

# Better Vision Through Manipulation

Giorgio Metta

Paul Fitzpatrick

LIRA-Lab, DIST

Artificial Intelligence Lab

University of Genova

Massachusetts Institute of Technology

Genova, Italy

Cambridge, MA, USA

`pasa@dist.unige.it`

`paulfitz@ai.mit.edu`

March 13, 2003

Address correspondence to Giorgio Metta, LIRA-Lab, DIST – University of Genova

Viale Causa, 13 – I-16145, Genova, Italy

(phone: +39 0103532791, fax: +39 0103532948).

## Abstract

Vision and manipulation are inextricably intertwined in the primate brain. Tantalizing results from neuroscience are shedding light on the mixed motor and sensory representations used by the brain during reaching, grasping, and object recognition. We now know a great deal about *what* happens in the brain during these activities, but not necessarily *why*. Is the integration we see functionally important, or just a reflection of evolution's lack of enthusiasm for sharp modularity? We wish to instantiate these results in robotic form to probe their technical advantages and to find any lacunae in existing models. We begin with a precursor to manipulation, simple poking and prodding, and show how it facilitates object segmentation, a long-standing problem in machine vision. The robot can familiarize itself with the objects in its environment by acting upon them. It can then recognize other actors (such as humans) in the environment through their effect on the objects it has learned about. We argue that following causal chains of events out from the robot's body into the environment allows for a very natural developmental progression of visual competence, and relate this idea to results in neuroscience.

**keywords:** humanoid robotics, active segmentation, epigenesis

**running title:** Vision and Manipulation

# Better Vision Through Manipulation

Giorgio Metta

Paul Fitzpatrick

LIRA-Lab, DIST

Artificial Intelligence Lab

University of Genova

Massachusetts Institute of Technology

Genova, Italy

Cambridge, MA, USA

pasa@dist.unige.it

paulfitz@ai.mit.edu

March 13, 2003

## Abstract

Vision and manipulation are inextricably intertwined in the primate brain. Tantalizing results from neuroscience are shedding light on the mixed motor and sensory representations used by the brain during reaching, grasping, and object recognition. We now know a great deal about *what* happens in the brain during these activities, but not necessarily *why*. Is the integration we see functionally important, or just a reflection of evolution's lack of enthusiasm for sharp modularity? We wish to instantiate these results in robotic form to probe their technical advantages and to find any lacunae in existing models. We believe it would be missing the point to investigate this on a platform where dextrous manipulation and sophisticated machine vision are already implemented in their mature form, and instead follow a developmental approach from simpler primitives.

We begin with a precursor to manipulation, simple poking and prodding, and show how it facilitates object segmentation, a long-standing problem in machine vision. The robot can familiarize itself with the objects in its environment by acting upon them. It can then recognize other actors (such as humans) in the environment through their effect on the objects it has learned about. We argue that following causal chains of events out from the robot's body into the environment allows for a very natural developmental

# 1 Vision, action, and development

Robots and animals are actors in their environment, not simply passive bystanders. They have the opportunity to examine the world using causality, by performing probing actions and learning from the response. Tracing chains of causality from motor action to perception (and back again) is important both to understand how the brain deals with sensorimotor coordination and to implement those same functions in an artificial system, such as a humanoid robot (Sperber et al., 1995). In this paper, we propose that such causal probing can be arranged in a developmental sequence leading to a manipulation-driven representation of objects. We present results for many important steps along the way, and describe how they fit in a larger scale implementation. And we discuss in what sense our artificial implementation is substantially in agreement with neuroscience.

Table 1 shows four levels of causal complexity that we address in the paper. The simplest causal chain that an actor – whether robotic or biological – may experience is the perception of its own actions. The temporal aspect is immediate: visual information is tightly synchronized to motor commands. Once this causal connection is established, we can go further and use it to actively explore the boundaries of objects. In this case, there is one more step in the causal chain, and the temporal nature of the response may be delayed since initiating a reaching movement doesn’t immediately elicit consequences in the environment. Finally we argue that extending this causal chain further will allow the actor to make a connection between its own actions and the actions of another. This is reminiscent of what has been observed in the response of the primate’s premotor cortex.

Taken together these observations from neuroscience suggest a critical role for motor action in perception. Certainly vision and action are intertwined at a very basic level. While an experienced adult can interpret visual scenes perfectly well without acting upon them, linking action and perception seems crucial to the developmental process that leads to that competence. We can construct a working hypothesis: that action is

required for object recognition in cases where an agent has to develop categorization autonomously. Of course in standard supervised learning action is not required since the trainer does the job of pre-segmenting the data by hand. In an ecological context, some other mechanism has to be provided. Ultimately this mechanism is the body itself that through action (under some suitable developmental rule) generates informative percepts. This notion is related to the “toil versus theft” distinction used by Harnad (Harnad, 2002). Harnad points out that although the meaning of concepts must eventually be traced back to experience (“toil”), evolution and communication provide a way to bypass this through genetic or social “theft”. Human infants, for example, exhibit significant perceptual abilities before their motor skills have developed fully. Nevertheless, they clearly “detect object properties with increasing specificity in relation to their own emerging action capabilities” (Adolph et al., 1993). In our robotic experiments, we seek to trace a causal path all the way from the perception and exploitation of object affordances back to a very minimal set of sensor and motor primitives. When seeking analogues between this process and development in humans or other primates, it is important to bear in mind that some logically required steps may be subsumed by the animal’s evolutionary legacy.

We can distinguish three main conceptual functions in the developmental process that leads to object representation (similar to the schema of Arbib et al. (Arbib, 1981)): reaching, grasping (manipulation), and object recognition. These functions correspond to the levels of causal understanding introduced in Table 1. They also form an elegant progression of abilities which emerge out of very few initial assumptions. All that is required is the interaction between the actor and the environment, and a set of appropriate developmental rules specifying what information is retained during the interaction, the nature of the sensory processing, the range of motor primitives, etc. If we consider the actual localization of functions in the brain we can observe a developmental sequence roughly following a dorsal to ventral gradient. Unfortunately this is a question which has not yet been investigated in detail by neuroscientists, and there is very little empirical support for this claim (beside the work of Kovacs et al. (Kovacs, 2000)).

What is certainly true is that the three modules/functions can be clearly identified. If our hypothesis is correct then the first developmental step has to be that of transporting the hand close to the object (that

we numbered level 1 in our concise description of table 1). In humans, this function is accomplished mostly by the circuit VIP/7b-F4-F1 (see also figure 2). Reaching requires at least the detection of the object and hand, and the transformation of their positions into appropriate motor commands. Parietal neurons seem to be coding for the spatial position of the object in non-retinotopic coordinates by taking into account the position of the eyes with respect to the head. According to (Pouget et al., 2002) and to (Flanders et al., 1999) the gaze direction (the eye motor plant) seems to be the privileged reference system used to code reaching. Relating to the description of causality, the link between an executed motor action and its visual consequences can be easily formed by a subsystem that can detect causality in a short time frame (the immediate aspect).

Once reaching is reliable enough, we can start to move our attention outwards onto objects (identified as level 2 in table 1). Area AIP (parietal lobe) and F5 (frontal cortex) are involved in the control of grasping and manipulation. F5 talks to the primary motor cortex for the fine control of movement. The AIP-F5 system responds to the “affordances” of the observed object with respect to the current abilities. Arbib and coworkers (Fagg and Arbib, 1998) proposed the FARS model as a possible description of the computation in AIP/F5. They did not however consider how affordances can be actually learned during interaction with the environment. Learning and understanding affordances requires a slightly longer time frame since the initiation of an action (motor command) does not immediately elicit all relevant sensory consequences. In this example, the initiation of reaching requires a mechanism to detect when an object is actually touched, manipulated, and whether the collision/touch is causal to the initiation of the movement.

The next step along this hypothetical developmental route is to acquire the F5 mirror representation. We might think of canonical neurons as an association table of grasp/manipulation (action) types with object (vision) types. Mirror neurons can then be thought of as a second-level associative map which links together the observation of a manipulative action performed by somebody else with the neural representation of one’s own action. Mirror neurons bring us to an even higher level of causal understanding (level 3). In this case the action execution has to be associated with a similar action executed by somebody else. The two events do not need to be temporally close to each other. Arbitrary time delays might occur.

The conditions for when this is feasible are a consequence of active manipulation. During a manipulative act there are a number of additional constraints that can be factored in to simplify perception/computation. For example, detection of useful events is simplified by information from touch, by timing information about when reaching started, and from a knowledge of the location of the object.

The last subsystem to develop is object recognition (level 4). Object recognition can build on manipulation in finding the boundaries of objects and segmenting them from the background. More importantly, once the same object is manipulated many times the brain can start learning about the criteria to identify the object if it happens to see it again. These functions are carried out by the infero-temporal cortex (IT). The same considerations apply to the recognition of the manipulator (either one's own, or another's). In fact, the STs region is specialized for this task. Information about object identity is also sent to the parietal cortex and contributes to the formation of the affordances. However object recognition is performed, at a minimum all information (visual in this case) pertaining to a certain object needs to be grouped during development so that a model of the object can be constructed.

[Table 1 about here.]

For the robotic implementation we endeavor to follow the same developmental pathway and exploit the same sort of causal links between actions and sensory feedback. Also, we wish to instantiate these results in robotic form to probe their technical advantages and to find any lacunae in existing models.

We wished to keep the actions implemented on our robotic system as simple as possible, to avoid obscuring the core issue of development behind an elaborately engineered dextrous system. We found that simple poking gestures (prodding, tapping, swiping, batting, etc.) were rich enough to evoke object affordances such as rolling. They also provided exactly the kind of training data needed to bootstrap perception, since they facilitated “active segmentation”, where the motion of the object generated by the robot served to identify its boundaries.

[Figure 1 about here.]

## 2 Object or illusion?

Following (Manzotti and Tagliasco, 2001), we can ask whether macroscopic objects exist completely in their own right, or instead owe something of their existence to their interaction with an observer. How the world is divided up, and what parts of it we grant status as objects, says as much about us as about the world around us (Hendriks-Jansen, 1996). For example, would a chair still be a chair if we had a completely different embodiment? Further, even if a part of the physical world could be separated out from the background in an objective manner, its function still depends on our body and skills – for example, a floppy disk is of little use to one who is computer illiterate, and perhaps can be just regarded as a clumsy frisbee or ugly drink coaster.

Consider the example in Figure 1. It is clear that the cross on the left is a cross and does not seem to owe its existence to us as observers. The array in the middle for many of us is still a cross. This would still be the case even if we had not developed the concept of number or these particular graphic symbols to identify numbers. What can we say about the array on the right? On a first examination it looks like a random collection of numbers. But if we are told that the criterion is “prime numbers vs. non-prime” then a cross can still be identified.

On the very right of figure 1 we show a cube sitting on the table. While humans are very good in analyzing scenes like this one, there are many features that can fool a computer vision system. The edges of the cube and table happen to be aligned, the color is poorly separated, and the surface pattern of the cube does not really tell much about the object itself. Is the internal dark square a different object lying on top of the cube? Another possibility is that the cube is extremely heavy or even part of the table and thus it is not manipulable or movable. Does it make sense then to speak about objects in images, as if there were a unique correspondence between the two? As early as 1734, Berkeley observed that:

...objects can only be known by touch. Vision is subject to illusions, which arise from the distance-size problem... (Berkeley, 1972)

Vision is indeed subject to many illusions. But touch also can be fooled since it has been shown that vision



and touch combine optimally with respect to a maximum likelihood criterion (Ernst and Banks, 2002). Which sensory modality dominates depends on the experimental conditions and apparently we shouldn't always "blindly" trust our senses. The key to resolving ambiguity is to take action, rather than remain a passive observer. In the remainder of the paper we argue that in the presence of manipulation – even a simple form of manipulation – vision becomes more powerful and many of its illusions fade away.

### 3 Objects and action in humans

The example of the cross composed of prime numbers is a novel (albeit unlikely) type of segmentation in our experience as adult humans. We might imagine that when we were very young, we had to initially form a set of such criteria to solve the object identification/segmentation problem in more mundane circumstances. That such abilities develop and are not completely innate is suggested by results in neural science. For example Kovacs (Kovacs, 2000) has shown that perceptual grouping is slow to develop and continues to improve well beyond early childhood (14 years). Long-range contour integration was tested and this work elucidated how this ability develops to enable extended spatial grouping.

A useful concept to understand how such capabilities could develop is the well-known theory of Ungerleider and Mishkin (Ungerleider and Mishkin, 1982) who first formulated the hypothesis that the brain's visual pathways split into two main streams: the dorsal and the ventral. The dorsal is the so-called "where" pathway, concerned with the analysis of the spatial aspects of motor control. The ventral is related with the "what", i.e. the identity of objects.

Goodale and Milner (Milner and Goodale, 1995) refined the theory by proposing that objects are represented differently during action than they are for a purely perceptual task. The dorsal deals with the information required for action, while the ventral is important for more cognitive tasks such as maintaining an object's identity and constancy. Although the dorsal/ventral segregation is emphasized by many commentators, it is significant that there is a great deal of cross talk between the streams. Observation of agnosic patients (Jeannerod, 1997) shows a much more complicated relationship than the simple dor-

sal/ventral dichotomy would suggest. For example, although some patients could not grasp generic objects (e.g. cylinders), they could correctly preshape the hand to grasp known objects (e.g. a lipstick): interpreted in terms of the two pathways, this implies that the ventral representation of the object can supply the dorsal stream with size information.

[Figure 2 about here.]

Grossly simplifying, the brain circuitry responsible for object oriented actions is thought to consist of at least four interacting regions (Figure 2), namely the primary motor cortex (F1), the premotor cortex (F4, F5), the inferior parietal lobule (AIP, VIP), and the temporal cortex (TE, TEO) (see (Rizzolatti et al., 1997; Fadiga et al., 2000; Jeannerod, 1997) for a review). While this is a useful subdivision, it is worth bearing in mind that the connectivity of the brain is much more complex, that bidirectional connections are present, and that behavior is the result of a population activity of these areas. The example about the grasping of known objects in agnosic patients testifies to the *abundance of anatomical connections* between different regions (Jeannerod et al., 1995).

Another way of looking at the same connectivity is in terms of the main function of each area. For example F4, VIP, and 7b are involved in the control of reaching, F5 and AIP contain the majority of grasp related neurons, while TE and TEO are thought to subserve object recognition. These regions together form a network of parallel and yet interacting processes. In fact, at the behavioral level, it has been observed that reaching and grasping need to interact to correctly orient and preshape the hand (Jeannerod et al., 1995).

Neurons responsive to reaching are present in the inferior parietal lobule. For example, Jeannerod et al. reported that the temporary inactivation of the caudal part (VIP) of the intraparietal sulcus by injecting a GABA agonist disrupts reaching (Jeannerod et al., 1995). Conversely, injection in the more rostral part (area AIP) interferes with the preshaping of the hand.

Some of the VIP neurons have bimodal visual and somatic receptive fields (RF). About 30% of them have a RF which does not vary with movement of the head (Rizzolatti et al., 1997). The tactile and visual RF often overlap (e.g. a central visual RF corresponds to a tactile RF in the nose or mouth). The parietal cortex also contains cells related to eye position/movements that appear to be involved in the visuo-motor

transformation required for reaching. VIP projects to area F4 in the premotor cortex. Area F4 contains neurons that respond to objects and are related to the description of the peripersonal space with respect to reaching (Graziano et al., 1997b; Fogassi et al., 1996). A subset of the F4 neurons have a somatosensory, visual, and motor receptive field. The visual receptive field extends in 3D from a given body part, such as the forearm. The somatosensory RF is usually in register with the visual one (as in VIP neurons). Motor information is integrated into the representation by maintaining the receptive field anchored to the correspondent body part (the forearm in this example) irrespective of the relative position of the head and arm.

Also, Graziano et al. (Graziano et al., 1997a) described neurons that maintain a memory of the position of objects for the purpose of reaching. They found neurons that change their firing rate after an object is illuminated briefly within reaching distance. The neurons return to their baseline firing rate only after the monkey is shown that the object have been taken away or moved to a different position.

Sakata and coworkers (Sakata et al., 1997) investigated the response of neurons in the parietal cortex and in particular in area AIP (anterior intra-parietal). They found cells responsive to complex visual stimuli. Neurons in AIP responded during grasping/manipulative actions and when an object was presented to the monkey but no reaching was allowed. Neurons were classified as motor dominant, visual dominant or visuo-motor type depending on how they fired in the dark. Of the visual dominant neurons, some responded to the presentation of the object alone and often they were very specific to the size and orientation of the object, others to the type of object, while yet others responded indifferently to the presentation of a broad class of objects. Area AIP is interesting because it contains both motor and visually responsive cells intermixed in various proportions; it can be thought of as a visuo-motor vocabulary for controlling object directed actions. It is also interesting because projections from AIP terminate in the agranular frontal cortex. For many years, because of the paucity of data, this part of the cortex was considered a unitary motor control area. Recent studies (see (Jeannerod, 1997; Fadiga et al., 2000)) have demonstrated that this is not the case. Particularly surprising was the discovery of visual responsive neurons. A good proportion of them have both visual/sensory and motor responses. Area F5, one of the main targets of the projection from AIP (to which

it sends back recurrent connections), was thoroughly investigated by Rizzolatti and colleagues (Gallese et al., 1996).

F5 neurons can be classified in at least two different categories: canonical and mirror. Canonical and mirror neurons are indistinguishable from each other on the basis of their motor responses; their visual responses however are quite different. The canonical type is active in two situations: i) when grasping an object and ii) when fixating that same object. For example, a neuron active when grasping a ring also fires when the monkey simply looks at the ring. This could be thought of as a neural analogue of the “affordances” of Gibson (Gibson, 1977). However, given the heavy projection from AIP, it is not entirely true that the affordances are fully described/computed by F5 alone. A more conservative stance is that the system of AIP, F5, and other areas (such as TE) participate in the visual processing and motor matching required to compute the affordances of a given object.

The second type of neuron identified in F5, the mirror neuron (Fadiga et al., 2000), becomes active under either of two conditions: i) when manipulating an object (e.g. grasping it, as for canonical neurons), and ii) when watching someone else performing the same action on the same object. This is a more subtle representation of objects, which allows and supports, at least in theory, mimicry behaviors. In humans, area F5 is thought to correspond to Broca’s area; there is an intriguing link between gesture understanding, language, imitation, and mirror neurons (Rizzolatti and Arbib, 1998).

The superior temporal sulcus region (STs) and parts of TE contain neurons that are similar in response to mirror neurons (Perrett et al., 1990). They respond to the sight of the hand; the main difference compared to F5 is that they lack the motor response. It is likely that they participate in the processing of the visual information and then communicate with F5 (Gallese et al., 1996) most likely via the parietal cortex.

A possible developmental explanation of the acquisition of these functions can be framed in terms of tracing/interpreting chains of causally related events. The ability to probe longer chains triggers the emergence of new functionality and/or a new set of behaviors. The next sections delve deeper into this proposal for the ontogenesis of object oriented action and provides experimental results of many steps towards this goal.

## 4 The experimental platform

This work is implemented on the robot Cog, an upper torso humanoid (Brooks et al., 1999; Adams et al., 2000). The robot has previously been applied to tasks such as visually-guided pointing (Marjanović et al., 1996), and rhythmic operations such as turning a crank or driving a slinky (Williamson, 1998). Cog has two arms, each of which has six degrees of freedom – two per shoulder, elbow, and wrist. The joints are driven by series elastic actuators (Williamson, 1995) – essentially a motor connected to its load via a spring (think strong and torsional rather than loosely coiled). The arm is not designed to enact trajectories with high fidelity. For that a very stiff arm is preferable. Rather, it is designed to perform well when interacting with a poorly characterized environment, where collisions are frequent and informative events.

[Figure 3 about here.]

The following sections 5 through 9 explore the four levels of causation which are at the core of our working hypothesis. The rationale of the experiments is to show that one possible route to object recognition goes through the “understanding” of longer chains of cause-effects relationships. In particular section 5 describes the simplest causal chain where motion of the robot causes immediate visual effects. Simple cross-correlation over time of motor and visual signals allows localizing the robot’s end-point. Reaching is seen as an extension of the same mechanism. Subsequently, we show in section 6 how the robot explores its peripersonal space and get to explore physical objects. Exploiting causality leads to object segmentation (figure/ground separation). In this case there is a potentially delayed effect because initiating the reaching action does not automatically lead to the interaction with the object. The experiments described in section 7 and 8 build on top of the segmentation to learn object “affordances”. Exploring further a complex causal chain where the actions of others are considered moves us naturally to a “mirror neuron” like response. Eventually, object recognition and an empirical definition of “objecthood” are presented in section 9. Our definition relies on a combination of object affordances and the acquisition of data through multiple instances of the same manipulative act.

## 5 Perceiving direct effects of action

Motion of the arm may generate optic flow directly through the changing projection of the arm itself, or indirectly through an object that the arm is in contact with. While the relationship between the optic flow and the physical motion is likely to be extremely complex, the correlation in time of the two events will generally be exceedingly precise. This time-correlation can be used as a “signature” to identify parts of the scene that are being influenced by the robot’s motion, even in the presence of other distracting motion sources. In this section, we show how this tight correlation can be used to localize the arm in the image without any prior information about visual appearance. This is in fact why we chose to detect the arm using optic flow rather than by searching for a predetermined color or shape. In the next section we will show that once the arm has been localized we can go further, and identify the boundaries of objects with which the arm comes into contact.

A similar procedure was described by Piaget as “circular reaction” (Piaget, 1963). In Piaget’s observations the circular reaction is the mechanism by which the loop between vision and action is closed. In the child, this seemingly random activity mediates the discovery of contingent activation of visual, motor, and somatosensory areas. Other researchers (Bullock et al., 1993) applied a similar model in learning visuomotor transformations. The “motor babbling” activity was used to self-train the sensori-motor transformations required for reaching. We instantiate here a similar mechanism to learn to localize the robot effector.

### Reaching out

The first step towards manipulation is to reach objects within the workspace. If we assume targets are chosen visually, then ideally we need to also locate the end-effector visually to generate an error signal for closed-loop control. Some element of open-loop control is necessary since the end-point may not always be in the field of view (for example, when it is in its the resting position), and the overall reaching operation can be made faster with a feed-forward contribution to the control.

[Figure 4 about here.]

The simplest possible open loop control would map directly from a fixation point to the arm motor commands needed to reach that point (Metta et al., 1999) using a stereotyped trajectory, perhaps using postural primitives (Mussa-Ivaldi and Giszter, 1992). If we can fixate the end-effector, then it is possible to learn this map by exploring different combinations of direction of gaze vs. arm position (Marjanović et al., 1996; Metta et al., 1999). So locating the end-effector visually is key both to closed-loop control, and to training up a feed-forward path. We shall demonstrate that this localization can be performed without knowledge of the arm’s appearance, and without assuming that the arm is the only moving object in the scene.

### Localizing the arm visually

The robot is not a passive observer of its arm, but rather the initiator of its movement. This can be used to distinguish the arm from parts of the environment that are more weakly affected by the robot. The arm of a robot was detected in (Marjanović et al., 1996) by simply waving it and assuming it was the only moving object in the scene. We take a similar approach here, but use a more stringent test of looking for optic flow that is correlated with the motor commands to the arm. This allows unrelated movement to be ignored. Even if a capricious engineer were to replace the robot’s arm with one of a very different appearance, and then stand around waving the old arm, this detection method will not be fooled.

The actual relationship between arm movements and the optic flow they generate is complex. Since the robot is in control of the arm, it can choose to move it in a way that bypasses this complexity. In particular, if the arm rapidly reverses direction, the optic flow at that instant will change in sign, giving a tight, clean temporal correlation. Since our optic flow processing is coarse (computed by a generic correlation-based approach over a  $16 \times 16$  grid over a  $128 \times 128$  image at 15 Hz), we simply repeat this reversal a number of times to get a strong correlation signal during training. With each reversal the probability of correlating with unrelated motion in the environment goes down. This probability could also be reduced by higher resolution (particularly in time) visual processing.

Figure 4 shows an example of this procedure in operation, comparing the velocity of the arm’s wrist with

the optic flow at two positions in the image plane. A trace taken from a position away from the arm shows no correlation, while conversely the flow at a position on the wrist is strongly different from zero over the same period of time. Figure 4 shows examples of detection of the arm and rejection of a distractor.

### Localizing the arm using proprioception

The localization method for the arm described so far relies on a relatively long “signature” movement that would slow down reaching. This can be overcome by training up a function to estimate the location of the arm in the image plane from proprioceptive information (joint angles) during an exploratory phase, and using that to constrain arm localization during actual operation.

The response of such a filter is not too distant from that of the monkey’s parietal and frontal cortices. In particular we already described, in section 3, neurons that respond to the sight of a body part (e.g. the hand) irrespective of the relative position of the eyes, head, and arms.

[Figure 5 about here.]

As a function approximator we simply fill a look-up table, implemented as a list of nodes allocated dynamically. This implementation was chosen to reduce memory consumption; the input space is six dimensional and even a coarse discretization of this space would require memory in the order of several Mbytes. Rather than using all the joint angles the current direction of gaze is first coded in terms of only two angles representing the global pan ( $\theta$ ) and tilt ( $\phi$ ) of one of the cameras. This is easily computed from the kinematics of the head and the joint angles. The end-point position is coded considering only the first four joints ( $q_1 \dots q_4$ ). The position of joint  $q_5$  and  $q_6$  is not employed because the wrist does not significantly contribute to the end-point position. The output of the approximator is the position of the end-point (the forearm) on the image plane. Figure 5 shows the resulting behavior after about twenty minutes of real-time learning.

### Reaching for the object

Reaching is implemented as a direct mapping between the direction of gaze ( $\theta, \phi$ ) and the command required to reach the fixation point. Inspiration is drawn from the experiments of (Flanders et al., 1999) and (Pouget



et al., 2002) that have shown that both humans and monkeys employ gazing as a reference for reaching. This procedure is consistent because we are interested in reaching a point on a plane in front of the robot (a table): i.e. each point on the table is identified by one and only one gaze value. The resulting map is thus  $2D \rightarrow 6D$  (the arm has 6 degrees of freedom). The same argument could be extended to the 3D case by augmenting the encoding of gaze with, for example, the vergence angle. The arm motor commands are represented in terms of joint positions, and the mapping is linear:

$$\begin{pmatrix} \hat{q}_1 \\ \vdots \\ \hat{q}_6 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ \vdots & \vdots \\ a_{61} & a_{62} \end{pmatrix} \cdot \begin{pmatrix} \theta \\ \phi \end{pmatrix} \quad (1)$$

where  $\hat{q}$  are the desired joint positions. The coefficients  $a_{nm}$  are estimated following a brief calibration procedure from a small number of training pairs of the form  $(\hat{\mathbf{q}}, (\theta, \phi))$ . The linear approximation is justified in our case because of the relatively small region of the workspace where the reaching is expected to operate. The complete robot workspace is much bigger because the torso can also move to keep the operational point of the linear approximation within reasonable limits.

Redundancy is not an issue since the vector of joint angles  $\hat{q}$  spans only a 2D subspace of the full 6D space. This subspace is determined by the  $a_{nm}$  and it is uniquely indexed by the input vector  $(\theta, \phi)$ . Hence the mapping allows the robot to reach a point on a particular plane in space consistently. We verified this to be true empirically for a large number of configurations.

At a lower level a low-stiffness position control and a simple trajectory generator interpolate the motion of the arm from the current position to the commanded one. Gravity compensation for the shoulder joint has been implemented to further improve accuracy.

[Figure 6 about here.]

## 6 Perceiving indirect effects of action

We have assumed that the target of a reaching operation is chosen visually. As discussed in the introduction, visual segmentation is not easy, so we should not expect a target selected in this way to be correctly segmented. For the example scene in Figure 1 (a cube sitting on a table), the small inner square on the cube’s surface pattern might be selected as a target. The robot can certainly reach towards this target, but grasping it would prove difficult without a correct estimate of the object’s physical extent. In this section, we develop a procedure for refining the segmentation using the same idea of correlated motion used earlier to detect the arm.

When the arm enters into contact with an object, one of several outcomes are possible. If the object is large, heavy, or otherwise unyielding, motion of the arm may simply be resisted without any visible effect. Such objects can simply be ignored, since the robot will not be able to manipulate them. But if the object is smaller, it is likely to move a little in response to the nudge of the arm. This movement will be temporally correlated with the time of impact, and will be connected spatially to the end-effector – constraints that are not available in passive scenarios (Birchfield, 1999). If the object is reasonably rigid, and the movement has some component in parallel to the image plane, the result is likely to be a flow field whose extent coincides with the physical boundaries of the object.

[Figure 7 about here.]

Figure 6 shows how a “poking” movement can be used to refine a target. During a poke operation, the arm begins by extending outwards from the resting position. The end-effector (or “flipper”) is localized as the arm sweeps rapidly outwards, using the heuristic that it lies at the highest point of the region of optic flow swept out by the arm in the image (the head orientation and reaching trajectory are controlled so that this is true). The arm is driven outward into the neighborhood of the target which we wish to define, stopping if an unexpected obstruction is reached. If no obstruction is met, the flipper makes a gentle sweep of the area around the target. This minimizes the opportunity for the motion of the arm itself to cause confusion; the motion of the flipper is bounded around the endpoint whose location we know from tracking during the

extension phase, and can be subtracted easily. Flow not connected to the end-effector can be ignored as a distractor. For simplicity, the head is kept steady throughout the poking operation, so that simple image differencing can be used to detect the presence of motion at a higher resolution than optic flow. Figure 7 shows an example of the kind of results that are possible (see Section 9 for further examples).

The poking operation gives clear results for a rigid object that is free to move. What happens for non-rigid objects and objects that are attached to other objects? Here the results of poking are likely to be more complicated to interpret – but in a sense this is a good sign, since it is in just such cases that the idea of an object becomes less well-defined. Poking has the potential to offer an operational theory of “objecthood” that is more tractable than a vision-only approach might give, and which cleaves better to the true nature of physical assemblages. The idea of a physical object is rarely completely coherent, since it depends on where you draw its boundary and that may well be task-dependent. Poking allows us to determine the boundary around a mass that moves together when disturbed, which is exactly what we need to know for manipulation. As an operational definition of object, this has the attractive property of breaking down into ambiguity in the right circumstances – such as for large interconnected messes, floppy formless ones, liquids, and so on.

## 7 Experimenting with object affordances

Poking moves us one step outwards on a causal chain away from the robot and into the world, and gives a simple experimental procedure for segmenting objects. There are many possible elaborations of this method, all of which lead to a vision system that is tuned to acquiring data about an object by seeing it manipulated by the robot.

This kind of active segmentation will nevertheless be inconvenient in many situations if not coupled with a mechanism to learn from experience. For example, it would be terribly inefficient to always have to poke an object first before it could be grasped. It would be much better if the robot could learn about objects and, in particular, how to identify a previously encountered object. A further difficulty, at least for a robot with a simple manipulator (such as Cog’s flipper), is that “affordances” are scarce: most of the time the

object will simply move from one position to another if we are willing to discount when it falls from the table.

However, for objects that roll there is a cue the robot can exploit to understand their behavior. An object that rolls tends to do so even if it is not poked precisely. We selected a small set of objects to experiment with: a cube, a toy car, an orange juice bottle, and a ball. Affordances are not only a property of the mechanics of the object, but rather a combination of visual appearance, of the object’s physical composition, and of the ability of the actor. We selected a measure of the principal axis of the object (easily obtained from the segmentation) as a visual component of the affordance. Table 2 shows the expected behavior.

[Table 2 about here.]

We need to group the data belonging to the same object obtained across many poking acts into coherent clusters. We adopted simple color histogram similarity as our clustering criterion. After each poking action, a color histogram of the pixels in the segmented region is built and used to judge whether the object belongs to an existing group (for example, if it is mostly yellow, it is likely to be the toy car). This works well for a small set of objects but more sophisticated methods would be required for a more general case with a large set of objects (Schiele and Crowley, 2000). The data structure that simulates the AIP-F5 affordance computation maintains all the instances of poking grouped by object, all the prototypes of the segmented object, the direction of movement, and the action applied by the robot in each trial.

An alternative to the vision-based clustering procedure would be to try to classify the behavior of an object after a *single* encounter, and use the behavior itself as a clustering criterion. So how an object rolls could be used as a feature to recognize that object. Adopting this strategy would have made our results much more sensitive to the performance of the motor and vision system, since we cannot average over the noise they generate. Nevertheless, this would be a perfectly reasonable strategy for a next-generation system to adopt.

Figure 8 shows the results of the segmentation, clustering and estimation of the affordance of the same set of four objects. The training set consists of about 100 actions per object. The motor vocabulary of the robot consists of four possible directions of poking. We labeled them for convenience as: pull in, side tap,

push away, and back slap, depending on the effect they have on the object from the point of view of the robot. Actions were generated at random during this training stage. During a poking action, the object is tracked for 12 frames after the time of contact and the overall displacement is computed.

[Figure 8 about here.]

This description of the affordances shows clear differences between the objects. But it is not yet an *effective* description since it does not by itself tell the robot how to take action once an object is observed. For this purpose a description of the geometry of poking is required. This information can be derived from the same training set we collected for learning about rolling. Figure 9 shows the histograms of the direction of movement averaged over all objects for each possible action. For example, the back slap moves an object mostly upward (about  $-100^\circ$  on average,  $0^\circ$  being the direction parallel to the image  $x$  axis) and away from the robot. A similar consideration applies to the other poking gestures. Figure 9 was obtained from the data of about 500 poking events.

The last step is to connect all these elements together. If a known object is presented to Cog, the object is recognized, localized, and its orientation estimated (by finding its principal axis). Recognition is based on the color histograms. The same procedure used to form the clusters is employed here. Localization is simply implemented by histogram back-projection and a search across the image. The current orientation of the object is then estimated by comparing the current image with all the prototypes contained in the cluster. The whole procedure has an error on the estimation of the principal axis in the range of  $10^\circ$  to  $25^\circ$  depending on the object.

To actually exploit the understanding of the affordance we need to connect vision to behavior. The robot looks for the preferred rolling direction of the object (see figure 8) and adds it to its current orientation. The action whose effects are closer (on average) to the combination of the orientation and affordance is selected.

We performed a simple qualitative test of the robot's behavior presenting randomly two of the objects (the toy car and the bottle) - note that the ball and the cube do not have a well defined principal axis so there is no point in running the experiment. Out of 100 trials the robot made 15 mistakes. Analysis of the

errors reveals that they are mainly due to imprecise control (12) and to a less extent to misinterpretation of the orientation of the object (3).

[Figure 9 about here.]

## 8 Developing mirror neurons

An interesting question then is whether the system could extract useful information from seeing an object manipulated by someone else. In the case of poking, the robot needs to be able to estimate the moment of contact and to track the arm sufficiently well to distinguish it from the object being poked. We are interested in how the robot might learn to do this. One approach is to chain outwards from an object the robot has poked. If someone else moves the object, we can reverse the logic used in poking – where the motion of the manipulator identified the object – and identify a foreign manipulator through its effect on the object. The next experiment was designed to explore this aspect.

In fact, the same processing used for analyzing an active poking can be used to detect a contact and segment the object from the manipulator. This is not different from what we used for learning. While one might argue then that learning can be carried out just by mere observation, it is worth noting that: i) this situation is not as well defined as the active one, and ii) there is no connection to the motor aspects of the action and consequently it is difficult to link the observation to the behavior. There is no physical contact, thus there is plenty of room for getting confused by false positives. The temporal aspect, so well constrained during active manipulation, is more vague here – the robot, for example, does not know when the foreign manipulator starts or stops the action. If missing a contact event or getting a false or mistaken segmentation is not much of a problem in “observation mode”, it is much more troublesome if we corrupt the training data with unreliable/noisy observations. Further, we should not assume the human “teacher” is truly collaborative. There is no guarantee that actions suited to the robot perceptual system and/or goal are performed at all. More seriously, the link to behavior is completely missing. Even if visual information about objects can be collected as before, tracing back which action causes a particular consequence cannot

be autonomously learned by the robot. Conversely, in the case the robot has already learned about objects, as e.g. we have shown in the previous section, this information can be factored in to help the observation of somebody else’s action. Touch (not in Cog) and physical contact are additional bits of information about the ongoing activity.

In our case, if any activity is detected close to the object – measured by the amount of motion in a neighborhood of the fixation point corresponding to the robot’s foveal camera – reaching is inhibited and the whole action observed (assuming there is one at all). An example of human poking is shown in figure 10.

[Figure 10 about here.]

The first obvious thing the robot can do is to identify the action just observed with respect to its motor vocabulary. It is easily done, in this case, by comparing the displacement of the object with the four possible actions and by choosing the action whose effects are closer to the observed displacement. Indeed it allows – even if in this limited setting – recognizing a complex action by interpreting its consequences on the environment. This is orders of magnitude simpler than trying to completely characterize the action in terms of the observed kinematics of the movement. Here, the complexity of the data we need to obtain from the observations is somehow proportional to the complexity of the goal rather than that of the structure/skills of the foreign manipulator. In our case, because the action, the goal, and the object are relatively simple, the only information required is about the displacement of the object.

Therefore, the next question is whether we can use this “understanding” of observed actions to implement mimicry behavior. It would be easy now to try to replicate the action just observed if the same object were presented again. However, there is still a bit of ambiguity in that we can choose to mimic either the observed displacement of the object or the way the object was poked with respect to its rolling affordance.

We chose to implement the latter. It is clear that poking along a particular observed direction requires trivial modifications. In practice, after an action is observed the angle between the affordance (see table 2) and the actual displacement is measured and stored. If it happens to see the same object again, the robot chooses the action that has the greatest probability of poking the object along the previously stored angle. Figures 10 and 11 show examples of such mimicry.

[Figure 11 about here.]

This response is exactly what we would expect from a “mirror-type” representation. The observed action is interpreted on the basis of the robot own motor code. The same data structure is also used/activated when performing an action in response to the sight of a known object. The causal link between the two events that could be separated by several seconds is the object, the goal, and the object’s affordances. There is considerable precedent in the literature for a strong connection between viewing object manipulation performed by either oneself or another (Wohlschläger and Bekkering, 2002). There is also a growing evidence that imitation is goal-directed (Bekkering and Wohlschläger, 2000) and that the object of the action is explicitly coded (e.g. during reaching) (Woodward, 1998).

## 9 Towards object recognition

[Figure 12 about here.]

Although poking is a very crude and primitive form of manipulation we have shown that it can help to bootstrap more complex behaviors without relying on an external teacher. With only minimal assumptions (using motion as segmentation cue) we were able to build a system that exploits its environment to learn novel behaviors. If Cog had a dextrous hand, it could further exploit temporal constraints (e.g. an object remains the same unless it is dropped) to collect tightly/temporally correlated data. There are already examples in robotics of the acquisition of object categorization based on this kind of temporally correlated information (Scheier and Lambrinos, 1996). This form of “object constancy” could be exploited for instance to learn about an object with confusing visual features such as many different colors, different geometric patterns, and so forth (see the example of the cube in figure 13). A finer form of manipulation can be used also to group objects on the basis of their behavior rather than purely by visual appearance: e.g. the class of “bottle” or of “toy cars”. This, in some future implementation, can help the robot to attain a goal by using a suitable tool (among many) rather than exactly the same tool it used when initially learned the task.

A possible and obvious extension is to use the object segmentation provided by poking (and manipu-



lation in general) to build models of the appearance of objects beyond the color histogram we used in our experiments (think again about the colored cube shown in figure 13). Also in this case the robot could work autonomously on learning. Furthermore, the interaction between manipulator and object provides another element that can be used to learn about the manipulator itself (see figure 14). The robot can then learn about the appearance of its own hand or, equally, about the human hand. It is remarkable that the complexity of the robot manipulator does not necessarily have to match that of the human manipulator. We can envision a similar procedure to learn about any object that functions as manipulator.

[Figure 13 about here.]

[Figure 14 about here.]

## 10 Discussion and Conclusions

In this paper, we showed how causality can be probed at different levels by the robot. Initially the environment was the body of the robot itself, then later a carefully circumscribed interaction with the outside world. This is reminiscent of Piaget’s distinction between primary and secondary circular reactions (Ginsburg and Oppen, 1978). Objects are central to interacting with the outside world. We raised the issue of how an agent can autonomously acquire a working definition of objects.

In computer vision there is much to be gained by bringing a manipulator into the equation. Many variants and extensions to the experimental “poking” strategy explored here are possible. For example, a robot might try to move an arm around *behind* the object. As the arm moves behind the object, it reveals its occluding boundary. This is a precursor to visually extracting shape information while actually manipulating an object, which is more complex since the object is also being moved and partially occluded by the manipulator. Another possible strategy that could be adopted as a last resort for a confusing object might be to simply hit it firmly, in the hopes of moving it some distance and potentially overcoming local, accidental visual ambiguity. Obviously this strategy cannot always be used! But there is plenty of room to be creative here. There are also limitations in our current implementation that could usefully be addressed.

The robot itself is not mobile, so its workspace is limited. There are also many constraints on the arm that make fine motor control impossible – it cannot maintain all reachable poses indefinitely, and there is significant noise and some hysteresis in its analog sensors. The robot will only attempt to reach towards a target that is actually accessible to its arm – not too close, not too far, as determined using visual disparity. In practice, this means that the ideal workspace is a table in front of the robot, and the motor control of the robot has been specifically tuned to work well in that situation. A simple attention system and tracking mechanism are used to bring the robot’s attention to a target. This phase can fail if the robot gets distracted by some more salient (but unreachable) part of the scene. Objects that move together are not individually segmented. And segmentation does not always succeed, due to shadows, or strong nearby edges.

In spite of some limitations, the robotic experiments support the view that reaching, grasping, and recognition can be learned by following a particular ontogenetic pathway without the intervention of an external teacher. This pathway is consistent with and inspired by what is known of this process in biological systems (primates/mammals). We have endeavored to build from as few innate components as possible, to elucidate the visual and motor challenges faced by a learning robot rather than simply solving them by fiat. Although newborns show amazing abilities (Spelke, 2000) such as early imitation (Meltzoff and Moore, 1977), face detection, etc, there is also evidence that the maturation of the brain is far from complete at birth and complex perceptual abilities require a long time to emerge (Kovacs, 2000). We have given a simple existence proof that object segmentation, recognition and localization can develop without any prior knowledge of visual appearance. We have also shown that, without any prior knowledge of the human form, the robot can identify episodes when a human is manipulating objects that are familiar to the robot purely by the operational similarity of the human arm and its own manipulator in this situation. We believe such demonstrations are important both in their own right, and in their elucidation of a concrete series of steps that lead to a desired behavior.

Many researchers have shown now examples of the application of developmental principles in the design of autonomous systems, for example (??) and (Metta et al., 1999). This approach may provide novel directions to robotics. Besides, it may also serve as a useful reference point from which to investigate the biological

solution to the same problem – although it can’t provide the answers, it can at least suggest useful questions.

## Acknowledgements

This work benefited from discussions with Charles Kemp, Giulio Sandini, and Luciano Fadiga. Many people have contributed to developing the Cog platform (Brooks et al., 1999). Funds for this project were provided by DARPA as part of the “Natural Tasking of Robots Based on Human Interaction Cues” project under contract number DABT 63-00-C-10102, and by the Nippon Telegraph and Telephone Corporation as part of the NTT/MIT Collaboration Agreement.

## References

- Adams, B., Breazeal, C., Brooks, R., and Scassellati, B. (2000). The Cog project. *IEEE Intelligent Systems*. To appear.
- Adolph, K. E., Eppler, M. A., and Gibson, E. J. (1993). Development of perception of affordances. *Advances in Infancy Research*, 8:51–98.
- Arbib, M. A. (1981). *Handbook of Physiology*, chapter Perceptual Structures and Distributed Motor Control. American Physiological Society.
- Bekkering, H. and Wohlschlager, A. (2000). Imitation in children is goal-directed. *The quarterly journal of experimental psychology*, 53A(1):153–164.
- Berkeley, G. (1972). *A new theory of vision and other writings*. Dent, London. First published in 1734.
- Birchfield, S. (1999). *Depth and Motion Discontinuities*. PhD thesis, Dept. of Electrical Engineering, Stanford University.
- Brooks, R. A., Breazeal, C., Marjanovic, M., and Scassellati, B. (1999). The Cog project: Building a humanoid robot. *Lecture Notes in Computer Science*, 1562:52–87.
- Bullock, D., Grossberg, S., and Guenther, F. H. (1993). A self-organizing model of motor equivalent reaching and tool use by a multijoint arm. *Journal of Cognitive Neuroscience*, (5):408–435.
- Ernst, M. O. and Banks, M. S. (2002). Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415:429–433.
- Fadiga, L., Fogassi, L., Gallese, V., and Rizzolatti, G. (2000). Visuomotor neurons: ambiguity of the discharge of ‘motor’ perception? *International Journal of Psychophysiology*, 35:165–177.
- Fagg, A. H. and Arbib, M. A. (1998). Modeling parietal-premotor interaction in primate control of grasping. *Neural Networks*, 11(7–8):1277–1303.

- Flanders, M., Daghestani, L., and Berthoz, A. (1999). Reaching beyond reach. *Experimental Brain Research*, 126(1):19–30.
- Fogassi, L., Gallese, V., Fadiga, L., Luppino, G., Matelli, M., and Rizzolatti, G. (1996). Coding of peripersonal space in inferior premotor cortex (area F4). *Journal of Neurophysiology*, pages 141–157.
- Gallese, V., Fadiga, L., Fogassi, L., and Rizzolatti, G. (1996). Action recognition in the premotor cortex. *Brain*, 119:593–609.
- Gibson, J. J. (1977). The theory of affordances. In Shaw, R. and Bransford, J., editors, *Perceiving, acting and knowing: toward an ecological psychology*, pages 67–82. Hillsdale NJ: Lawrence Erlbaum Associates Publishers.
- Ginsburg, H. and Oppen, S. (1978). *Piaget’s theory of intellectual development*. Prentice-Hall, Englewood Cliffs, NJ. 2nd edition.
- Graziano, M. S. A., Hu, X., and Gross, C. G. (1997a). Coding the location of objects in the dark. *Science*, 277:239–241.
- Graziano, M. S. A., Hu, X., and Gross, C. G. (1997b). Visuo-spatial properties of ventral premotor cortex. *Journal of Neurophysiology*, 77:2268–2292.
- Harnad, S. (2002). *Computationalism: New Directions*, chapter Symbol grounding and the origin of language, pages 143–158. MIT Press.
- Hendriks-Jansen, H. (1996). *Catching Ourselves in the Act*. MIT Press, Cambridge, Massachusetts.
- Jeannerod, M. (1997). *The Cognitive Neuroscience of Action*. Blackwell Publishers Inc., Cambridge Massachusetts and Oxford UK.
- Jeannerod, M., Arbib, M. A., Rizzolatti, G., and Sakata, H. (1995). Grasping objects: the cortical mechanisms of visuomotor transformation. *Trends in Neurosciences*, 18(7):314–320.
- Kovacs, I. (2000). Human development of perceptual organization. *Vision Research*, 40(10-12):1301–1310.

- Manzotti, R. and Tagliasco, V. (2001). *Coscienza e realtà: una teoria della coscienza per costruttori di menti e cervelli*. il Mulino.
- Marjanović, M. J., Scassellati, B., and Williamson, M. M. (1996). Self-taught visually-guided pointing for a humanoid robot. In *From Animals to Animats: Proceedings of 1996 Society of Adaptive Behavior*, pages 35–44, Cape Cod, Massachusetts.
- Meltzoff, A. N. and Moore, M. K. (1977). Imitation of facial and manual gestures by human neonates. *Science*, 198:75–78.
- Metta, G., Sandini, G., and Konczak, J. (1999). A developmental approach to visually-guided reaching in artificial systems. *Neural Networks*, 12:1413–1427.
- Milner, A. D. and Goodale, M. A. (1995). *The visual brain in action*. Oxford University Press.
- Mussa-Ivaldi, F. A. and Giszter, S. F. (1992). Vector field approximation: a computational paradigm for motor control and learning. *Biological Cybernetics*, 67:491–500.
- Perrett, D. I., Mistlin, A. J., Harries, M. H., and Chitty, A. J. (1990). Understanding the visual appearance and consequence of hand action. In *Vision and action: the control of grasping*, pages 163–180. Ablex, Norwood, NJ.
- Piaget, J. (1963). *The origin of intelligence in children*. Norton, New York, NY.
- Pouget, A., Ducom, J.-C., Torri, J., and Bavelier, D. (2002). Multisensory spatial representation in eye-centered coordinates for reaching. *Cognition*, 83:B1–B11.
- Rizzolatti, G. and Arbib, M. A. (1998). Language within our grasp. *Trends in Neurosciences*, 21:188–194.
- Rizzolatti, G., Fogassi, L., and Gallese, V. (1997). Parietal cortex: from sight to action. *Current Opinion Neurobiology*, 7(4):562–567.
- Sakata, H., Kusunoki, M., Taira, M., Murata, M., and Tanaka, Y. (1997). The tins lecture - the parietal association cortex in depth perception and visual control of action. *Trends in Neurosciences*, 20(8):350–358.

- Scheier, C. and Lambrinos, D. (1996). Categorization in a real-world agent using haptic exploration and active perception. In *Proceedings of SAB96 (FROM ANIMALS TO ANIMATS, Fourth International Conference on Simulation of Adaptive Behavior*, Cape Cod, MA, USA.
- Schiele, B. and Crowley, J. L. (2000). Recognition without correspondence using multidimensional receptive field histograms. *International Journal of Computer Vision*, 36(1):31–50.
- Spelke, E. (2000). Core knowledge. *American Psychologist*, 55(11):145–160.
- Sperber, D., Premack, D., and Premack, A. J. (1995). *Causal cognition: a multisdisciplinary debate*. Oxford University Press, New York, NY.
- Ungerleider, L. G. and Mishkin, M. (1982). Two cortical visual systems. In *Analysis of visual behavior*, pages 549–586. MIT Press, Cambridge, Massachusetts.
- Williamson, M. (1995). Series elastic actuators. Master’s thesis, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA.
- Williamson, M. (1998). Neural control of rhythmic arm movements. *Neural Networks*, 11(7-8):1379–1394.
- Wohlschläger, A. and Bekkering, H. (2002). Is human imitation based on a mirror-neurone system? Some behavioural evidence. *Experimental Brain Research*, 143:335–341.
- Woodward, A. L. (1998). Infant selectively encode the goal object of an actor’s reach. *Cognition*, 69:1–34.

List of Tables

1    Degrees of causal indirection, localization and function in the brain. There is a natural trend  
from simpler to more complicated tasks. The more time-delayed an effect, the more difficult  
it is to model. . . . . 35

2    Behavior of a small set of objects when poked at random by the robot manipulator. . . . . 36



## List of Figures

- 1    On the left are three examples of crosses, following (Manzotti and Tagliasco, 2001). The human ability to segment objects is not general-purpose, and improves with experience. On the right is an image of a cube on a table, illustrating the ambiguities that plague machine vision. The edges of the table and cube happen to be aligned (dashed line), the colors of the cube and table are not well separated, and the cube has a potentially confusing surface pattern. 37
- 2    Monkey brain with indication of the main areas participating in object oriented actions (adapted from (Fagg and Arbib, 1998)). As described in the text, three main functions can be identified: object recognition, reaching, and grasping. These form three parallel yet connected streams of processing. The circuit connecting the visual cortex to the inferior parietal lobule VIP, F4 and F1 is thought to compute the visuomotor transformations required to control reaching. Some evidence also suggests a possible role in the organization of reaching played by the posterior parietal cortex PO and dorsal premotor area F2, reciprocally connected. AIP and F5 are responsible for grasping. Temporal areas (TE, TEO) and STs are correlated to the semantic of object recognition. . . . . 38
- 3    The robot Cog, an upper-torso humanoid. The ultimate goal of this work is for our robot to follow chains of causation outwards from its own simple body into the complex world. Such an incremental process suggests that perception and action develop together, supporting each other. The head, torso, and arms together contain 22 degrees of freedom. . . . . 39
- 4    (a) An example of the correlation between optic flow and arm movement. The traces show the movement of the wrist joint (upper plot) and optic flow sampled on the arm (middle plot) and away from it (lower plot). (b) The robot's point of view and the optic flow generated are shown on the left. On the right are the results of correlation. Large circles represent the results of applying a region growing procedure to the optic flow. The small circle marks the point of maximum correlation, identifying the regions that correspond to the robot's own arm. 40

- 5 Predicting the location of the arm in the image as the head and arm change position. The rectangle represents the predicted position of the arm using the map learned during a twenty-minute training run. The predicted position just needs to be sufficiently accurate to initialize a visual search for the exact position of the end-effector. . . . . 41
- 6 The upper sequence shows an arm extending into a workspace, tapping an object, and retracting. This is an exploratory mechanism for finding the boundaries of objects, and essentially requires the arm to collide with objects under normal operation, rather than as an occasional accident. The lower sequence shows the shape identified from the tap using simple image differencing and flipper tracking. . . . . 42
- 7 An example of the power of active segmentation. The images marked “scene” show two presentations of a yellow toy car sitting on a yellow table. The robot extends its arm across the table. In the upper sequence it strikes from below, in the lower sequence it strikes from the side (“action” images). Once the arm comes in contact with the car, it begins to move, and it can be segmented from the stationary background (“object”). On the left of the figure, a zoomed view of the car/table boundary is shown – the difference between the two is very subtle. . . . . 43
- 8 Probability of observing a roll along a particular direction for the set of four objects used in our experiments. Abscissae represent the difference between the principal axis of the object and the observed direction of movement. Ordinates the estimated probability. . . . . 44
- 9 Histogram of the direction of movement of object for each possible poking action. For each of the four plots the abscissa is the direction of motion of the object where the  $0^\circ$  direction is parallel to the x axis, and  $-90^\circ$  to the y axis. The ordinate is the empirical probability distribution of the direction of motion of the objects. . . . . 45

- 10 Basic mimicry. The first step in mimicking an action is to actually be able to observe it. The first sequence shows a human demonstration of a poking operation. Frames around the moment of contact are shown. The object, after segmentation, is tracked for 12 frames using a combination of template matching and optic flow. The big circles represent the tracked position of the bottle in successive frames. The arrow displayed on the frame of contact ( $3^{rd}$  from the left) projects from the position at the time of contact and at the  $12^{th}$  frame respectively. In the second sequence, the bottle is presented to the robot in the same orientation it had in the demonstrated action and the robot repeats the observed action, a “side tap”. In the third sequence, the car is presented at a different angle. The appropriate action to exploit the affordance and make the bottle roll is now a “back slap”. . . . . 46
- 11 An extended mimicry example using the toy car. The sequences on the left show the robot mimicking a human exploiting the car’s rolling affordance. The sequences on the right show what happens when the human hits the car in a contrary fashion, going against its preferred direction of motion. The robot mimics this “unnatural” action, suppressing its usual behavior of trying to evoke rolling. . . . . 47
- 12 Once objects have been segmented from the background, they are much easier to distinguish from each other since the irrelevant similarity of their shared environment is eliminated. To build object models, the robot clusters all the segmented views it receives based on similarity of their color histogram. This figure shows samples from four of the clusters found, corresponding to the four objects used in Section 7. Note the baseball cap classified with the ball, lower right – a young child wandered by the robot while we were collecting data and got it to poke his cap. . . . . 48
- 13 Poking also gives the robot the opportunity to collect many views of a single object, and so we can hope to deal with recognizing objects like this toy cube that has a different appearance from every side (the segmentations shown here were collected automatically). . . . . 49

	34
14 Early experiments on segmenting the robot arm, or a human hand poking an object the robot	
is familiar with, by working backwards from a collision event. . . . .	50

	<i>nature of causation</i>	<i>main path</i>	<i>function and/or behavior</i>	<i>time profile</i>
1	<b>direct causal chain</b>	VC-VIP/7b-F4-F1	reaching	strict synchrony
2	<b>one level of indirection</b>	VC-AIP-F5-F1	poking, prodding, grasping	fast onset upon contact, potential for delayed effects
3	<b>complex causation involving multiple causal chains</b>	VC-AIP-F5-F1+STs+IT	mirror neurons, mimicry	arbitrarily delayed onset and effects
4	<b>complex causation involving multiple instances of manipulative acts</b>	STs+TE-TEO+F5-AIP(?)	object recognition	arbitrarily delayed onset and effects

Table 1

<i>object</i>	<i>angle between principal axis and preferred direction of rolling</i>	<i>behavior</i>
<b>cube</b>	n.a.	no principal axis, does not roll
<b>car</b>	0°	rolls along the principal axis
<b>bottle</b>	90°	rolls at right angle
<b>ball</b>	n.a.	no principal axis, does roll

Table 2

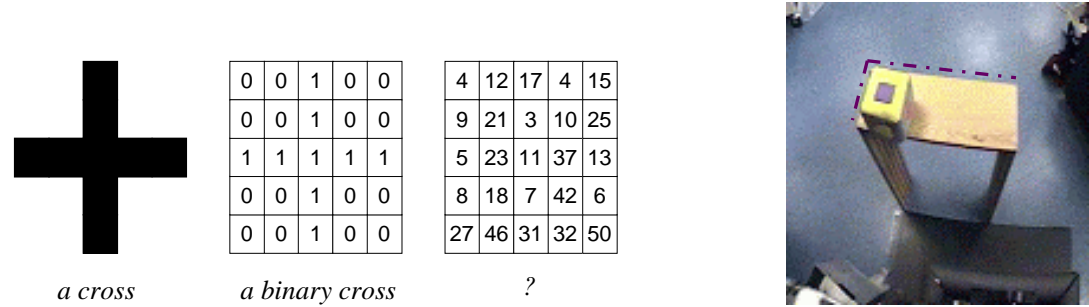


Figure 1

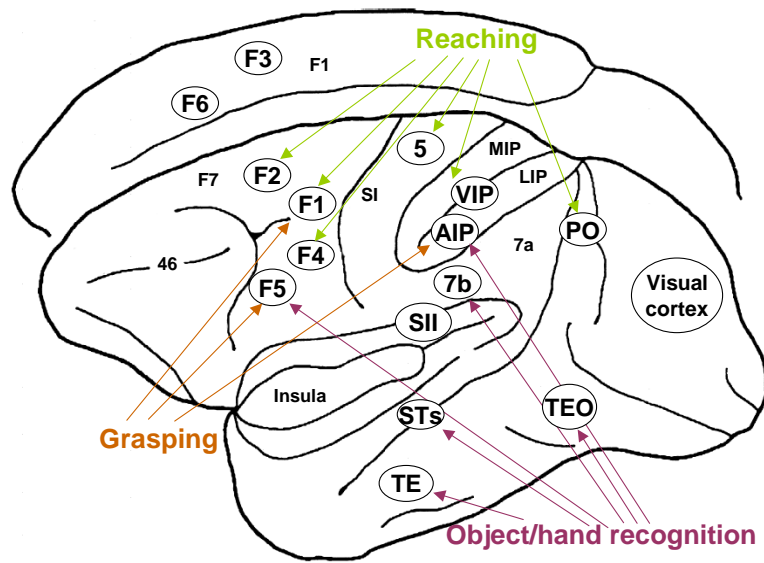


Figure 2



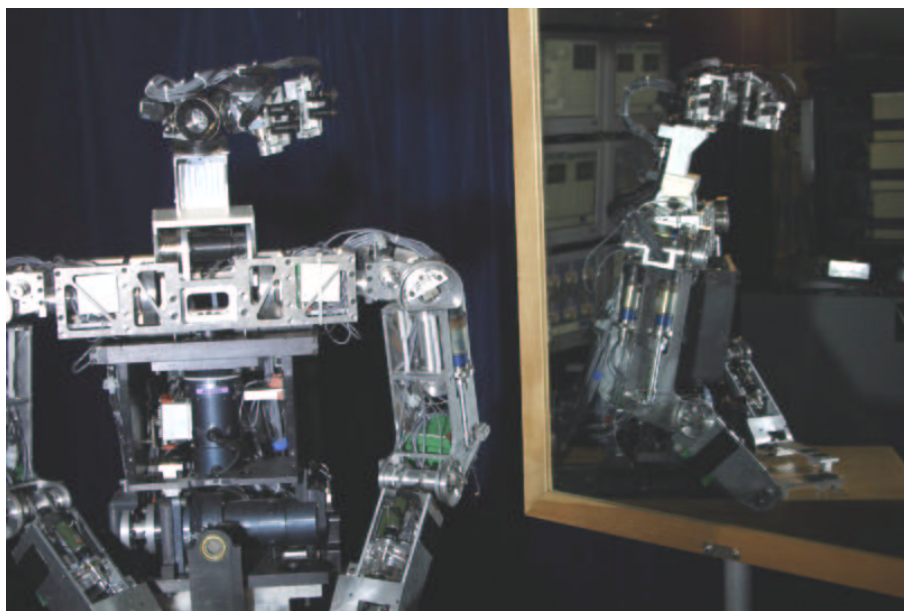


Figure 3

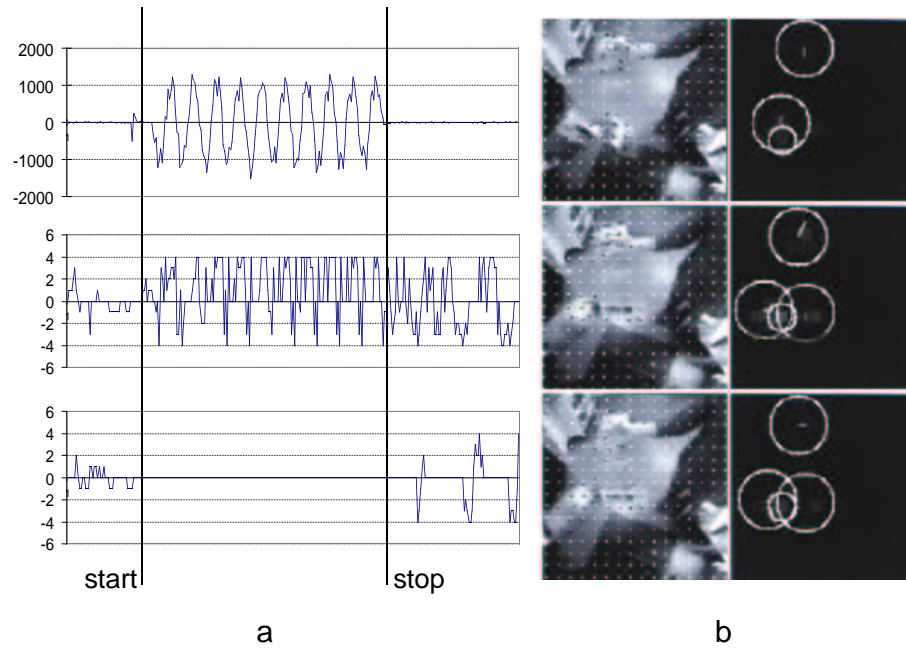


Figure 4



Figure 5

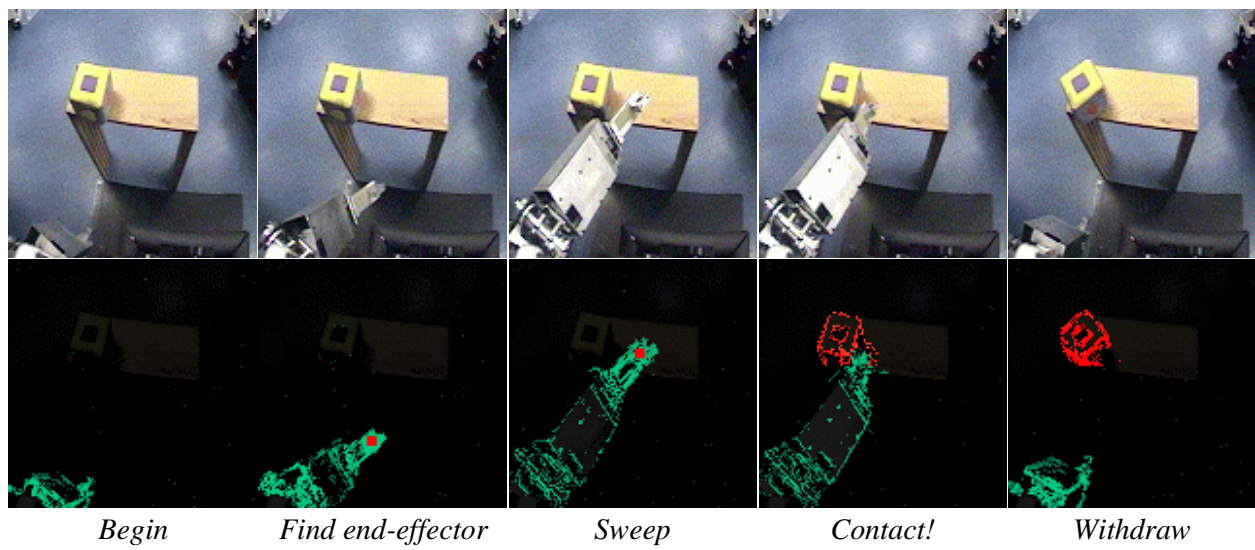


Figure 6

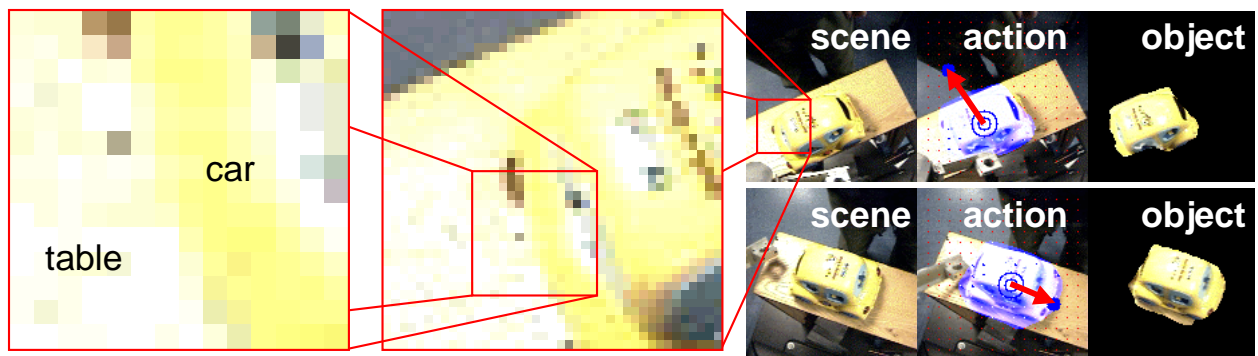


Figure 7

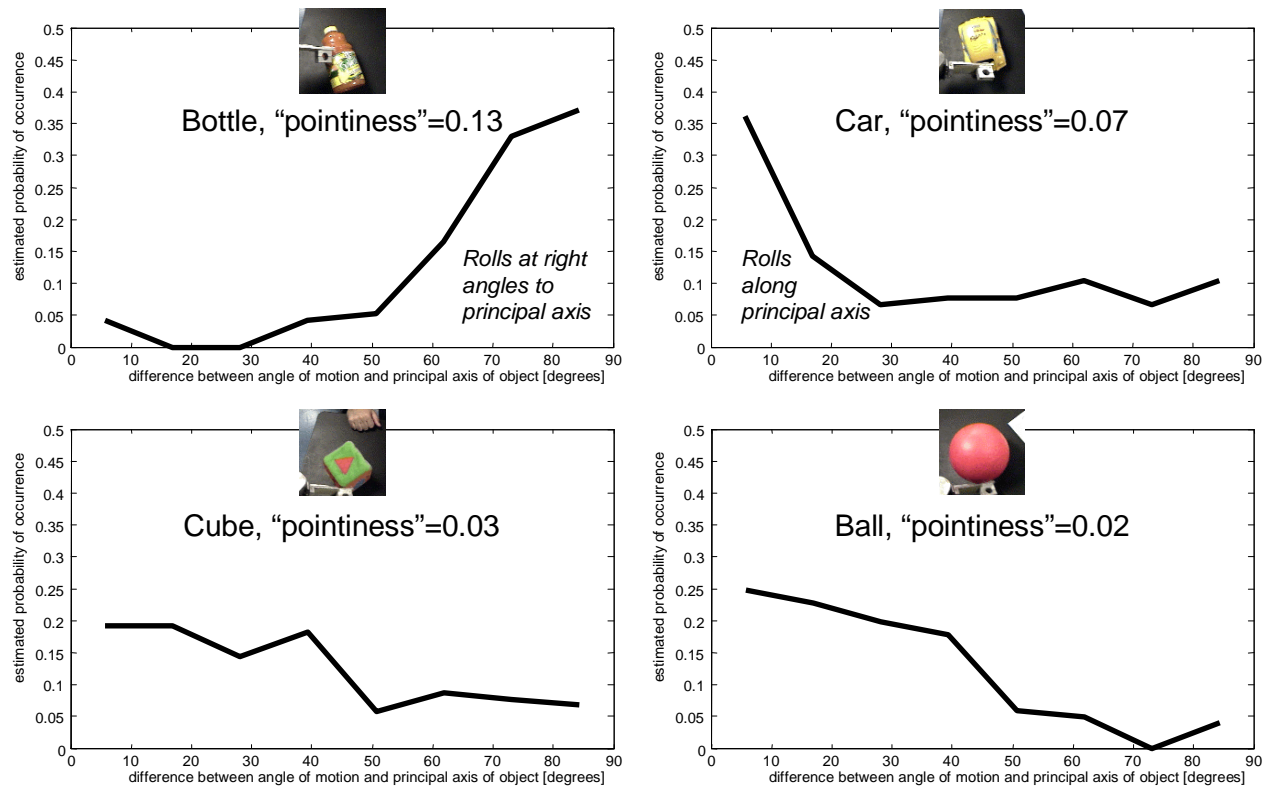


Figure 8

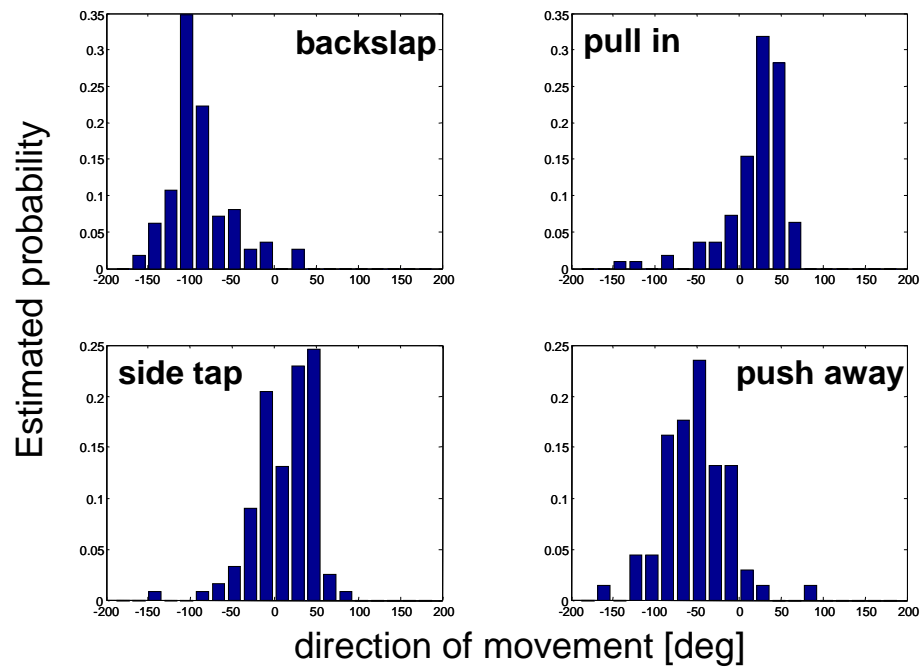


Figure 9

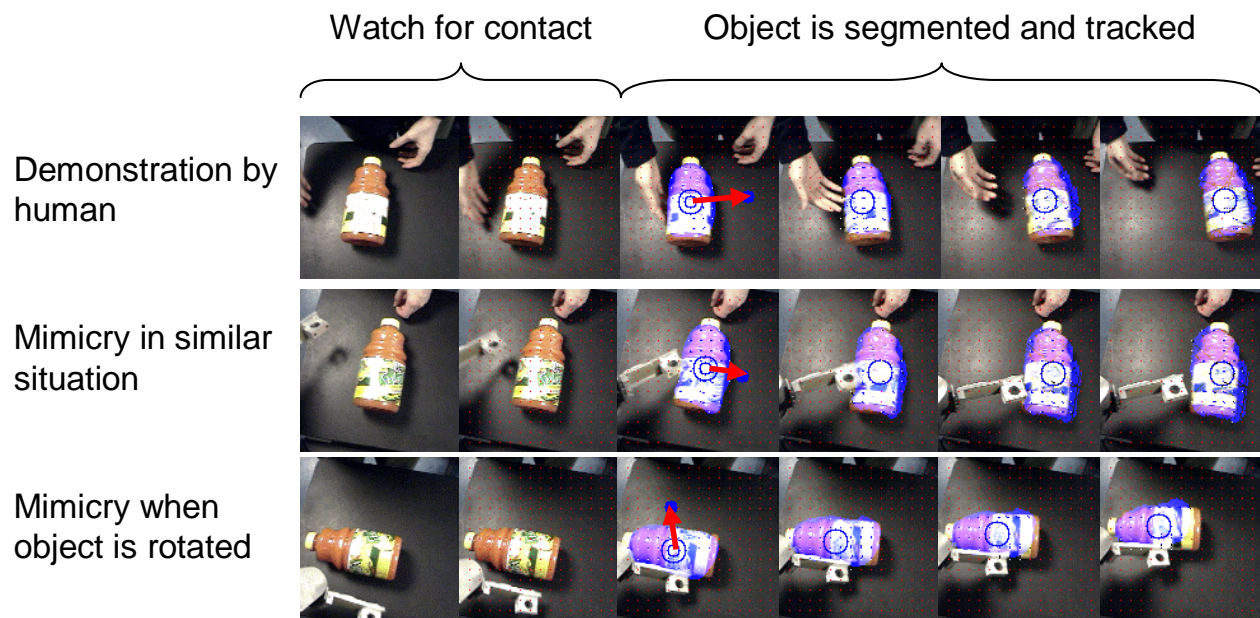


Figure 10



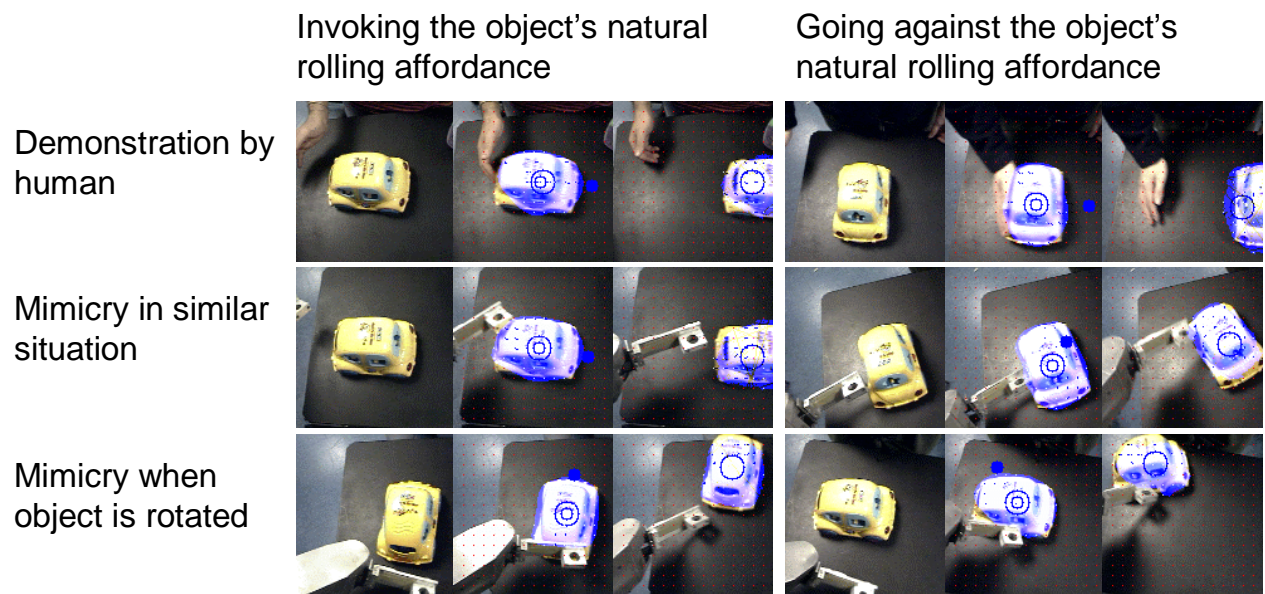


Figure 11

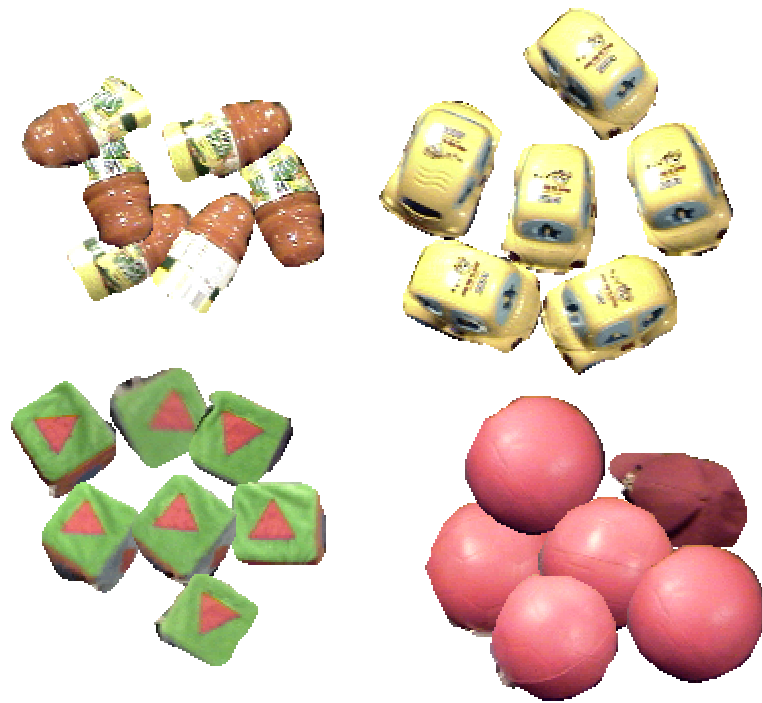


Figure 12

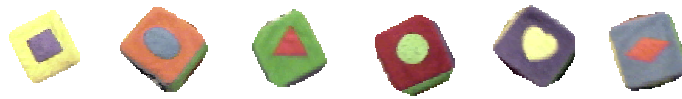


Figure 13

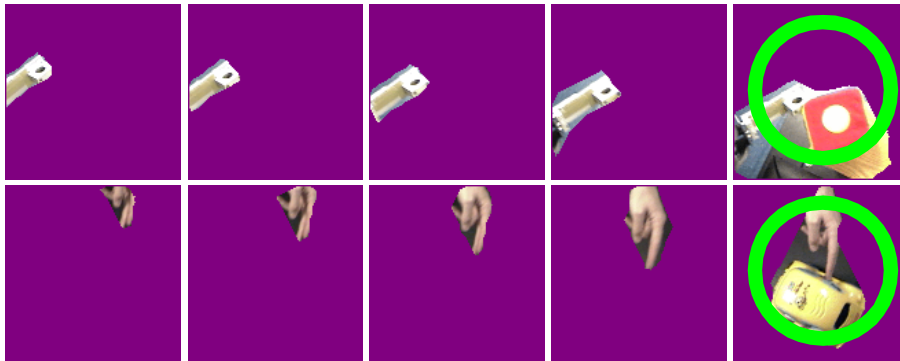


Figure 14