

# Early Integration of Vision and Manipulation

Giorgio Metta  
LIRA-Lab, DIST  
University of Genova  
Genova, Italy  
pasa@dist.unige.it

Paul Fitzpatrick  
Artificial Intelligence Lab  
Massachusetts Institute of Technology  
Cambridge, MA, USA  
paulfitz@ai.mit.edu

October 7, 2002

## Abstract

Vision and manipulation are inextricably intertwined in the primate brain. Tantalizing results from neuroscience are illuminating 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 progression of visual competence, and relate this idea to results in neuroscience.

## 1 Vision, action, and development

Robots and animals are actors in their environment, not simply passive observers. 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. 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.

We address three levels of causal complexity. The simplest causal chain that the actor experiences 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 her own actions and the actions of another. This is reminiscent of what has been observed in the response of the monkey's premotor cortex.

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 elaborate dextrous system. We found that simple poking gestures (prodding, tapping, swiping, batting, etc.) were rich enough to evoke object affordances such as rolling and to provide the kind of training data needed to bootstrap perception.

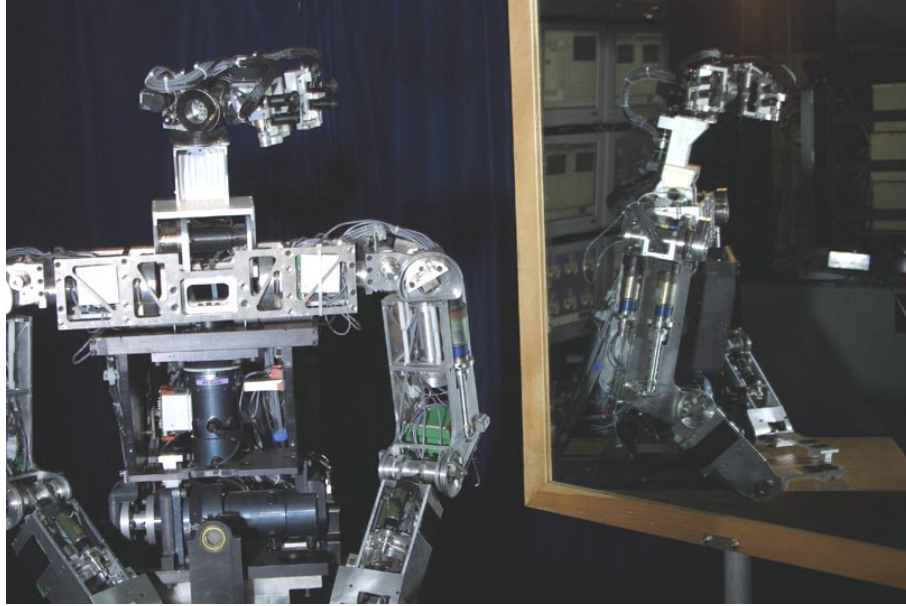


Figure 1: The world is complicated, but contingent. The ultimate goal of this work is for our robot to follow chains of causation outwards from its own simple body into the complex world. Such an incremental process suggests that perception and action develop together, supporting each other.

## 2 Objects and action in humans

Humans are experts at solving the figure/ground problem, and segmenting objects out from their background. 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.

Key to understanding how such capabilities could develop is the well-known result by Ungerleider and Mishkin Ungerleider and Mishkin (1982) who first formulated the hypothesis that objects are represented differently during action than they are for a purely perceptual task. Briefly, they argue that the brain's visual pathways split into two main streams: the dorsal and the ventral Milner and Goodale (1995). 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 dorsal/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.

Grossly simplifying (see also figure 2), the brain circuitry responsible for object oriented actions is thought to consist of at least four interacting regions: the primary motor cortex (F1), the premotor cortex (F4, F5), the inferior parietal lobule (AIP, LIP), and the temporal cortex (TE, TEO) (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

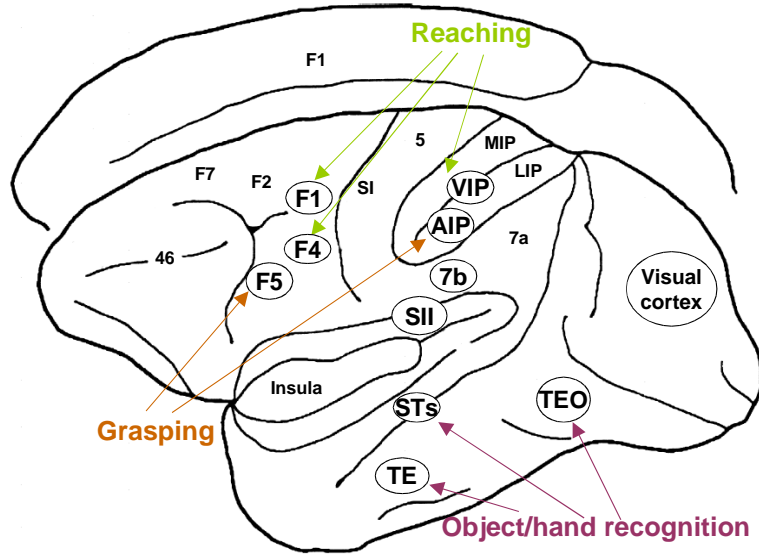


Figure 2: Monkey brain with indication of the main areas participating to object oriented actions (adapted from Fagg and Arbib (1998)). As described in the text, three main functions can be identified: object recognition, reaching, and grasping. These form three parallel yet connected streams of processing. The circuit connecting the visual cortex to the inferior parietal lobule (VIP/LIP), F4 and F1 subserves reaching. AIP and F5 are responsible for grasping. Temporal areas (TE, TEO) and STs are correlated to the semantic of object recognition. [It's likely we need to ask permission to use this figure].

example F4, LIP, 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 Jeannerod et al. (1995) that reaching and grasping need to interact to correctly orient and preshape the hand.

Neurons responsive to reaching are present in the inferior parietal lobule. For example, Jeannerod et al. (1995) reported that the temporary deactivation (suppression?) of the caudal part (VIP) of the intraparietal sulcus by injecting a GABA agonist disrupts reaching. 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. (1997); Fogassi et al. (1996). A subset of the F4 neurons has 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. (2000) 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 had been actually 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 just another big motor 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 (see Figure 3). 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 (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.

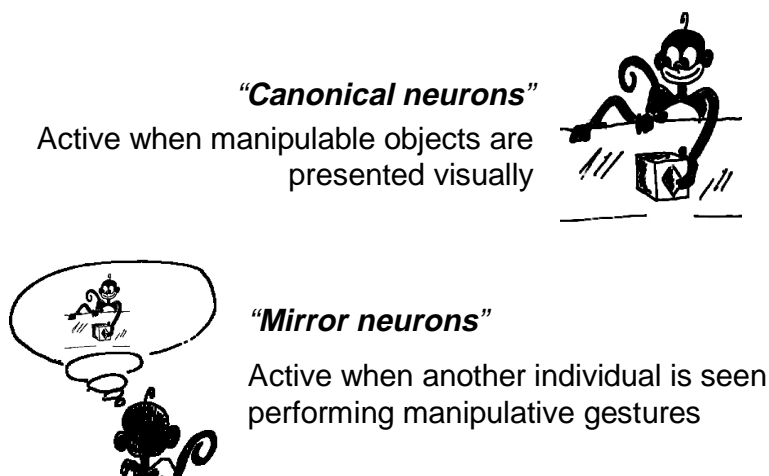


Figure 3: Canonical and mirror neurons.

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).

A possible developmental explanation of the acquisition of these functions can be framed in terms of tracing/interpreting chains of causally related events. Although it is still speculative, this analysis predicts that i) development of functions roughly follows a dorsal to ventral temporal gradient (i.e. reaching, grasping, recognition); ii) the ability to probe longer chains triggers the emergence of a new functionality and/or a new set of behaviors. The next section delves deeper into this proposal for the ontogenesis of object oriented action and provides a hypothesis amenable to implementation.

### 3 A working hypothesis

Taken together these results from neuroscience suggest a very basic role for motor action. 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.

We can distinguish three main conceptual functions (similar to the schema of Arbib et al. Arbib (1981)): reaching, grasping (manipulation), and object recognition. These functions correspond, as mentioned in the previous sections, to three levels of causal understanding (see ??). They form also a nice 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 the abovementioned developmental rules: i.e. specifying what information is retained during the interaction, what sort of sensory processing, what are the motor primitives, etc.

The neuroscience results outlined in the previous section can be streamlined into a developmental sequence roughly following a dorsal to ventral gradient. Unfortunately this is a question which was not 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. In humans, this function is accomplished mostly by the circuit VIP/LIP/7b-F4-F1. Reaching requires at least the detection of the object and the transformation of its position 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 subsystems 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. Area AIP and F5 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 the 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 a sensory consequence. 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. In this case the action execution has to be associated to 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 in the first place.

The last subsystem to develop is object recognition. Object recognition can build on manipulation in finding the boundaries of objects and segment 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. This functions are carried out by the infero-temporal cortex (IT). The same considerations apply to the recognition of the manipulator (either self or foreign). In fact, STs specialized to this task. Information about object identity is also sent to the parietal cortex and contributes to the formation of the affordances. For the actual recognition we can resort to a fuzzier definition of causality where multiple instances of manipulation on a certain object need to be grouped together. That is, all the information (visual in this case) pertaining to a certain object has to be grouped (and stored somewhere/somehow) to build a model of some sort of the object.

<i>nature of causation</i>	<i>main path</i>	<i>function and/or behavior</i>
<b>direct causal chain</b>	VC-VIP/LIP/7b-F4-F1	reaching
<b>one level of indirectness</b>	VC-AIP-F5-F1	poking, prodding, grasping
<b>complex causation involving multiple causal chains</b>	VC-AIP-F5-F1+STs+IT	mirror neurons, mimicry
<b>complex causation involving multiple instances of manipulative acts</b>	STs+TE-TEO	object recognition

Table 1: Degrees of causal indirectness, localization and function in the brain.

For the robotic implementation we followed the same developmental pathway and exploited the same sort of causal links between actions and sensory feedback.

## 4 The experimental platform

This work is implemented on the robot Cog, an upper torso humanoid Brooks et al. (1999). 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.

## 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

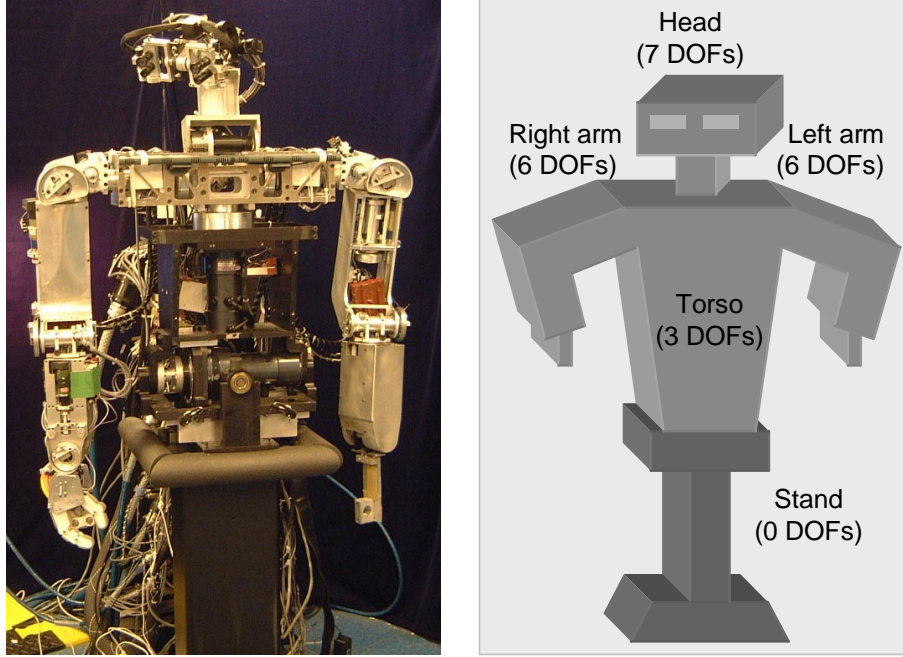


Figure 4: Degrees of freedom (DOFs) of the robot Cog. The arms terminate either in a primitive “flipper” or a four-fingered hand. The head, torso, and arms together contain 22 degrees of freedom.

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. 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.

### 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.

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



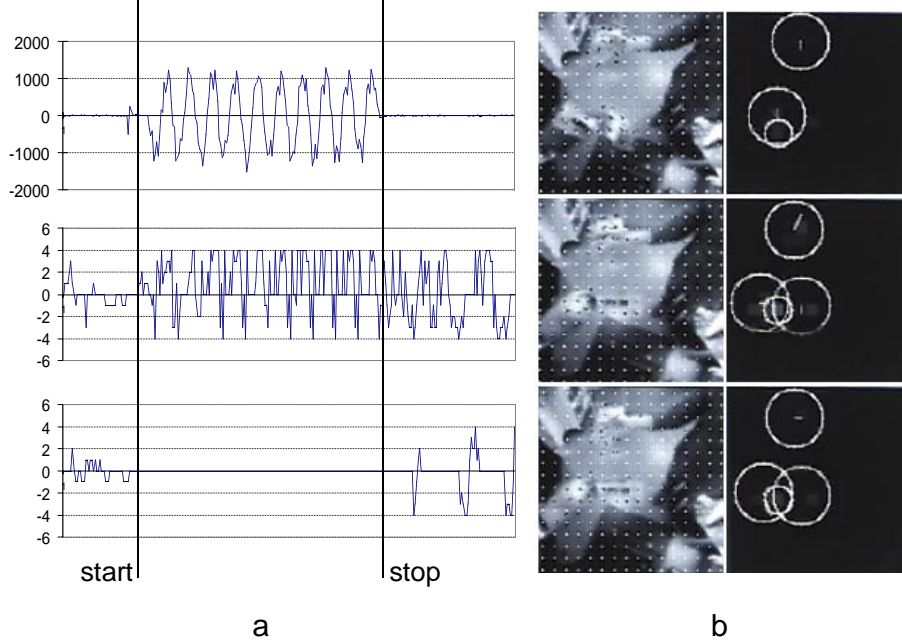


Figure 5: (a) An example of the correlation between optic flow and arm movement. The traces show the movement of the wrist joint (upper plot) and optic flow sampled on the arm (middle plot) and away from it (lower plot). (b) The robot’s point of view and the optic flow generated are shown on the left. 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.

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 (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 5 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 5 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.

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.





Figure 6: 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.

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 6 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. This procedure is consistent because we are interested in reaching a point on a plane in front of the robot (a table). The resulting map is thus  $2D \rightarrow 2D$ . 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 coefficient  $a_{nm}$  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.

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.

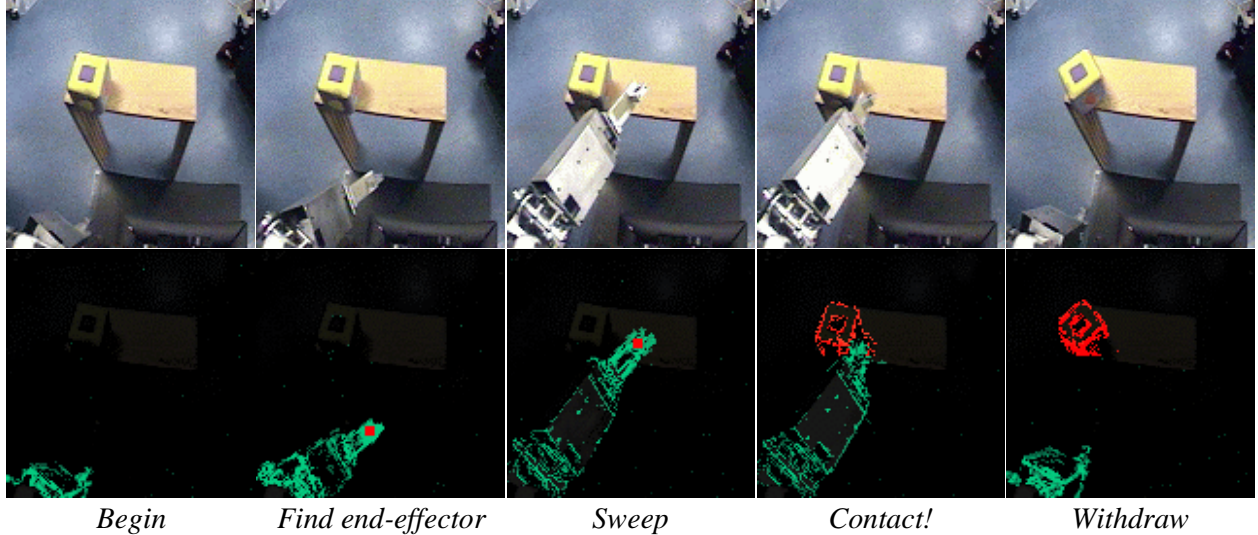


Figure 7: The upper sequence shows an arm extending into a workspace, tapping an object, and retracting. This is an exploratory mechanism for finding the boundaries of objects, and essentially requires the arm to collide with objects under normal operation, rather than as an occasional accident. The lower sequence shows the shape identified from the tap using simple image differencing and flipper tracking.

## 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 a correctly segmented. For the example scene in Figure ?? (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 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 motion at a higher resolution than optic flow. Because a poking operation currently

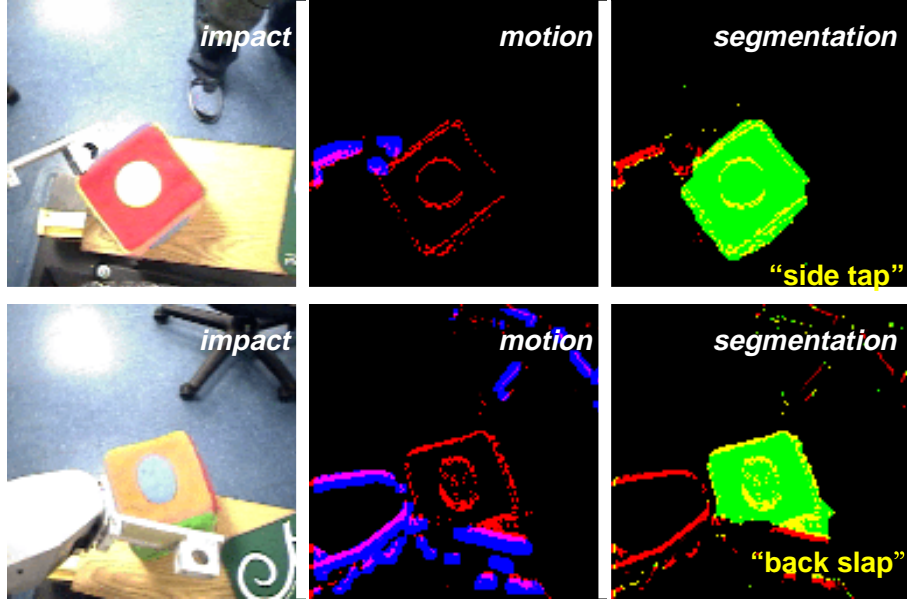


Figure 8: Cog batting a cube around. The top two rows show the flipper poking the object repeatedly from the side, turning it slightly. The third row shows Cog batting an object away. The images in the first column are frames prior to a collision. The second column shows the actual impact. The third column shows the motion signal at the point of contact. The bright regions in the images in the final column show the segmentations produced for the object.

always starts from the same location, the arm is localized using a simple heuristic rather than the procedure described in the previous section – the first region of optic flow appearing in the lower part of the robot’s view when the reach begins is assumed to be the arm.

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 can 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 (e.g. as COG’s flipper), is that “affordances” are scarce: most of the time the object

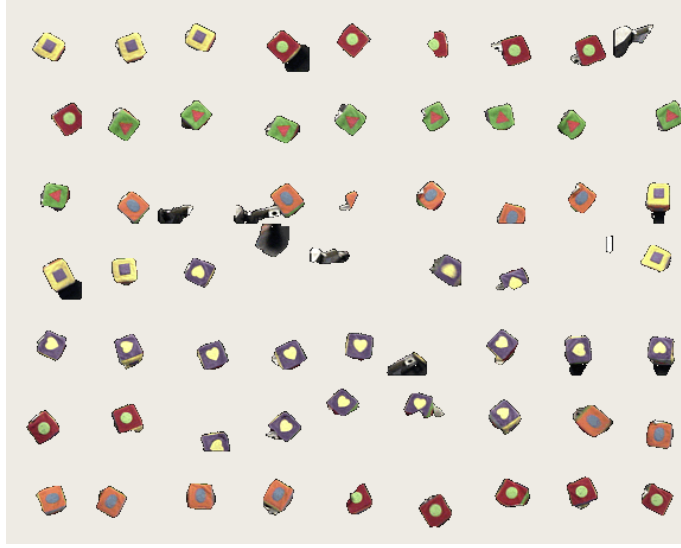


Figure 9: Sample results

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 constituent, 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.

<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: Behavior of a small set of objects when poked at random by the robot manipulator.

A further elaboration is required to group the data belonging to the same object as obtained from many poking acts into coherent clusters. As clustering mechanism we employed color histograms. After each poking action, a color histogram of the pixels of the segmented region is built and used as criterion to judge whether the object belongs to an existing group (e.g. if it is mostly yellow, it is likely to be the toy car). This works well for a small set of objects but sophisticated methods are required for a more general case with a large set of objects. 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 is to try to come to grips with the behavior of the objects after a single encounter, and use the behavior itself as clustering criterion. This is more difficult because of noise: e.g. there is still a non-zero probability that the object would not roll at all.

Figure 10 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.

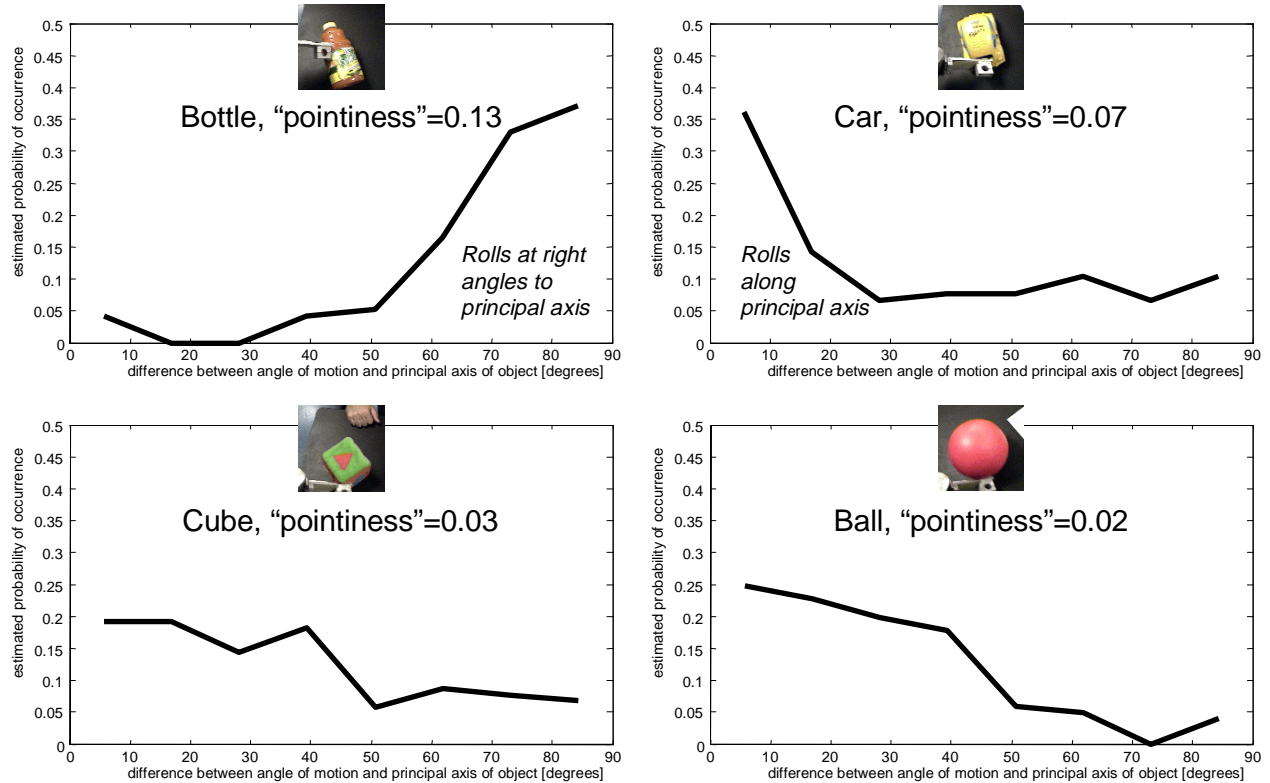


Figure 10: Probability of observing a roll along a particular direction for the set of four objects used in our experiments. Abscissae represent the difference between the principal axis of the object and the observed direction of movement. Ordinates the estimated probability.

Yet this description of the affordances does not have any usable quantity to take action once an object is observed. For this purpose a description of the geometry of poking is required: i.e. the description of the properties of objects (figure 10) has to be connected to a description of the behavior of the object. This information can be derived from the same training set we collected for learning about rolling. Figure 11 shows the histograms of the direction of movement of the object for each possible action. For example, the back slap moves the 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 11 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 (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 backprojection 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 10) 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.



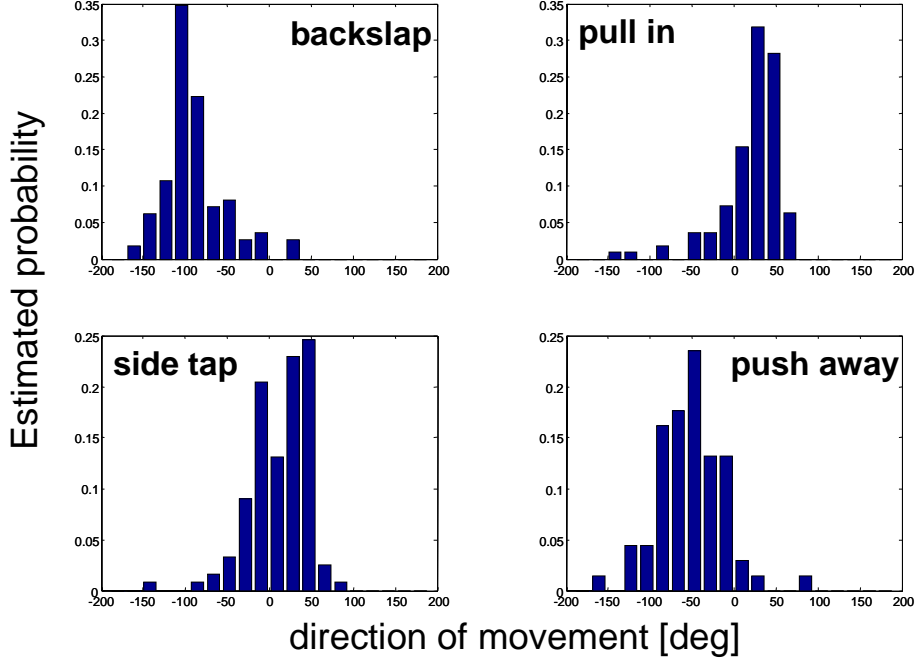


Figure 11: Histogram of the direction of movement of object for each possible poking action.

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 unprecise control (12) and to a less extent to misinterpretation of the orientation of the object (3).

## 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 [circumscribed] 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 learnt 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 12.

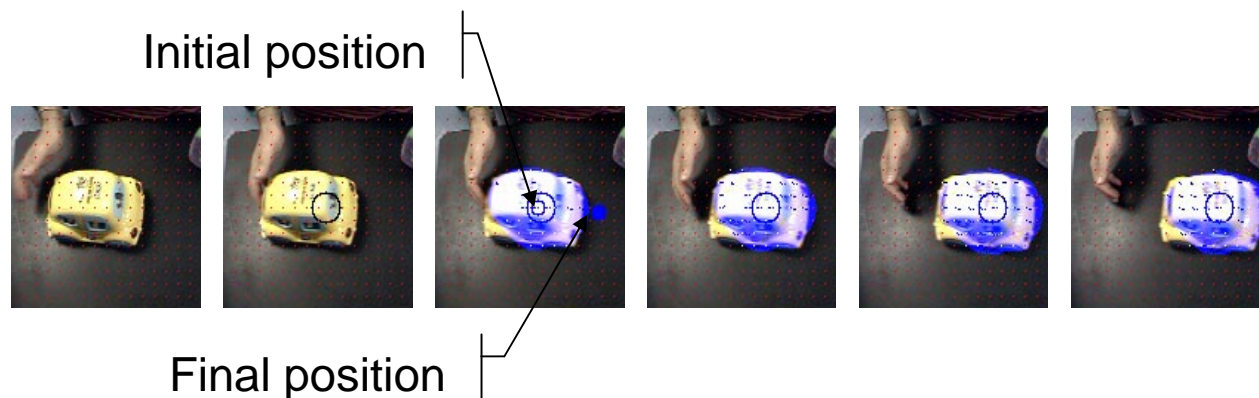


Figure 12: An example of observed sequence with tracking superimposed. Frames around the moment of contact are shown. The object, after segmentation, is tracked for 12 frames using a combination of template matching and optic flow. The big circles represent the position of the toy car in successive frames. The two small circles (outline and solid) displayed on the frame of contact ( $3^{rd}$  from the left) are the position at the time of contact and at the  $12^{th}$  frame respectively.

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 works well and 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 mimick 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. Figure 13 shows to examples of mimicry as a consequence of the action shown in figure 12.

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



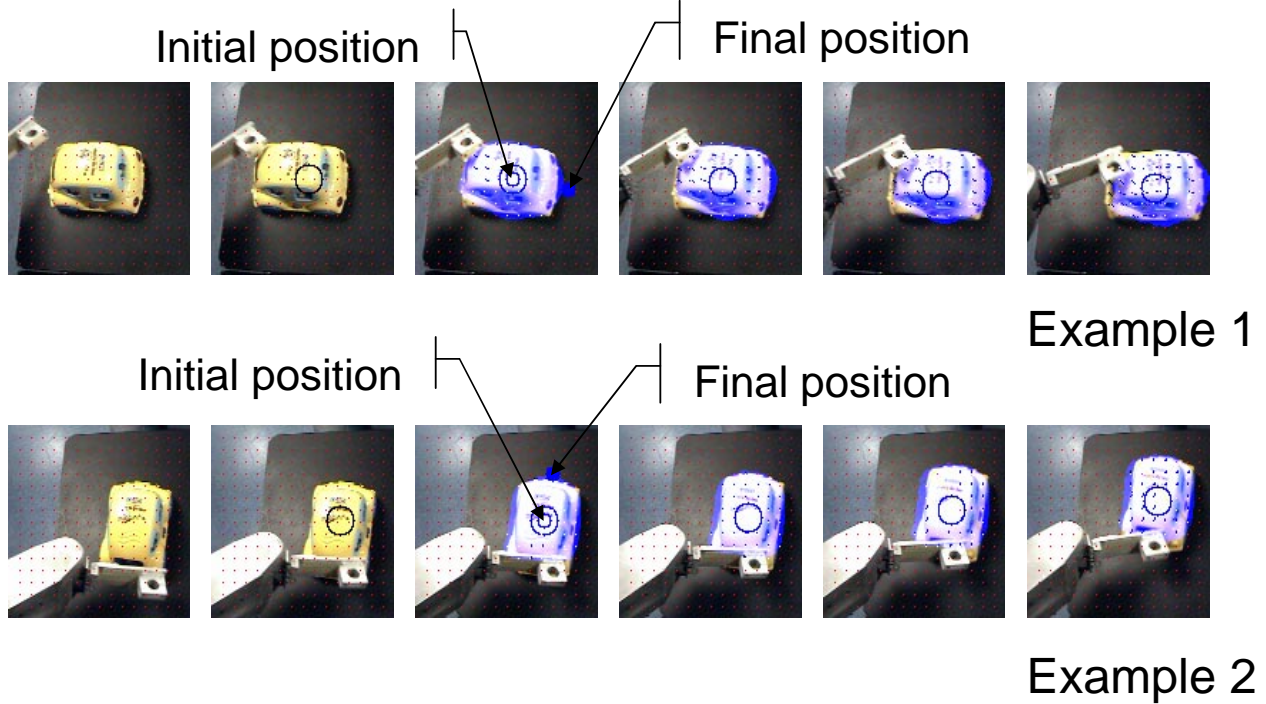


Figure 13: Two examples of mimicry following the observation in figure 12 where a human manipulator pokes the toy car exploiting the affordance (the car rolls). In example 1 (top row), the toy car has the same orientation it had in the demonstrated action and the robot repeats the observed action. In example 2 (bottom), the car is  $90^\circ$  with respect to example 1. The appropriate action to exploit the affordance and make the car roll is thus a back slap. )

explicitly coded (e.g. during reaching) Woodward (1998).

## 9 Towards object recognition

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 dexterous hand, it could further exploit temporal constraints (e.g. an object remains the same unless it is dropped) to collect tightly/temporally correlated data. 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 9). 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 the exactly same tool it used when initially learnt the task.

A possible and obvious extension is to use the object segmentation provided by poking (and manipulation 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 9). 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 15). The robot can then learn about the appearance of its own hand or, indifferently, 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.



Figure 14: Object clusters.

[We can envision a similar procedure to learn about any object that functions as manipulator.]

## Learning about manipulators

How could a robot find human arms and hands in the environment without any prior knowledge of their appearance? We could imagine segmenting any moving objects in the scene, and relying on the heuristic that hands are often the fastest moving objects around [cite]. Another approach is possible in our situation. If the robot can detect when an impact event occurs, it can collect segmentations of the object that caused the impact. The set of objects that habitually trigger the motion of other objects is not a bad operational definition of a manipulator, and should include the human hand/arm, and the robot's own arm. See Figure 15.

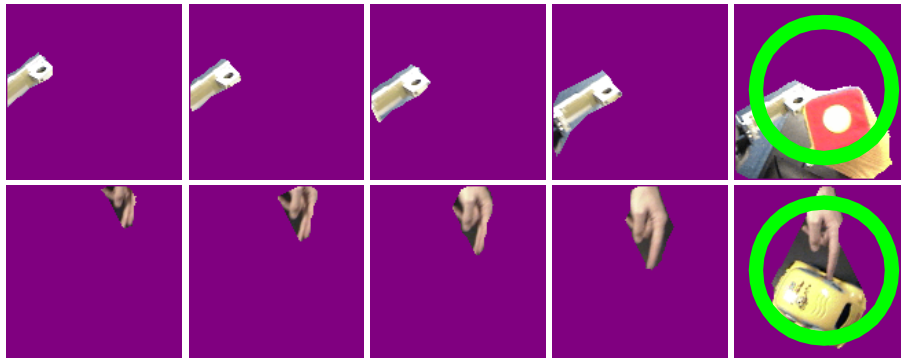


Figure 15: Early experiments on segmenting the robot arm, or a human hand poking an object the robot is familiar with, by working backwards from a collision event.

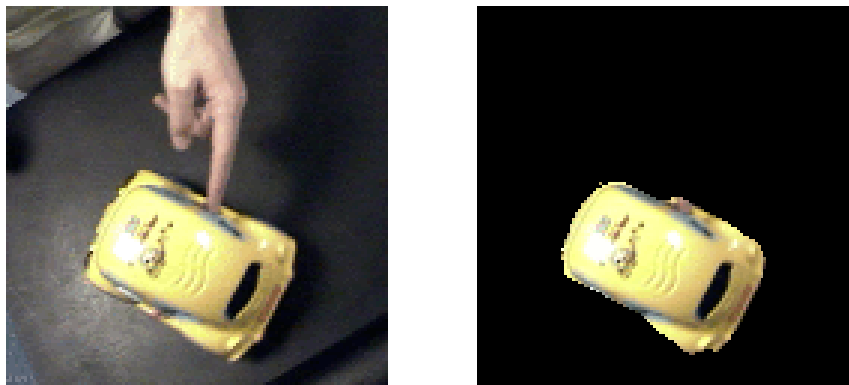


Figure 16: A poke by hand

## 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.

The robotic experiments support the view that reaching, grasping, and recognition can be learnt following a particular ontogenic pathway without the intervention of an external teacher [this is not strictly true for the present implementation because of the missing pieces]. This same sequences might be exploited by biological systems (primates/mammals) although the support up to date is rather tenuous. However, the sequence of events leading to object manipulation/recognition cannot take an arbitrary form unless we assume that some/many of its components are innate. 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 requires time) and complex perceptual abilities require a long time to emerge Kovacs (2000). We cannot claim that this is the only possible view but it is certainly one worth investigating. Rephrasing Berkeley we can say:

...objects can only be known by *action*. Vision is subject to illusions, which arise from *many different* problems...

[that AI guys know far too well]

Could relate some of this to the embodied intelligence ideas of Brooks... particularly the working hypothesis.

## Acknowledgements

This work benefited from discussions with Charles Kemp and Giulio Sandini. 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

- Arbib, M. A. (1981). *Handbook of Physiology*, chapter Perceptual Structures and Distributed Motor Control. American Physiological Society.
- Bekkering, H. and Wohlschläger, A. (2000). Imitation in children is goal-directed. *The quarterly journal of experimental psychology*, 53A(1):153–164.
- 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.
- 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., Cooke, D. F., and Taylor, C. S. R. (2000). Coding the location of the arm by sight. *Science*, 290(December):1782–1786.
- Graziano, M. S. A., Hu, X., and Gross, C. G. (1997). Visuo-spatial properties of ventral premotor cortex. *Journal of Neurophysiology*, 77:2268–2292.
- 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.
- 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.
- 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.
- Spelke, E. (2000). Core knowledge. *American Psychologist*, 55(11):145–160.
- 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.