

# Predictive Aspect Models for Multi-body Assembly

Grant Sherrick and Rod Grupen  
Laboratory for Perceptual Robotics  
Computer Science Department  
University of Massachusetts Amherst  
{sherrick, grupen}@cs.umass.edu

*Abstract—*

**The modeling of objects as a dynamic Bayesian network describing relationships between groups of related action possibilities allows for grasping, manipulation, and multi-body assemblies to all be reasoned about in the same representation across a large variety of tasks. We demonstrate the use of this representation in the context of a re-grasping task with a bimanual humanoid robot.**

## I. INTRODUCTION

This investigation extends past work [1] in functional object models (composed of affordances and their invariant spatial relationships) to multi-body assemblies and articulated objects. Functional object models that leverage the invariant spatial relationships between affordances, common to all rigid body objects, provide a powerful method for the inference of existence and locations of temporally co-occurring affordances. However, because affordances between objects in multi-body assemblies do not possess these invariant spatial properties, we propose an extension upon our current representation. We postulate that the result of each assembly can be described in a consistent manner with rigid body objects, in terms of their constituent affordances, as a single articulated object. In addition, we show that the region of assembly can be most compactly and intuitively described as an affordance in the force domain by a proper symmetry group in  $SE(3)$ . To ensure quality assemblies, we examine those assemblies that are statically stable while allowing for the expressiveness of our models to ensure robust behavior. We demonstrate the effectiveness of this approach on a re-grasping experiment involving several different assemblies in a variety of contexts with more extensive experiments and analysis to be provided in future investigations.

## II. RELATED WORK

Representing knowledge about the world in terms of controllable interactions provides a powerful and computationally efficient way for an agent to encode its experiences. Psychologist J. J. Gibson introduced the term affordance [2] as all *action possibilities* present in the environment that are objectively measurable in relation to an actor dependent on that actor's capabilities. The use of the theory of affordances within autonomous robotics is mostly confined to behavior-based control; consequently, its use in deliberation and planning remains a largely unexplored area. However, there have been several

recent studies that investigate the use of the theory of affordances in modeling and subsequent reasoning about a robot's environment. Specifically, Sun [3] reasoned about a mobile robot's environment to determine whether or not object's in its environment will support a variety of affordances including push, roll, traverse, support, etc. Their work utilized groups of perceptual features as a cue for the existence of related manipulation and locomotion affordances. Similarly, in this work, we utilize features from actions available to the robot in its current state to predict the success/failure and location of related affordances. However, in our representation, we treat all actions in a homogeneous manner and allow for any type of action (whether related to perception or manipulation) to be a cue for not only existence but also the location of related affordances. Stoytchev [4], [5] and Fitzpatrick [6] showed that affordance learning can be used to differentiate objects in the course of interaction with the environment. Stoytchev's and Fitzpatrick's work uses affordance as a higher level concept, which a developing cognitive agent learns about by interacting with objects in the environment. Montesano [7] presented an affordance based model based on Bayesian networks that linked actions and their effects to object features. In their work, functional information is stored as a feature set on the whole object so that it can then be applied in later recognition or interaction scenarios with the same or similar objects. In this paper, we concentrate on modeling probable groups of affordances and the spatial relationships between them. The modeling of multiple groups of affordances for each object provides support for robust behavior by allowing for multiple competing hypotheses about the location and existence of each affordance on an object while requiring no information about the pose of the object itself.

Uncertainty is a key issue when determining object and action parameters. Ek et al. [8] presented a system that is able to infer the location parameters for a commanded task and reason about action selection given information derived from partial observations. In their work, an optimal perceptual action is defined to be the action that will maximally disambiguate (reduce entropy over) the state-space. In Algorithm 1 of this paper, we present an action selection algorithm that first disambiguates the state of the robot defined over the equilibrium states of control actions. Following this step, our algorithm can then plan to a task through simulating the task by using our model in a generative manner to predict the sequence of actions and states that will reach the goal

with highest probability. In [9], Dragiev utilizes Gaussian processes to represent uncertainty for object pose and shape given sensor data from several different modalities in the context of reaching and grasping tasks. This representation results in distributions over pose and shape that are directly dependent upon observed data without making any simplifying assumptions other than requiring the pose distribution be unimodal across sensor input data. Similarly, in our current work, we use a Dynamic Bayesian Network to represent the distribution over the state of the robot’s environment that is described in terms of affordances, the relationships between them, and their properties in perceptual and motor spaces. In [10], Hsiao, Kaelbling, and Lozano-Perez developed a decision theoretic framework for task-driven exploration with POMDPs in which their system iteratively minimizes uncertainty in object pose by probing an object. In contrast to this system, we propose a system that suggests new actions based upon the expected decrease in uncertainty with respect to a task, represented by the successful completion of another action. This formulation is able to exploit past interactions from multiple objects, environments, and tasks, as well as, reason about the predicted effects of selected actions.

### III. BACKGROUND

#### A. Affordances in the Control Basis

Primitive control actions in this work,  $c \equiv \phi_\tau^\sigma$ , are described in the control basis [11]. These actions are closed-loop feedback controllers constructed by combining potential functions,  $\phi$ , with feedback signals,  $\sigma$ , and motor resources,  $\tau$ . The sensitivity of the output of the potential function to changes in the motor variables provides a control gradient that is used to derive reference motor inputs ( $\mathbf{u}_\tau$ ). Events in the error dynamics of each controller provide a natural discrete abstraction of the underlying continuous state space. In this work, we employ a four-valued control state,  $p(c) \in \{0, 1\}$ , where ‘0’ indicates the transient control response and ‘1’ denotes convergence/quiescence. Given a collection of  $n$  distinct primitive control actions, a discrete state space  $X$  can be formed, where  $\mathbf{x} \in X$  is defined by  $x = (p_1, \dots, p_n)$ .

#### B. Control Programs - SEARCHTRACK

In our representation of actions, we define two classes of control actions that share potential functions and effector resources: TRACK and SEARCH. These classes are distinguished by their differences in input signals; TRACK actions,  $\phi_\tau^\sigma$  preserve a reference value in the feedback signal e.g., the position of a feature on the image plane or the value of a contact force on a fingertip whereas SEARCH actions,  $\phi_\tau^{\tilde{\sigma}}$ , have an input signal,  $\tilde{\sigma}$ , that is derived from probabilistic models describing distributions over effector reference inputs ( $\mathbf{u}_\tau$ ), where TRACKing actions have converged in the past,  $p(\phi_\tau^\sigma) = 1$ . For example, such a controller can be used to direct the field of view of a robotic system to look at places on a table top where a specific signal has been found e.g. the color blue in an outdoors environment is typically found in the sky, above an agent. Initially the distribution  $Pr(\mathbf{u}_\tau | p(\phi_\tau^\sigma) = 1)$  is

uniform; however, over the course of many learning episodes, this distribution reflects the long term statistics of the runtime environment. The current investigation extends recent work [1] that improved upon existing signal-specific SEARCH distributions by developing a more robust object-specific representation for SEARCH distributions. In this work, we formalize these SEARCH distributions in the idea of predictive aspect models as well as extend this form of modeling to multi-object interactions.

TABLE I  
BASIC SOLIDS AND THEIR CORRESPONDING SYMMETRY GROUPS

Basic Shapes	Symmetry Groups
Half Plane	$G_{plane}$
Prism	$D_{2n}$
Cylinder	$G_{cyl}$
Sphere	$SO(3)$
Screw	$G_{screw}(p)$
Gear	$D_{2n}$
Cone	$SO(2)$
Pyramid	$C_n$

#### C. Applications of Group Theory to Assembly Planning

Group theory is the standard method for describing symmetry in any system. We will use group theory to describe the space of translations and rotations possible in an assembly [12]. In this theory a group  $\langle S, * \rangle$  consists of a set  $S$  with an operation  $*$  that has the properties of being associative, containing an identity element, and containing an inverse for each element of the set  $S$ . For instance, the space of real numbers  $\mathbb{R}$  with the operation multiplication excluding the number 0 is a group, 0 must be excluded here because it does not have an inverse under multiplication. Similarly, a subgroup  $\langle S1, * \rangle$  is a group that conforms to these same properties with the additional property that the set  $S1$  defining the group must be a subset of the set  $S$  and use the same operation  $*$  from another well defined group  $\langle S, * \rangle$ . The Proper Euclidean Group  $\mathcal{E}^+$  is defined as the set of all isometries (distance preserving mappings) of  $\mathbb{R}^3$  with the functional composition of isometries specified as the operation on the group. A proper symmetry,  $g \in \mathcal{E}^+$ , is an isometry (which preserves handedness) that brings  $S \subseteq \mathbb{R}^3$  into coincidence with itself i.e.  $g(S) = S$ . This set of isometries contains all possible translations and rotations in Euclidean space. In the representation we will present, each assembly is defined by a symmetry group describing the possible rotations and translations that will maintain that assembly.

### IV. EXPERIMENTS

#### A. Objects - Aspects Related by Actions

Objects in this representation are modeled as collections of control affordances (aspects) and the actions that define the relationships between them. Each aspect defines a coordinate frame which allows for the spatial relationships between affordances to be directly specified. However, because the

TABLE II  
SYMMETRY GROUPS ON EUCLIDEAN SPACE ( $\mathcal{E}$ )

Groups	Members
$G_{identity}$	identity transformations
$T_1$	translations along a line
$T_2$	translations in the plane
$T_3$	translations in Euclidean Space
$SO(3)$	rotations in Euclidean Space
$SO(2)$	rotations about a line in Euclidean Space
$SO(1)$	rotations in a plane in Euclidean Space
$G_{cyl}$	translation along the axis of the cylinder and rotation about the axis
$G_{plane}$	translation in the plane and rotation about the normal
$G_{screw}$	translation along the axis of the screw and rotation about the axis of the screw by $\frac{2z\pi}{p}$ , where $z$ is the amount of rotation about the axis
$D_{2n}$	rotations by $\frac{2\pi}{n}$ and flipping about the normal to these rotations
$C_n$	rotations by $\frac{2\pi}{n}$
$\mathcal{E}^+$	rotations and translations in Euclidean Space, excludes reflections

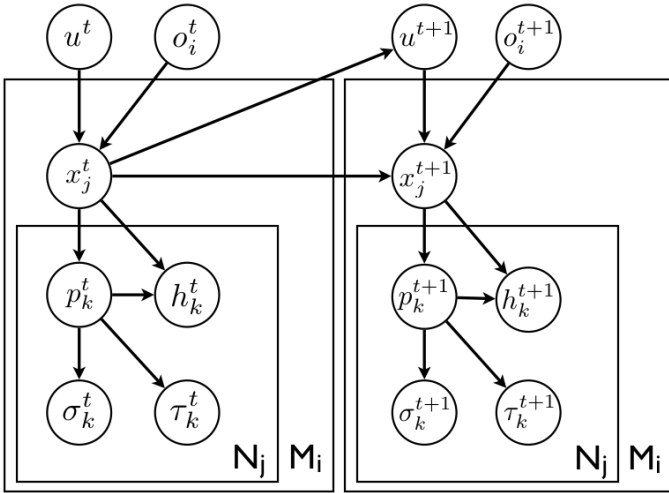


Fig. 1. A dynamic Bayesian network representing object  $o_i$  at time  $t$  as a spatial distribution over  $M_i$  aspects, each represented by Bernoulli random variable  $x_j$ . The action taken (whether for perception or manipulation) is represented by random variable  $u$ . An aspect induces a distribution over the state of controllable interactions ( $p_k$ ) afforded by the object. The random variables  $\sigma_k$  and  $\tau_k$  model the distributions over signal (e.g., color, force) and effector (e.g. hand, pan-tilt unit) as gaussian distributions. The gaussian random variable  $h_k$  models the position and orientation of the affordance instance in the object frame.

relationships between aspects are specified in terms of manipulation action, there is no one object frame. We take inspiration for this type of qualitative modeling of spatial relationships from work on grounded cognition in the psychology literature in which it is hypothesized that humans use simulation as a primary mechanism for reasoning about actions and objects [13]. Figure 1 shows a graphical model that encodes the logical dependencies between the variables of the environment affordance model. An object,  $o_i$ , at any instant of time affords a set of controlled interactions with the robot. Each stable set of spatially distributed interactions define an object aspect ( $a_j$ ). For each aspect, there exist  $N_j$  affordances that have a non-zero probability of occurring. There can be multiple

instances of each aspect within an object. Each affordance is represented by a Bernoulli random variable  $c_k$  describing the state of each associated SEARCHTRACK action. ( $c_k = 1$ , if the action converges, and 0 otherwise.) The difference between affordances and taking actions is that affordances are passive control action possibilities. Affordances in this representation never represent the taking of an action, but they can represent whether or not an action in the current state of the robot interacting with the environment has converged (succeeded) or failed to converge. Each possible affordance is modeled by its position and orientation ( $p_k$ ) in the object's frame, the feature values of the signal ( $f_k$ ), and the shape ( $s_k$ ) defined by the eigenvalues of the signal. The resulting generative model describes objects in terms of affordances and the spatial relationships between them.

Utilizing past experience encoded as a prior, this model is able to aid in accomplishing a variety of tasks by telling the robot which affordances are likely to co-occur and where they occur with objects that are similar to those that have been previously encountered. However, by introducing a temporal dependence, we are able to describe how taking actions affect the existence and location of affordances. For example, if the robot would like to pick an object up off of a table, then the state of the model must afford lifting i.e. the robot must first grasp the object in order to be able control the object in the direction of lifting. In this case, “grasping” changes the aspect of the object in a manner that supports the goal of lifting. We encode the aspect-action dependencies of an object as a dynamic Bayesian network shown in Figure 1, where the instance of an aspect being observed is a hidden variable that the robot can infer from state of the affordances,  $P^t = \{p_1^t, \dots, p_k^t\}$  and the actions,  $u^t$ . This temporal model consists of a finite number of states (given by the aspect,  $x^t$ ), a finite number of actions  $U = \{u_1, u_2, \dots, u_k\}$  and a set of possible observations. For every aspect instance,  $x^t$ , the transition probability  $\Pr(x_j^t | u^{t-1}, x_j^{t-1})$  describing the probability that an aspect,  $x_j^{t-1}$ , transitions to another aspect instance,  $x_j^t$ , by taking manipulation action,  $u^{t-1}$ .

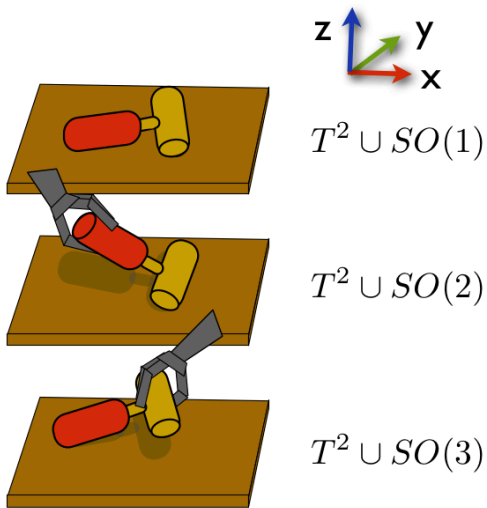


Fig. 2. Above are drawn the three possible symmetry groups for a table and mallet. (The robotic hand is included to help disambiguate the three cases due to lack of artistic talent of the authors.) In the first case, at the top of the figure, the mallet is able to be rotated around the  $z$  axis, or moved in the  $x$ - $y$  plane in order to maintain the current assembly. In the middle of the figure, the mallet is contacting the table along a line that runs the length of the cylindrical mallet head. In this case, the mallet is able to rotate freely around this line as well as around the  $z$  axis. In the lowest part of the figure, the cylindrical handle of the mallet is contacting the table at a point. This frees up the mallet to rotate around this point as well as move in the plane of the table.

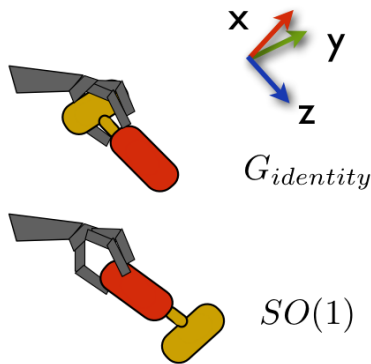


Fig. 3. Above are drawn two possible symmetry groups for a three finger hand and mallet. In the first case, at the top of the figure, the mallet is not able to rotate with respect to the hand while maintaining the current assembly (which expects force in all 6 dimensions of  $SE(3)^*$ ). In the lower part of the figure, the cylindrical handle of the mallet is contacting the hand at two points. This frees up the mallet to rotate around the  $x$  axis of the coordinate frame.

### B. Assemblies - Collections of Objects (Meta-objects)

Assemblies are modeled in a manner consistent with the previously mentioned modeling of rigid body objects, see Figure 1 and Section IV-A. Each object in an assembly is described by its own aspect, a stable collection of affordances that co-occur in the same region of space and time. Because each object in the assembly is denoted by a separate aspect and each aspect is defined by its own coordinate frame, the region of assembly is defined twice (once with respect to each aspect in each object of the assembly). In addition, for an assembly to be properly defined, there must exist an

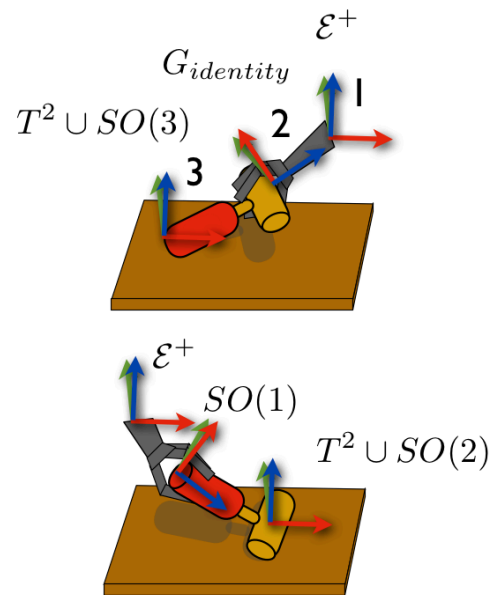


Fig. 4. Above are drawn two examples of assemblies composed of three two-body subassemblies and four objects: subassembly 1, between arm and hand with no expected force signature i.e. complete freedom of movement, subassembly 2 between hand and mallet, subassembly 3 between mallet and table. In each assembly, the symmetry group associated with each subassembly is given. These groups denote the only permissible actions which will maintain each subassembly. In addition, the complement of each group, in wrench space, denotes the set of forces that should be expected at the region of assembly.

affordance in each aspect whose reference values reside in the force domain, which is located at the junction of the two objects. This affordance describes the expected forces which will occur between the two objects with a signal value that consist of a symmetry group in  $SE(3)$ , see Table II for a list of characteristic symmetry groups. In general, these groups are derived from the proper Euclidean group, which describes the space of possible translations and rotations in Euclidean space. However, as has been done in past work [14], we will use these signal values both to describe possible motions of one object with respect to another in an assembly, and also to describe the dual wrench space of expected force signatures that are present as a result of the constraints imposed by the assembly. We are currently working to be able to derive this force-based affordance for describing assemblies from already existing affordances on single rigid body objects by taking the intersection of each separate wrench space, see Mason and Salisbury [15] and Popplestone, Liu, and Weiss [16] for examples of using group theory in this manner.

Grasping and manipulation can also be modeled as an assembly in this representation. As was discussed previously, each of the objects in an assembly (the robot's end effector and a rigid body object in this case) are modeled as separate aspects of one meta-object. Because each aspect is defined functionally in terms of affordances, this representation at first seems ill suited for the modeling of the robot's end effector. However, a robot's end effector can be described by

affordances associated with other sensor and effector modalities (not the force sensors or proprioception that are most likely a part of the end effector). In addition, when describing assemblies between the robot’s end effector and rigid body objects (grasping and manipulation), the robot’s end effector can be described in terms of the actions that it affords with respect to the other object that it is interacting with. This makes most sense when thinking of the object that the robot is manipulating or grasping as being the object taking action and the the robot’s end effector being a static rigid body. For example, when a robot is grasping the head of a mallet with a palmar grasp, as can be seen in the lower part of Figure 3, the robot’s hand is fixed and affords the force-domain-based action that immobilizes the cylindrical head of the mallet associated with the identity group,  $G_{identity}$ .

Control motions can be described through thinking of the end effector of the robot as being in an assembly with some object that affords no forces i.e. complete freedom to move and orient in Euclidean space. Therefore, any object that is an assembly with the end effector will always be described as if there is one extra virtual object in the assembly. This allows for the space of possible motion that will maintain each assembly to be described, see Figure 4.

---

**Algorithm 1** ACTIONSELECTION( $u_g, Z, U$ )

---

```

1:  $u_g$  - goal action
2:  $Z$  - sequence of observations
3:  $U$  - sequence of actions taken
4: repeat
5:   Compute  $\Pr(a_j^t | z^t, u^{t-1}, a_j^{t-1})$  of observation set  $z^t \in Z$  for predictive aspect models  $a_j^t \in A, j = 1, \dots, M_i$ 
6:   if  $\Pr(a_j^t | z^t, u^{t-1}, a_j^{t-1}) > \Pr(a_k^t | z^t, u^{t-1}, a_k^{t-1})$  then
7:     place  $a_j^t$  in set  $C$ , candidates
8:   end if
9:   Infer the pose distribution,  $\Pr(x_g^t | a_j^t, z^t)$ , of goal affordance,  $u_g$ , for each candidate aspect  $a_j^t \in C$ .
10:  if  $\sum_{j,k:j \neq k} D_{KL}(\Pr(x_g^t | a_j^t, z^t) || \Pr(x_g^t | a_k^t, z^t)) < \alpha$  then
11:    if goal pose,  $x_g^t$ , is a valid reference for action  $u_g$  then
12:      execute action  $u_g$ 
13:    else
14:      CHANGEASPECTTOREACHGOAL( $u_g, Z, U$ )
15:    end if
16:  else
17:    DECREASEASPECTUNCERTAINTY( $u_g, Z, U$ )
18:  end if
19: until  $p(a_g) = 1$  {goal action,  $u_g$ , succeeds}

```

---

### C. Task-based Action Selection

Reference values for each control action can be sampled from the Bayesian model. These actions can then be executed given knowledge of the object’s pose in the world frame. However, in the presence of partial information, choosing an action

---

**Algorithm 2** CHANGEASPECTTOREACHGOAL( $u_g, Z, U$ )

---

```

1:  $u_g$  - goal action
2:  $Z$  - sequence of observations
3:  $U$  - sequence of actions taken
4: Choose action,  $u^t$  from the model which changes the object aspect to one which affords the goal:
5: (This requires backchaining from an aspect that affords the goal to the current aspect)
6:  $A_n = \sum_{a \in A} \Pr(a_j^t | a_j^{t-1} = a, u^t = u_g)$ 
7:  $A_p = \sum_{a \in A_g} \Pr(a_j^{t-1} | u^t = u_g, a_j^t = a)$ 
8: while  $\Pr(A_p = C) < \beta$  do
9:    $A_n = \sum_{a \in A, u \in U} \Pr(a_j^t | a_j^{t-1} = a, u^t = u)$ 
10:   $A_p = \sum_{a \in A_p, u \in U} \Pr(a_j^{t-1} | u^t = u, a_j^t = a)$ 
11: end while
12:  $t = t + 1$ 
13:  $u^t = \arg \max_{u^t} \left[ \sum_{a_p \in A_p, a_n \in A_n} \Pr(a_j^{t-1} = a_p, u, a_j^t = a_n) \right]$ 
14: Gather new evidence,  $z^t$ 
15: Add new evidence and action taken to  $Z, U$ 
16: ACTIONSELECTION( $u^{t-1}, Z, U$ )

```

---



---

**Algorithm 3** DECREASEASPECTUNCERTAINTY( $u_g, Z, U$ )

---

```

1:  $u_g$  - goal action
2:  $Z$  - sequence of observations
3:  $U$  - sequence of actions taken
4: Choose action,  $u^t$  from the model which maximally reduces the uncertainty over aspects:
5:  $t = t + 1$ 
6:  $u^t = \arg \max_{u^t} \sum_{a_n, a_p \in A_i} H(a_n^{t+1} | u^t, z^t, a_p^t)$ 
7: {for object  $o_i$  with aspect set  $A_i$ }
8: Gather new evidence,  $z^t$ 
9: Add new evidence and action taken to  $Z, U$ 
10: ACTIONSELECTION( $u^{t-1}, Z, U$ )

```

---

given that it may be expensive or destructive (w.r.t. sensor measurements) requires safeguards to ensure that the robot chooses the next action that will lead towards successfully completing its intended task. It is not necessary for a robot to completely determine the pose of an object before it can take actions towards achieving its goal, it is only required that the pose relative to the goal action is completely determined. The procedure for taking such an action ( $u_g$ ) is described in Algorithm 1.

Given a task and an object model, the robot first takes action to reduce its uncertainty with respect to aspect, then attempts to reach a state in which the goal action can be achieved. The algorithm begins each iteration by finding the probability of each aspect given the latest observations of control actions,  $z^t$ , the last manipulation action,  $u^{t-1}$ , and the last aspect,  $a_j^{t-1}$  (Line 5). This is achieved by computing the probability that each aspect given the evidence can be generated by the model. If the same evidence is afforded in multiple regions of the object, a set of candidate aspects,  $C$ , are returned (of which only one is the aspect). Each of the candidate

aspects is combined with relative spatial information from the object model to produce candidate positions and orientations of the goal (Line 9). For example, if the goal is to grasp an object, this relative information is provided by SEARCH distributions for the arm(s) and hand(s) relative to the object frame that orient the system appropriately for TRACK-able forces comprising a grasp. If all the candidate aspects represent the same set of positions and orientations of the goal in the world frame i.e. the total KL divergence between all such distributions is low (Line 10), we say that the candidate goals are equivalent and unambiguous. Hence, even though there is uncertainty in the pose of the object, there is none in the goal affordance. Such cases arise in the case of symmetric objects, where the pose of the object may remain ambiguous even when grasping goals are not. If goals are ambiguous, then actions are selected that reduce the uncertainty over the distribution of aspects by calling Algorithm 3. In a sense, this algorithm tries to figure out “where” it is, before it tries to “drive” there. The action taken is that which maximally reduces the uncertainty (entropy) over the set of aspects,  $A_i$  for object  $o_i$ , given the latest observations,  $z^t$ , and aspects,  $a_p^t \in A_i$ . This is similar to the idea of an agent first increasing its belief that it is in a particular state before trying to achieve a task from the recent work by Platt, Tedrake, Kaelbling, and Lozano-Perez [17]. If the model has low uncertainty with respect to its aspect then Algorithm 2 is called. This algorithm back-chains from an aspect that affords the goal action until it has reached its current state. The action taken is then the action, which has the highest probability of reaching the aspect in which the goal action can succeed.

#### D. Quasi-static Multi-body Regrasping Experiments

We apply the representation described in Sections IV-A and IV-B in the context of a regrasp experiment with our humanoid robot Dexter. Dexter is a bimanual robot with two 7-DOF Whole-Arm Manipulators (WAMs) from Barrett Technologies, two 3-finger 4-DOF Barrett Hands equipped with one 6-axis force/torque load cell sensor on each fingertip, a stereo camera pair and a Kinect mounted on a pan/tilt head.

In this set of experiments, we show the efficacy of our representation for utility-driven action selection in the context of multi-body grasping and manipulation. Models of objects were hand-built spatial distributions of blobs (represented in terms of first and second moments) describing color space, range image clusters, and search distributions of force space goals represented in Cartesian space where “pincer grasp”, “palmar grasp”, and planar affordances can be found. We require that any objects in the environment maintain quasi-static stability i.e.  $\sum_{i=1,2,3} F_i = 0$  and  $\sum_{i=1,2,3} M_i = 0$ . Because different parts of the mallet support different assemblies with the hand, there are two sets of possible goals for the hand-mallet assembly. The mallet’s handle and head support a palmar grasp, whereas the entire body supports a pincer grasp. In certain regions of the workspace, the object does not afford haptic aspects, and additional manipulation actions have to be taken before grasp goals can be achieved. Figure 5 shows the

case when the object is presented in a region where the robot can grasp successfully. In such a case, Algorithm 1 computes the aspect of the object by matching observations to the model and returns the control action sequence necessary to achieve the grasp action. However, when the goal affordance is out of reach (and hence the object aspect doesn’t afford the goal - grasping in this case), the action selection algorithm chooses a manipulation action that can change the aspect to one that affords grasping. Figure 6 shows an example where the robot chooses to pull the object towards itself before executing the grasp action. Because of our quasi-static requirement, it can be seen that the robot is required to use the table in order to maintain force closure of the mallet due to the pincer grasp being the only achievable grasp.

#### V. CONCLUSIONS AND FUTURE WORK

We described in this study an example of using this representation for a regrasp task in which there were several different types of assembly both between hand and object (grasping/manipulation) and between two rigid body objects. We plan to utilize this framework in the performance of a larger variety of tasks including several more examples of assembly from those given in Table II.

#### ACKNOWLEDGMENT

This work was supported by the DARPA grant iRobot Corp. Prime Army W91CRB-10-C-0127. Grant Sherrick was supported by a Massachusetts Space Grant Consortium Summer Fellowship.

#### REFERENCES

- [1] S. Sen, G. Sherrick, D. Ruiken, and R. Grupen, “Choosing informative actions for manipulation tasks,” in *In Proceedings of 11th IEEE-RAS International Conference on Humanoid Robots*, 2011.
- [2] J. Gibson, “The theory of affordances,” in *Perceiving, acting and knowing: toward an ecological psychology*. Hillsdale, NJ: Lawrence Erlbaum Associates Publishers, 1977, pp. 67–82.
- [3] J. Sun, J. L. Moore, A. Bobick, and J. M. Rehg, “Learning Visual Object Categories for Robot Affordance Prediction,” *The International Journal of Robotics Research*, vol. 29, no. 2-3, pp. 174–197, 2010. [Online]. Available: <http://ijr.sagepub.com/content/29/2-3/174.abstract>
- [4] A. Stoytchev, “Toward learning the binding affordances of objects: A behavior-grounded approach,” in *Proceedings of the AAAI Spring Symposium on Developmental Robotics*, Stanford University, 2005.
- [5] —, “Behavior-grounded representation of tool affordances,” in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, Barcelona, Spain, 2005.
- [6] P. Fitzpatrick, G. Metta, L. Natale, S. Rao, and G. Sandini, “Learning about objects through action: Initial steps towards artificial cognition,” in *IEEE International Conference on Robotics and Automation*, Taipei, May 2003.
- [7] L. Montesano, M. Lopes, A. Bernardino, and J. Santos-Victor, “Modeling affordances using bayesian networks,” in *Proceedings of the IEEE International Conference on Intelligent Robots and Systems*, San Diego, CA, 2007.
- [8] C. Ek, D. Song, K. Huebner, and D. Kragic, “Task modeling in imitation learning using latent variable models,” in *In Proceedings of 10th IEEE-RAS International Conference on Humanoid Robots*, 2010.
- [9] S. Dragiev, M. Toussaint, and M. Gienger, “Gaussian process implicit surfaces for shape estimation and fluent grasping,” in *Proceedings of the IEEE International Conference on Robotics and Automation*. IEEE, 2011.
- [10] K. Hsiao, L. Kaelbling, and T. Lozano-Perez, “Task-driven tactile exploration,” in *Proceedings of Robotics: Science and Systems*, Zaragoza, Spain, June 2010.

- [11] J. Coelho and R. Gruben, "A control basis for learning multifingered grasps," *Journal of Robotic Systems*, vol. 14, no. 7, pp. 545–557, 1997.
- [12] Y. Liu and R. Popplestone, "A group theoretic formalization of surface contact," *The International Journal of Robotics Research*, vol. 13, no. 2, pp. 148–161, 1994.
- [13] L. W. Barsalou, "Grounded cognition: Past, present, and future," *Topics in Cognitive Science*, vol. 2, no. 4, pp. 716–724, 2010. [Online]. Available: <http://dx.doi.org/10.1111/j.1756-8765.2010.01115.x>
- [14] Y. Liu and R. Popplestone, "From characteristic invariants to stiffness matrices," in *Robotics and Automation, 1992. Proceedings., 1992 IEEE International Conference on*, may 1992, pp. 2375 –2380 vol.3.
- [15] M. Mason and J. K. Salisbury, *Robot Hands and the Mechanics of Manipulation*. Cambridge, MA: MIT Press, 1985.
- [16] R. Popplestone, Y. Liu, and R. Weiss, "A group theoretical approach to assembly planning," *Artificial Intelligence Magazine*, pp. 82 – 97, 1990.
- [17] R. Platt, R. Tedrake, L. Kaelbling, and T. Lozano-Perez, "Belief space planning assuming maximum likelihood observations," in *Proceedings of Robotics: Science and Systems*, Zaragoza, Spain, June 2010.

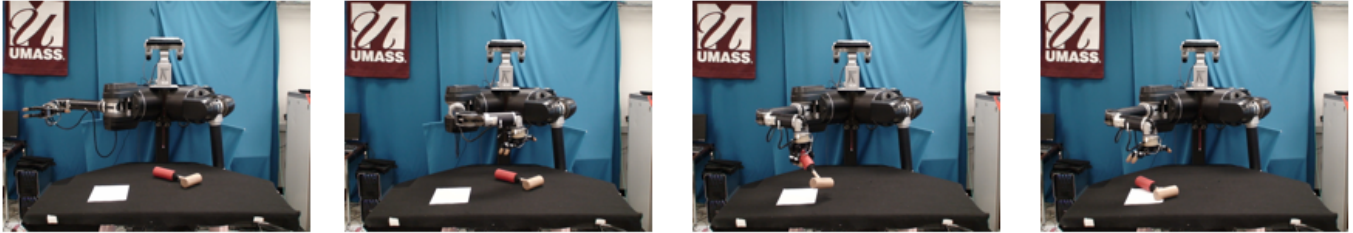


Fig. 5. The robot performing a top grasp on the mallet and placing it on the goal.

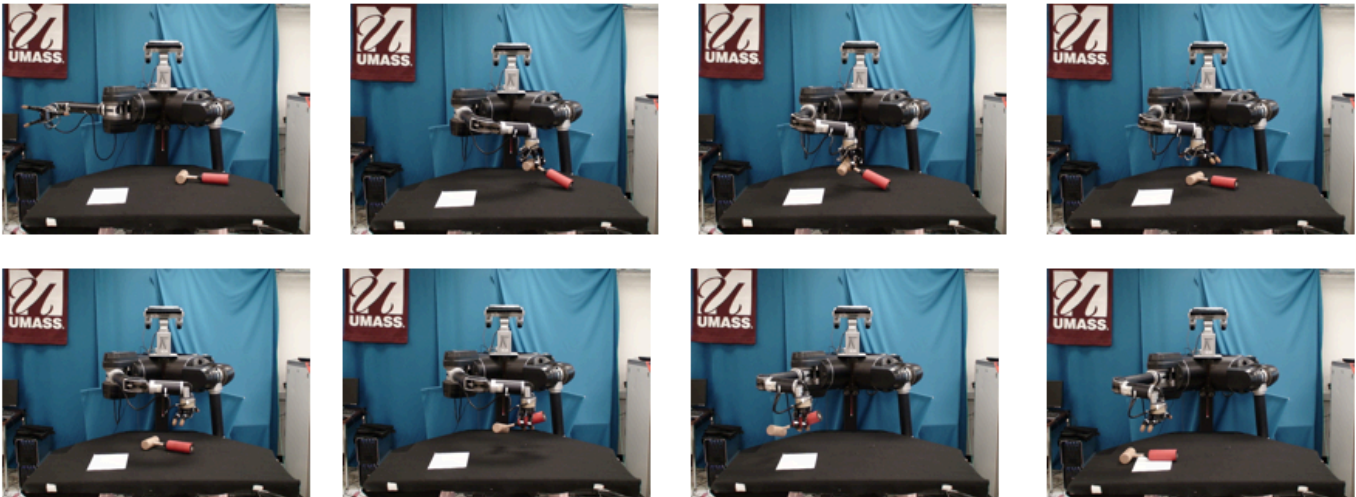


Fig. 6. The robot pulling the mallet towards itself before performing a top grasp on the object and placing it on the goal.