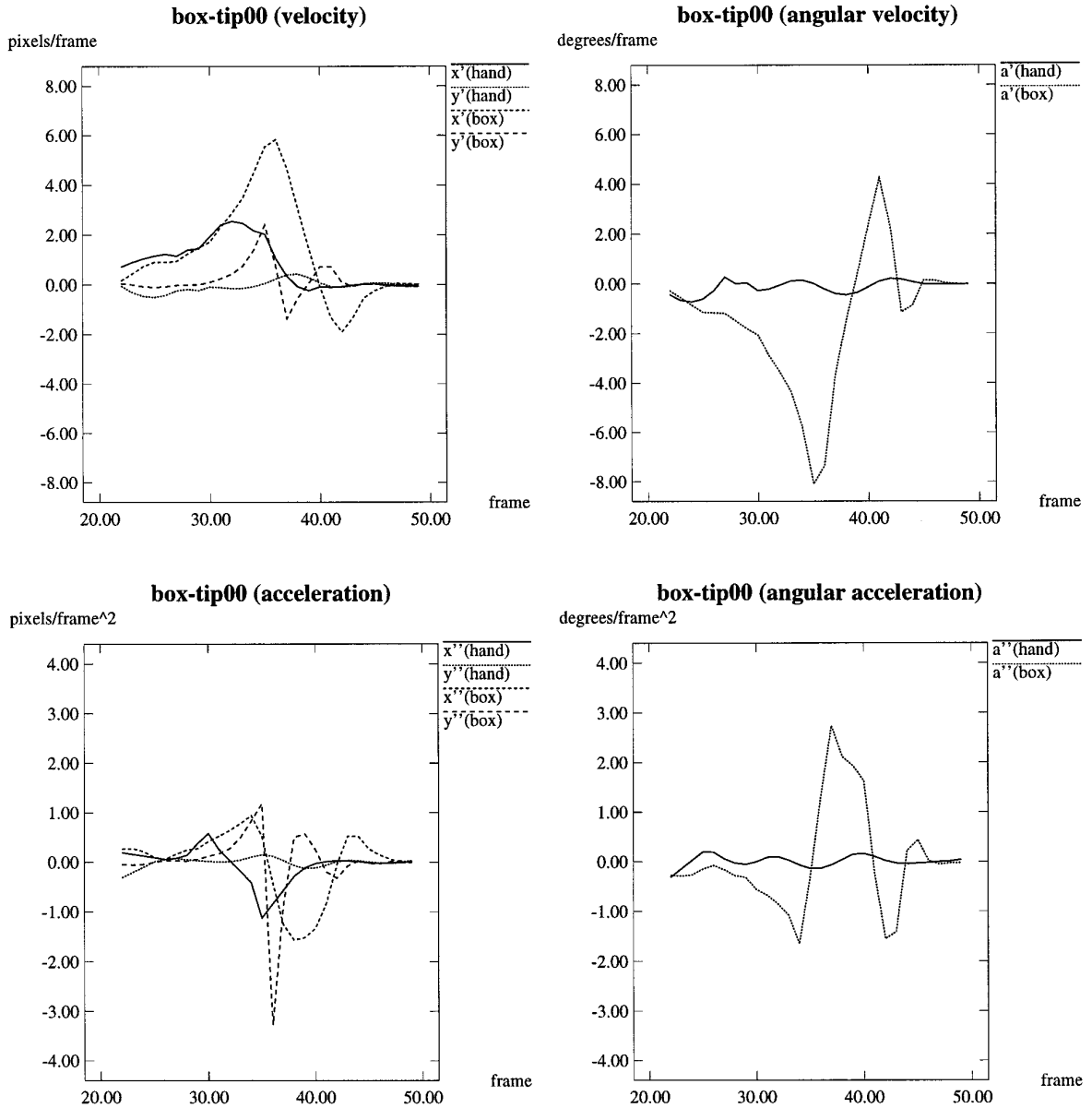


FIG. 4. Tracking results for the **tip** sequence using the view-based tracker followed by the robust interpolation algorithm.

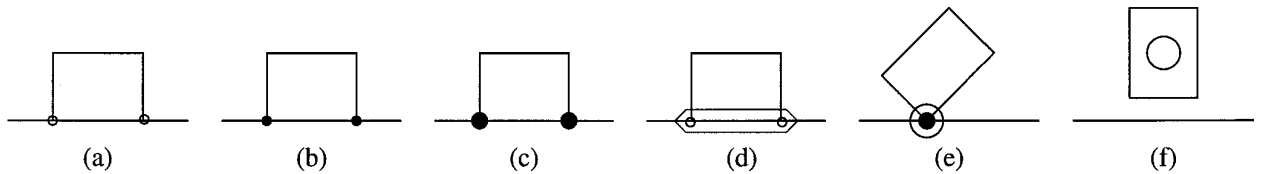
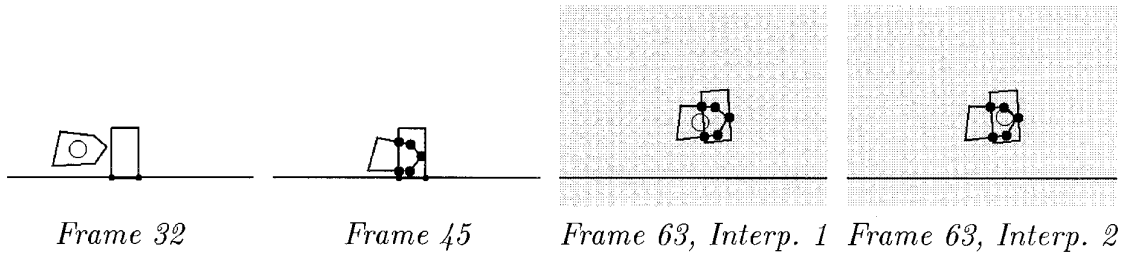
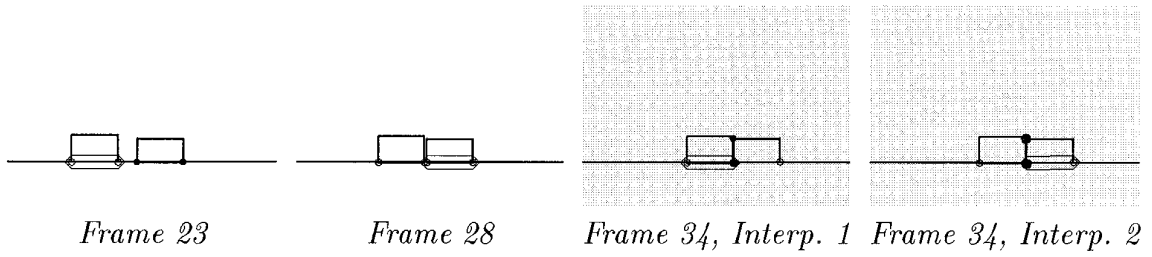


FIG. 5. In the presentation of results to follow, we use: (a) small circles to depict sliding contact; (b) small disks for nonsliding contact; (c) large disks for attachment; while (d), (e), and (f) depict $\text{LINEARMOTOR}(o_1, o_2, c)$, $\text{ANGULARMOTOR}(o_1, o_2, p)$, and $\text{BODYMOTOR}(o)$, respectively. For the first two motors, the closed curve surrounds the contact region over which the motors operate, while for body motors the large circle is placed at the object center.

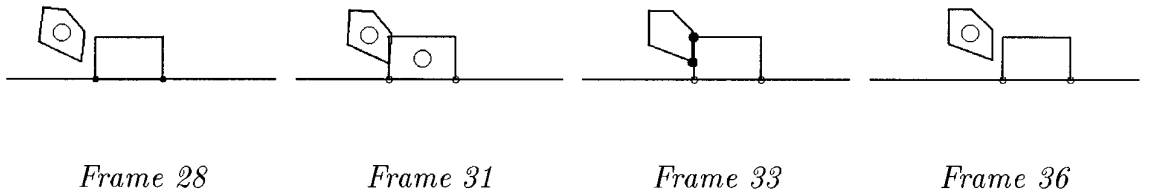
coke



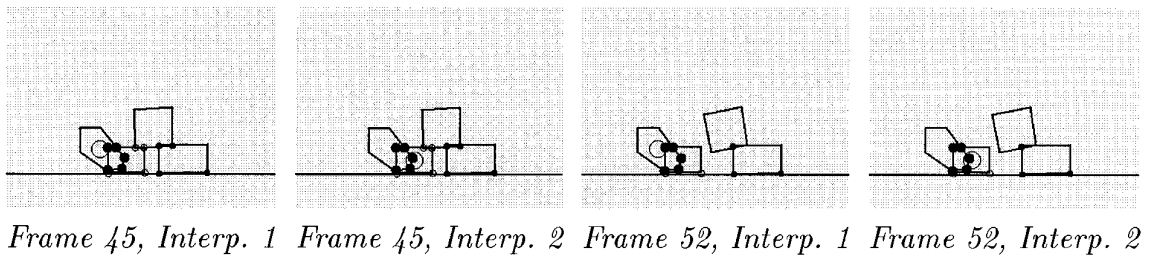
cars



hit



arch



tip

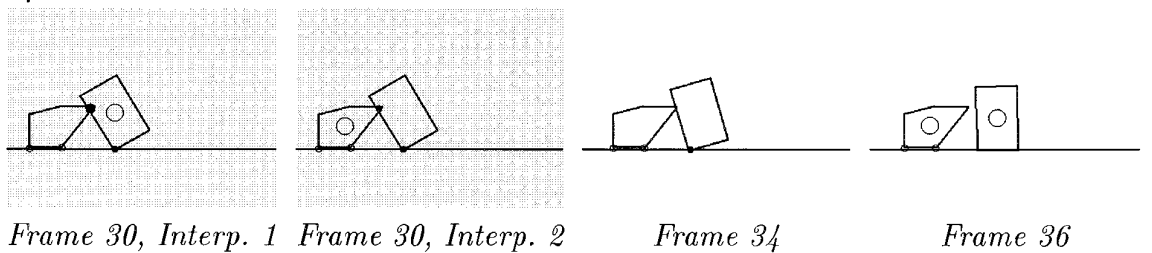


FIG. 6. Some preferred interpretations for **coke**, **cars**, **hit**, **arch**, and **tip**. Frames with a nonunique maximally preferred interpretation are shown with a grey background. Unique interpretations are shown on white.

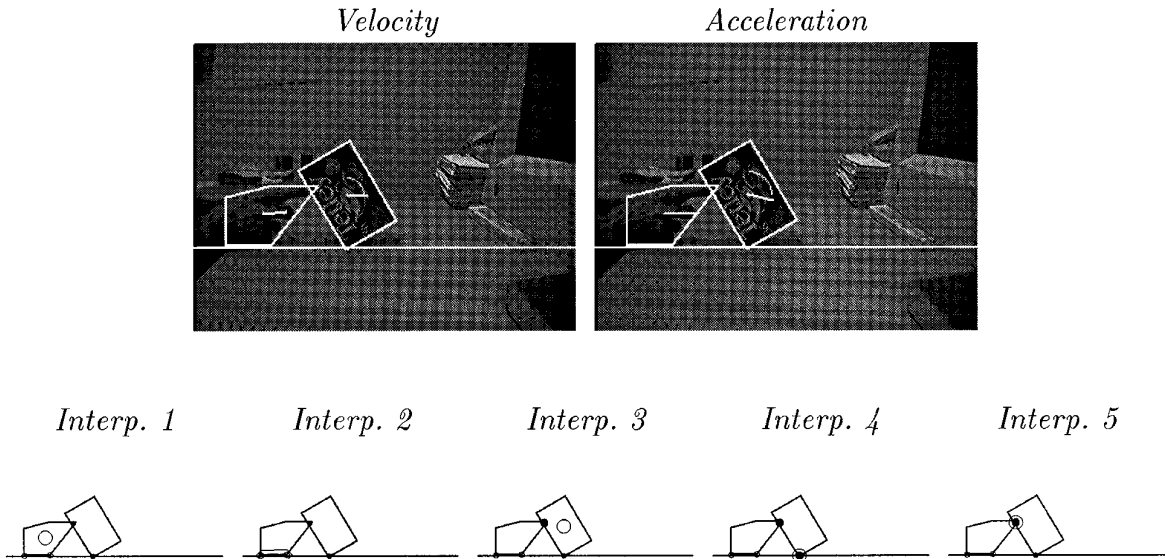


FIG. 7. How many ways are there to tip a box? The velocity and acceleration estimates for frame 30 of the **tip** sequence are depicted on the top row. (The line segments at the object centers indicate the linear velocity and linear acceleration. The arcs at the object centers indicate the angular velocity and angular acceleration.) All five maximally preferred interpretations are given below. Note that only the first two interpretations are plausible—the last three have a force balance, but are not consistent with rigid body motion. See the text for details.

though the left block is moving in frame 45, it is still partially supporting the top block. After the support is removed, the top block begins to move. Since the motion can be explained by gravity alone, however, the top block is still seen as a “passive” object. (Note that, as with the previous examples, the interpretations are ambiguous as to whether the hand is “pulling” the left block or whether the left block is “carrying” the hand.)

Finally, the **tip** sequence in the last row of Fig. 6 highlights the richness of our dynamic domain. In this sequence, a hand raises a box onto its corner and allows it to tip to an upright position. Precisely at frame 30, the box changes from a state of being supported by the hand to a state of tipping. There are two interpretations for frame 30 that involve only body motors: one where the hand “tips” the box and the other where the box “tips” itself and “drags” the hand. As with the previous examples, the latter interpretation requires an attachment between the box and the hand while the former does not. Frame 34 shows the situation shortly after the box is tipped. Note that since the center of gravity of the box is ahead of the support point there is significant angular acceleration. Our system correctly infers that this angular acceleration is due to gravity rather than some type of motor.

Figure 7 shows a more detailed analysis of frame 30 in which we allow all types of motors. The two images show the object polygons with their velocities and accelerations overlaid on the images. Note that the center of gravity of the box is directly above the pivot point, yet there is significant angular acceleration caused by the hand. Below

the images we show all of the preferred interpretations found by our system. The first two interpretations correspond to an active hand “pushing” the box while the last three correspond to an active box “pulling” the hand. Note that while all interpretations in Fig. 7 have a force balance, the last three are not consistent with rigid body motion. In particular, it is not *kinematically* feasible for the hand to be both attached to the box and in edge-to-edge sliding contact with the table. As discussed in Section 4.2, since our system considers only pairwise constraints between objects, it does not check for global kinematic consistency. Further tests would be required to rule out these interpretations.

While encouraging, our system exhibits a number of anomalies. These anomalies generally fall into two classes. In the first type of anomaly, slowly accelerating active objects are sometimes mistaken for passive objects. This results in interpretations in which objects with distinct motions are attached in order to reduce the number of active objects. Examples of this are shown in frame 45 of the **coke** sequence, where the hand attaches to the can, and in frame 33 of the **hit** sequence, where the hand attaches to the box. Such anomalies are to be expected, given that our system only examines single frames in isolation.

A second problem concerns the detection of changing contact relations between objects. In particular, when objects collide, the estimates for the velocity and acceleration at their contact points differ, resulting in the contact relation being deemed inadmissible. An example of this is shown in frame 28 of the **cars** sequence, where, during

a brief interval, the first car is decelerating while the second car is accelerating. While the two cars are actually in contact, this abutment is deemed to be inadmissible due to a large difference in the estimated accelerations. Note that the acceleration of the cars should be equal since they remain in contact after the collision. In our system, however, the interpolator has smoothed over this discontinuity and given unreliable estimates of the acceleration. Further examples of missed collisions are shown in frame 31 of the **hit** sequence and frame 36 of the **tip** sequence. Again, these anomalies are to be expected, since our ontology is restricted to continuous velocity and acceleration and is not designed to handle impulses and abrupt changes in contact.

8. RELATED WORK

The use of domain knowledge by a vision system has been studied extensively for both static and motion domains. Most work in motion understanding has focussed on extracting event descriptions from image sequences based on the spatio-temporal features of the input (see Badler [1], Tsotsos *et al.* [32], Neumann and Novak [24], Borchardt [5], and Kuniyoshi and Inoue [19] for examples). In contrast to these approaches, our work attempts to form descriptions based on a general physical model of the dynamics of the scene.

A number of systems have attempted to represent physical knowledge in static and dynamic scenes using qualitative physical models or rule-based systems (see Fahlman [7], Funt [10], Joskowicz and Sacks [18], Siskind [29, 31], and Brand *et al.* [6]). In contrast to these approaches, our system uses an explicit, quantitative, representation of the dynamics based on Newtonian mechanics.

A number of other systems have used physics-based representations for scenes in terms of forces in static scenes (see Blum *et al.* [4]) and changing kinematic relations in time-varying scenes (see Ikeuchi and Suehiro [13] and Siskind [30]). Our system extends these approaches to consider both kinematic and dynamic relations in time-varying scenes containing rigid objects. Shavit and Jepson [28] present a different approach to classifying motion based on the dynamic properties of nonrigid objects.

Our representation of an interpretation as a set of logical proportions about scene properties is similar to the approach presented by Reiter and Macworth [26]. Unlike our system, however, that system considered a simple domain (maps), where feasibility could be expressed as a set of logical constraints. In addition, while that system allowed the representation of uncertainty, preferences were not used to choose among feasible interpretations.

While the problem is quite different, our representation of geometric and dynamic models borrows heavily from the physical simulation and graphics communities [8, 2].

Finally, it is worth noting that there is evidence that humans generate qualitative physical descriptions of scenes. In particular, there is evidence that humans perceive the force-dynamic relationships among objects in static and dynamic scenes [9, 20]. However, it also appears that humans have a limited understanding of some dynamic events, such as those involving angular motion [25] and collisions between objects [11].

9. CONCLUSION

In this paper we have presented an implemented computational theory that can derive force-dynamic representations directly from camera input. Our system embodies a rich ontology that includes both kinematic and dynamic properties of the observed objects. Finally, the system provides a representation of uncertainty, along with a theory of preferences between multiple interpretations.

While encouraging, this work could be extended in several ways. First, in order to work in a general environment, 3D representations are required. As described in Section 2.3, our system is currently able to represent 3D scenes when provided with suitable input about 3D shape and motion. Further work will be required to determine what type of 3D representation is suitable and how accurate the shape and motion information will have to be.

Second, in order to deal with changing contact relations, a theory of *transitions* is required. Such a theory would require a treatment of the transfer of momentum, as well as forces, between objects. In addition, we need a way to determine the conditions under which the assumption of continuous motion holds. As described in Section 4.1, to determine changing contact relations we will require a representation of kinematics as well as dynamics.

Finally, as described in Section 2.3, to represent the causal structure of time-varying scenes we require a representation of object capabilities and how they are expected to change over time. We believe our current system provides the building blocks for such a representation, but additional work will be required to show how our ontology can be built into a more complex system. (See [22] for preliminary work on time-varying scenes.)

ACKNOWLEDGMENTS

The authors are grateful to IRIS and NSERC Canada for financial support. The authors thank Whitman Richards, Michael Black, and Chakra Chennubhotla for helpful comments on this work.

REFERENCES

1. N. I. Badler, *Temporal scene analysis: Conceptual descriptions of object movements*, Technical Report 80, University of Toronto Department of Computer Science, February 1975.
2. D. Baraff, Interactive simulation of solid rigid bodies, *IEEE Comput. Graphics Appl.* **15**(3), 1995, 63–75.

3. M. J. Black and A. D. Jepson, EigenTracking: Robust matching and tracking of articulated objects using a view-based representation, *Int. J. Comput. Vision*, in press.
4. M. Blum, A. K. Griffith, and B. Neumann, *A Stability Test for Configurations of Blocks*, A. I. Memo 188, MIT Artificial Intelligence Laboratory, February 1970.
5. G. C. Borchardt, Event calculus, in *Proc. Ninth International Joint Conference on Artificial Intelligence, Los Angeles, CA, August 1985*, pp. 524–527.
6. M. Brand, L. Birnbaum, and P. Cooper, Sensible scenes: Visual understanding of complex scenes through causal analysis, in *Proc. Eleventh National Conference on Artificial Intelligence, Washington, DC, July 1993*, pp. 588–593.
7. S. E. Fahlman, A planning system for robot construction tasks, *Artif. Intell.* **5**(1), 1974, 1–49.
8. R. Featherstone, *Robot Dynamics Algorithms*, Kluwer, Boston, 1987.
9. J. J. Freyd, T. M. Pantzer, and J. L. Cheng, Representing statics as forces in equilibrium, *J. Exp. Psychol. Gen.* **117**(4), 1988, 395–407.
10. B. V. Funt, Problem-solving with diagrammatic representations, *Artif. Intell.* **13**(3), 1980, 201–230.
11. D. L. Gilden and D. R. Proffitt, Understanding collision dynamics, *J. Exp. Psychol.: Human Percept. Perform.* **15**(2), 1989, 372–383.
12. H. Goldstein, *Classical Mechanics*, 2nd ed., Addison–Wesley, Reading, MA, 1980.
13. K. Ikeuchi and T. Suehiro, Towards an assembly plan from observation. part i. Task recognition with polyhedral objects, *IEEE Trans. Rob. Automat.* **10**(3), 1994, 368–385.
14. M. Jenkin and A. D. Jepson, Detecting floor anomalies, in *Proc. British Mach. Vision Conf., York, UK, 1994*, pp. 731–740.
15. A. D. Jepson and M. J. Black, Mixture models for optical flow, in *Proc. DIMACS Workshop on Partitioning Data Sets* (I. Cox, P. Hansen, and B. Julesz, Eds.), pp. 271–286, Am. Math. Soc., Providence, RI, 1993.
16. A. D. Jepson and W. Richards, *What is a percept?* Technical Report RBCV-TR-93-43, Department of Computer Science, University of Toronto, April 1993.
17. A. D. Jepson, W. Richards, and D. Knill, Modal structure and reliable inference, in *Perception as Bayesian Inference* (D. Knill and W. Richards, Eds.), pp. 63–92, Cambridge Univ. Press, Cambridge, 1996.
18. L. Joskowicz and E. P. Sacks, Computational kinematics, *Artif. Intell.* **51**(1–3), 1991, 381–416.
19. Y. Kuniyoshi and H. Inoue, Qualitative recognition of ongoing human action sequences, in *Proc. Thirteenth International Joint Conference on Artificial Intelligence, Chambéry, France, August 1993*, pp. 1600–1609.
20. A. M. Leslie and S. Keeble, Do six-month-old infants perceive causality? *Cognition* **25**, 1987, 265–288.
21. D. G. Luenberger, *Linear and Nonlinear Programming*, Addison–Wesley, Reading, MA, 1984.
22. R. Mann, *Computational Perception of Scene Dynamics*, Ph.D. thesis, Department of Computer Science, University of Toronto, in preparation.
23. J. McCarthy, Applications of circumscription to formalizing common sense reasoning, *Artif. Intell.* **28**, 1986, 89–116.
24. B. Neumann and H.-J. Novak, Event models for recognition and natural language description of events in real-world image sequences, in *Proc. Eighth International Joint Conference on Artificial Intelligence, Karlsruhe, August 1983*, pp. 724–726.
25. D. R. Proffitt, M. K. Kaiser, and S. M. Whelan, Understanding wheel dynamics, *Cog. Psychol.* **22**, 1990, 342–373.
26. R. Reiter and A. Mackworth, A logical framework for depiction and image interpretation, *Artif. Intell.* **41**, 1989, 125–155.
27. W. Richards, A. D. Jepson, and J. Feldman, Priors, preferences and categorical percepts, in *Perception as Bayesian Inference* (D. Knill and W. Richards, Eds.), pp. 93–122, Cambridge Univ. Press, Cambridge, 1996.
28. E. Shavit and A. D. Jepson, Qualitative motion from visual dynamics, in *IEEE Workshop on Qualitative Vision, New York, June 1993*.
29. J. M. Siskind, *Naive Physics, Event Perception, Lexical Semantics, and Language Acquisition*, Ph.D. thesis, MIT, January 1992.
30. J. M. Siskind, Axiomatic support for event perception, in *Proc. AAAI-94 Workshop on the Integration of Natural Language and Vision Processing, Seattle, WA, August 1994* (P. McKeivitt, Ed.), pp. 153–160.
31. J. M. Siskind, Grounding language in perception, *Artif. Intell. Rev.* **8**, 1995, 371–391.
32. J. K. Tsotsos, J. Mylopoulos, H. D. Covvey, and S. W. Zucker, A framework for visual motion understanding, *IEEE Trans. Pattern Anal. Mach. Intell.* **2**(6), 1980, 563–573.