

# Configuration Lattices for Planar Contact Manipulation Under Uncertainty

Michael C. Koval\*, David Hsu<sup>†</sup>, Nancy S. Pollard\*, and Siddhartha S. Srinivasa\*

[mkoval@cs.cmu.edu](mailto:mkoval@cs.cmu.edu), [dyhsu@comp.nus.edu.sg](mailto:dyhsu@comp.nus.edu.sg), [{nsp, siddh}@cs.cmu.edu">{nsp, siddh}@cs.cmu.edu](mailto)

\*The Robotics Institute, Carnegie Mellon University

<sup>†</sup>Department of Computer Science, National University of Singapore

**Abstract**—This work addresses the challenge of a robot using real-time feedback from contact sensors to reliably manipulate a movable object on a cluttered tabletop. We formulate contact manipulation as a partially observable Markov decision process (POMDP) in the joint space of robot configurations and object poses. The POMDP formulation enables the robot to actively gather information and reduce the uncertainty on the object pose. Further, it incorporates all major constraints for robot manipulation: kinematic reachability, self-collision, and collision with obstacles. To solve the POMDP, we apply DESPOT, a state-of-the-art online POMDP algorithm. Our approach leverages two key ideas for computational efficiency. First, it performs lazy construction of a configuration-space lattice by interleaving construction of the lattice and online POMDP planning. Second, it combines online and offline POMDP planning by solving relaxed POMDP offline and using the solution to guide the online search algorithm. We evaluated the proposed approach on a seven degree-of-freedom robot arm in simulation environments. It significantly outperforms several existing algorithms, including some commonly used heuristics for contact manipulation under uncertainty.

## I. INTRODUCTION

Our goal is to enable robots to use real-time feedback from contact sensors to reliably manipulate their environment. Contact sensing is an ideal source of feedback for a manipulator because contact is intimately linked with manipulation: the sense of touch is unaffected by occlusion and directly observes the forces that the robot imparts on its environment.

However, contact sensors suffer from a key limitation: they provide rich information while in contact with an object, but little information otherwise. Fully utilizing contact sensors requires *active sensing*. Consider the robot shown in fig. 1 that is trying to push a bottle into its palm. The robot is uncertain about the initial pose of the bottle, has only an approximate model of physics, and receives feedback from noisy contact sensors on its fingertips. The robot executes an *information-gathering action* by moving its hand laterally to force the bottle into one of its contact sensors. Once the bottle is localized, it moves to achieve the goal.

In this paper, our goal is to autonomously generate policies that allow a robot to use real-time feedback from contact sensors to reliably manipulate objects under uncertainty. We specifically consider the problem of quasistatic tabletop manipulation via pushing [33]. The robot’s goal is to push a movable object into a hand-relative goal region while avoiding collision with stationary obstacles in the environment.

We formulate the *planar contact manipulation problem* as a *partially observable Markov decision process* (POMDP) [23,

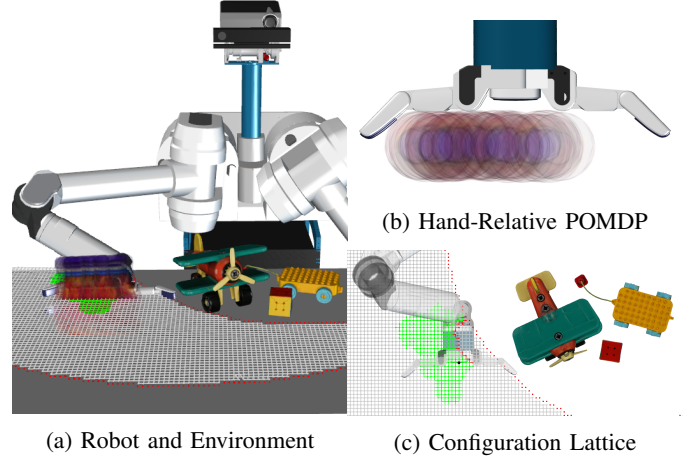


Fig. 1: A robot (a) uses real-time feedback from contact sensors to manipulate an bottle on a cluttered table. (b) Our approach uses policies from the obstacle free, hand-relative contact manipulation problem to a search in the full state space. (c) The configuration of the robot is represented as a point in a state-space lattice consisting of feasible (—), infeasible (—), and unevaluated (—) edges. The probability distribution over the pose of the bottle is shown as a collection of semi-transparent renders. Best viewed in color.

47] with a reward function that drives the robot towards the goal. Prior work has successfully used an offline POMDP solver [28] to find robust, closed-loop policies for manipulating an object relative to the hand [17, 27].

Unfortunately, these hand-relative policies perform poorly when executed on a robotic manipulator because of limited reachability, collision with the environment, and other types of *kinematic constraints*. These constraints occur frequently during execution, but cannot be represented in a *hand-relative POMDP* (Rel-POMDP). In this paper, we address this limitation by planning in a POMDP that includes both the fully-observable configuration of the manipulator and the partially-observable pose of the object.

Finding a near-optimal policy for this POMDP is challenging for two reasons. First, a manipulator’s configuration space is continuous and high-dimensional, typically having at least six dimensions. Second, it is difficult to perform pre-computation because the optimal policy may dramatically change when an obstacle is added to, removed from, or moved within the environment.

We leverage two key insights to overcome those challenges:

- 1) We define a *configuration lattice POMDP* (Lat-POMDP) that models both the configuration of the robot and the pose of the object. Lat-POMDP lazily constructs [8, 41] a lattice in the robot’s configuration space and uses the lattice to guarantee that the optimal policy does not take infeasible actions (section IV).
- 2) We prove that the optimal value function of Rel-POMDP is an upper bound on the optimal value function of Lat-POMDP. We use this fact to create a heuristic that guides DESPOT [48], a state-of-the-art online POMDP solver, to quickly find a near-optimal policy for Lat-POMDP with no environment-specific offline pre-computation (section V).

We validate the efficacy of the proposed algorithm on a simulation of HERB [49], a mobile manipulator equipped with a 7-DOF Barrett WAM arm [45] and the BarrettHand [51] end-effector, manipulating an object on a cluttered tabletop (fig. 1, section VI). First, we confirm that DESPOT successfully uses information-gathering actions to achieve good performance in the absence of kinematic constraints. Then, we show that the proposed algorithm outperforms five baselines in cluttered environments.

The proposed algorithm demonstrates that it is possible for an online POMDP planner to simultaneously reason about object pose uncertainty and kinematic constraints in contact manipulation. We are excited about the prospect of exploiting the scalability of recent developments in online POMDP solvers, like DESPOT, to solve more complex tasks in future work (section VII).

## II. RELATED WORK

Our work builds on a long history of research on enabling robotics reliably manipulating objects under uncertainty. Early work focused planning open-loop trajectories that can successfully reconfigure an object despite non-deterministic uncertainty [29] in its initial pose [5, 14]. More recently, the same type of worst-case, non-deterministic analysis has been used to plan robust open-loop trajectories for grasping [9, 10] and rearrangement planning [11, 25] under the quasistatic assumption [33, 34]. Our approach also makes the quasistatic assumption, but differs in two important ways: it (1) considers probabilistic uncertainty [29] and (2) produces a closed-loop policy that uses real-time feedback from contact sensors.

Another line of research aims to incorporate feedback from contact sensors into manipulator control policies. This approach have been successfully used to optimize the quality of a grasp [42] or servo towards a desired contact sensor observation [30, 52]. It is also possible to directly learn a feedback policy that is robust to uncertainty [38, 50]. These approaches have achieved impressive real-time performance, but require a higher-level planner to specify the goal.

One common approach is to plan a sequence of move-until-touch actions that localize an object, then execute an open-loop trajectory to complete the task [16, 21, 39]. Other approaches, like our own, formulate the problem as a POMDP [23, 47]

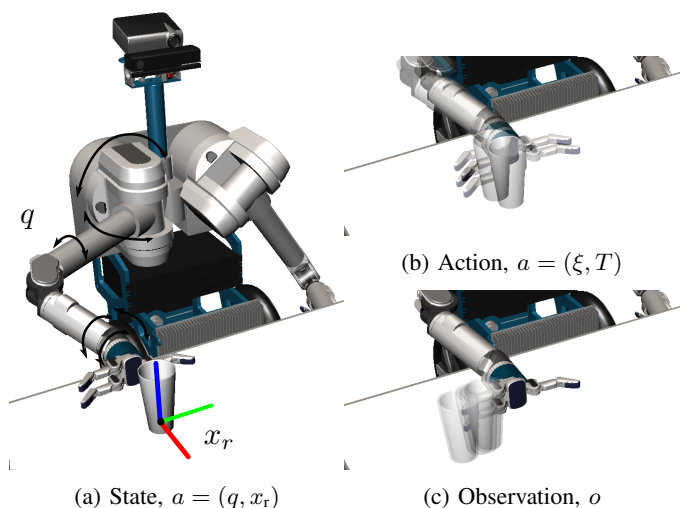


Fig. 2: We formulate planar contact manipulation as a POMDP with a state space (a) that contains the robot configuration  $q$  and the pose of the movable object  $x_r$ . The environment contains one movable object (white glass) and static obstacles. (b) action is a short joint-space trajectory  $\xi$  of duration  $T$ . After executing an action, the robot receives feedback from (c) binary contact sensors on its end effector.

and find a policy that only takes information-gathering actions when they are necessary to achieve the goal [18–20].

Unfortunately, most of this work assumes that the end-effector can move freely in the workspace and that objects do not significantly move when touched. More recent approaches have relaxed the latter assumption by incorporating a stochastic physics model into the POMDP [17, 27] and using SARSOP [28], an offline point-based [40] POMDP solver, to find a near-optimal policy for manipulating an object relative to the hand. Unfortunately, hand-relative policies often fail when executed on a manipulator due to kinematic constraints or collision with obstacles. We use a hand-relative policy to guide DESPOT [48], an online POMDP solver [44], in Lat-POMDP, a mixed-observable model [36] that includes these constraints.

Lat-POMDP represents the robot’s configuration space as a state lattice [41], a concept that we borrow from mobile robot navigation [31] and search-based planning for manipulators [8]. Similar to these approaches, and many randomized motion planners [3, 15], we use lazy evaluation to defer collision checking until an action is queried by the planner.

## III. PROBLEM FORMULATION

We formulate planar contact manipulation as a partially observable Markov decision process (POMDP) [47]. A POMDP is a tuple  $(S, A, O, T, \Omega, R)$  where  $S$  is the set of states,  $A$  is the set of actions,  $O$  is the set of observations,  $T(s, a, s') = p(s'|s, a)$  is the transition model,  $\Omega(s, a, o) = p(o|s, a)$  is the observation model, and  $R(s, a) : S \times A \rightarrow \mathbb{R}$  is the reward function.

The robot does not know the true state  $s_t$ . Instead, the robot tracks the *belief state*  $b(s_t) = p(s_t|a_{1:t}, o_{1:t})$ , which is a probability distribution over the current state  $s_t$  conditioned on the history of previous actions  $a_{1:t} = \{a_1, \dots, a_t\}$  and observations  $o_{1:t} = \{o_1, \dots, o_t\}$ . The set of all belief states is known as *belief space*  $\Delta$ .

Our goal is to find a policy  $\pi : \Delta \rightarrow A$  over belief space that optimizes the *value function*

$$V^\pi [b] = E \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$$

where the expectation  $E[\cdot]$  is taken over the sequence of states visited by  $\pi$ . We use  $V^*$  to denote the *optimal value function*, the value function of an optimal policy. The *discount factor*  $\gamma \in [0, 1)$  adjusts the relative value of present versus future reward.

In our problem, *state*  $s = (q, x_o) \in S$  is the configuration of the robot  $q \in Q$  and the pose of the movable object  $x_o \in X_o = SE(3)$  (fig. 2a). An *action*  $a = (\xi, T) \in A$  is a trajectory  $\xi : [0, T] \rightarrow Q$  that starts in configuration  $\xi(0)$  and ends configuration  $\xi(T)$  at time  $T$  (fig. 2b). We assume that uncertainty over the pose of the object dominates controller and proprioceptive error. Therefore, we treat  $q$  as a fully-observable state variable and neglect the dynamics of the manipulator; i.e. model it as being position controlled.

The robot executes a quasistatic, position controlled push if it comes in contact with the movable object. The *quasistatic assumption* states that friction is high enough to neglect the acceleration of the object; i.e. the object stops moving as soon as it leaves contact [34]. This approximation is accurate for many tabletop manipulation tasks [10, 11].

We define the stochastic *transition model*  $T(s, a, s')$  in terms of a deterministic quasistatic physics model [33] by introducing noise into the parameters [12]. We do not attempt to refine our estimate of these parameters during execution. After executing an action, the robot receives an *observation*  $o \in \{0, 1\}^{n_o} = O$  from its  $n_o$  binary contact sensors (fig. 2c). We assume that a stochastic *observation model*  $\Omega(s', a, o)$  is available, but make no assumptions about its form.

The robot's goal is to push the movable object into a hand-relative *goal region*  $X_{\text{goal}} \subseteq X$ . We encode this in the *reward function*

$$R(s, a) = \begin{cases} 0 & : [T_{\text{ee}}(q)]^{-1} x_o \in X_{\text{goal}} \\ -1 & : \text{otherwise} \end{cases}$$

where  $T_{\text{ee}} : Q \rightarrow SE(3)$  is the forward kinematics of the end effector. The reward function penalizes states where the movable object is outside of  $X_{\text{goal}}$ . Note that the choice of  $-1$  reward is arbitrary: any negative reward would suffice.

#### IV. CONFIGURATION LATTICE POMDP

Solving the POMDP formulated in section III is challenging because the state space is continuous, the action space is infinite dimensional, the transition model is expensive to

evaluate, and the optimal policy may perform long-horizon information gathering.

In this section, we simplify the problem by constraining the end effector to a fixed transformation relative to the support surface (section IV-A) and build a lattice in configuration space (section IV-B). Configurations in the lattice are connected by action templates that start and end on lattice points (section IV-C). Finally, we define a configuration lattice POMDP that penalizes infeasible actions to insure that the optimal policy will never take an infeasible action (section IV-E).

##### A. Planar Constraint

We assume that the robot's end effector is constrained to have a fixed transformation  ${}^{\text{sup}}T_{\text{ee}} \in SE(3)$  relative to the support surface  $T_{\text{sup}} \in SE(3)$ . A configuration  $q$  satisfies this constraint iff

$$T_{\text{ee}}(q) = T_{\text{sup}} {}^{\text{sup}}T_{\text{ee}} \text{Rot}(\theta_r, \hat{e}_z) \text{Trans}(x_r, y_r, 0) \quad (1)$$

where  $(x_r, y_r, \theta_r) \in X_r = SE(2)$  is the pose of the end effector in the plane.  $\text{Rot}(\theta, \hat{v})$  denotes a rotation about axis  $\hat{v}$  by angle  $\theta$  and  $\text{Trans}(v)$  denotes a translation by vector  $v$ .

We also assume that the movable object also has a fixed transformation  ${}^{\text{sup}}T_o \in SE(3)$  relative to the support surface; i.e. that it does not tip or topple. We parameterize its pose as

$$x_o = T_{\text{sup}} {}^{\text{sup}}T_o \text{Trans}(x_o, y_o, 0) \text{Rot}(\theta_o, \hat{e}_z)$$

where  $(x_o, y_o, \theta_o) \in X_o = SE(2)$  is its pose in the plane.

##### B. Configuration Lattice

We discretize the space of the end effector poses  $X_r$  by constructing a *state lattice*  $X_{r,\text{lat}} \subseteq X_r$  with a translational resolution of  $\Delta x_r, \Delta y_r \in \mathbb{R}^+$  and an angular resolution of  $\Delta \theta_r = 2\pi/n_\theta$  for some integer value of  $n_\theta \in \mathbb{N}$  [41]. The lattice consists of the discrete set of points

$$X_{r,\text{lat}} = \{(i_x \Delta x_r, i_y \Delta y_r, i_\theta \Delta \theta_r) : i_x, i_y, i_\theta \in \mathbb{Z}\}.$$

Each *lattice point*  $x_r \in X_{r,\text{lat}}$  may be reachable from multiple configurations. We assume that we have access to an *inverse kinematics function*  $q_{\text{lat}}(x_r)$  that returns a single solution  $\{q\}$  that satisfies  $T_{\text{ee}}(q) = x_r$  or  $\emptyset$  if no such solution exists. A solution may not exist if  $x_r$  is not reachable, the end effector is in collision, or the robot collides with the environment in all possible inverse kinematic solutions.

Instead of planning in  $S$ , we restrict ourselves to the state space  $S_{\text{lat}} = Q_{\text{lat}} \times X_o$  where  $Q_{\text{lat}} = \bigcup_{x_r \in X_{r,\text{lat}}} q_{\text{lat}}(x_r)$  is the discrete set of configurations returned by  $q_{\text{lat}}(\cdot)$  on all lattice points. Note that, despite the structure of the Cartesian lattice, planning occurs in configuration space  $Q$ , not the task space  $X_r$ .

##### C. Action Templates

Most actions do not transition between states in the lattice  $S_{\text{lat}}$ . Therefore, we restrict ourselves to action that are instantiated from one of a finite set  $A_{\text{lat}}$  of action templates. An action template  $a_{\text{lat}} = (\xi_{\text{lat}}, T) \in A_{\text{lat}}$  is a Cartesian trajectory  $\xi_{\text{lat}} : [0, T] \rightarrow SE(3)$  that specifies the relative motion of the



end effector. The template starts at the origin  $\xi_{\text{lat}}(0) = I$  and ends at some lattice point  $\xi_{\text{lat}}(T) \in X_{\text{r,lat}}$ . It is acceptable for multiple actions templates in  $A_{\text{lat}}$  to end at the same lattice point or have different durations.

An action  $a = (\xi, T) \in A$  satisfies template  $a_{\text{lat}}$  at lattice point  $x_r \in X_{\text{r,lat}}$  if it satisfies three conditions:

- 1) starts in configuration  $\xi(0) = q_{\text{lat}}(x_r)$
- 2) ends in configuration  $\xi(T) = q_{\text{lat}}(x_r \xi_{\text{lat}}(T))$
- 3) satisfies  $T_{\text{ee}}(\xi(\tau)) = x_r \xi_{\text{lat}}(\tau)$  for all  $0 \leq \tau \leq T$ .

These conditions are satisfied if  $\xi$  moves between two configurations in  $Q_{\text{lat}}$  and produces the same end effector motion as the Cartesian trajectory  $\xi_{\text{lat}}$ .

We define the function  $\text{Proj}(x_r, a) \mapsto a_{\text{lat}}$  to map an action  $a$  to the template  $a_{\text{lat}}$  it instantiates. The pre-image  $\text{Proj}^{-1}(x_r, a_{\text{lat}})$  contains the set of all possible instantiations of action template  $a_{\text{lat}}$  at  $x_r$ . We assume that we have access to a *local planner*  $\phi(q, a_{\text{lat}})$  that returns a singleton action  $\{a\} \subseteq \text{Proj}^{-1}(q, a_{\text{lat}})$  from this set or  $\emptyset$  to indicate failure. The local planner may fail due to kinematic constraints, end effector collision, or robot collisions.

#### D. Lattice POMDP

We use the lattice to define a *configuration lattice POMDP*, Lat-POMDP,  $(S_{\text{lat}}, A_{\text{lat}}, O, T_{\text{lat}}, \Omega_{\text{lat}}, R_{\text{lat}})$  with state space  $S_{\text{lat}} = X_{\text{r,lat}} \times X_o$ . The structure of the lattice guarantees that that all instantiations  $\text{Proj}^{-1}(q, a_{\text{lat}})$  of the action template  $a_{\text{lat}}$  execute the same Cartesian motion  $\xi_{\text{lat}}$  of the end effector. This motion is independent of the starting pose of the end effector  $x_r$  and configuration  $q_{\text{lat}}(x_r)$  of the robot.

If the movable object only contacts the end effector—no other parts of the robot or the environment—then the motion of the object is also independent of these variables. We refer to a violation of this assumption as *un-modelled contact*. The lattice transition model  $T_{\text{lat}}(s_{\text{lat}}, a_{\text{lat}}, s'_{\text{lat}})$  is identical to  $T(s, a, s')$  when  $a_{\text{lat}}$  is feasible and no un-modelled contact occurs. If either condition is violated, the robot deterministically transitions to an absorbing state  $s_{\text{invalid}}$ . Similarly, the lattice observation model  $\Omega_{\text{lat}}(s_{\text{lat}}, a_{\text{lat}}, o)$  is identical to  $\Omega(s, a, o)$  for valid states and is uniform over  $O$  for  $s_{\text{lat}} = s_{\text{invalid}}$ .

We penalize invalid states, infeasible actions, and un-modelled contact in the reward function

$$R_{\text{lat}}(s_{\text{lat}}, a_{\text{lat}}) = \begin{cases} -1 & : s_{\text{lat}} = s_{\text{invalid}} \\ -1 & : \phi(q_{\text{lat}}(x_r), a_{\text{lat}}) = \emptyset \\ -1 & : x_o \notin X_g \\ 0 & : \text{otherwise} \end{cases}$$

by assigning  $\min_{s \in S, a \in A} R(s, a) = -1$  reward to them.

#### E. Optimal Policy Feasibility

Our formulation of Lat-POMDP guarantees that an optimal policy  $\pi_{\text{lat}}^*$  will never take an infeasible action:

**Theorem 1.** *An optimal policy  $\pi_{\text{lat}}^*$  of Lat-POMDP will not execute an infeasible action in belief  $b$  if  $V_{\text{lat}}^*[b] > \frac{-1}{1-\gamma}$ .*

*Proof:* Suppose  $V_{\text{lat}}^*[b] > \frac{-1}{1-\gamma}$  and an optimal policy  $\pi_{\text{lat}}^*$  executes the invalid action in belief state  $b$ . The robot receives

a reward of  $-1$  and transitions to  $s_{\text{invalid}}$ . For all time after that, regardless of the actions that  $\pi_{\text{lat}}^*$  takes, the robot receives a reward of  $R_{\text{lat}}(s_{\text{invalid}}, \cdot) = -1$  at each timestep. This yields a total reward of  $V_{\text{lat}}^{\pi_{\text{lat}}^*} = \frac{-1}{1-\gamma}$ , which is the minimum reward possible to achieve.

The value function of the optimal policy satisfies the Bellman equation  $V_{\text{lat}}^*[b] = \arg \max_{a_{\text{lat}} \in A_{\text{lat}}} Q^*[b, a_{\text{lat}}]$ , where  $Q^*[b, a]$  denotes the value of taking action  $a$  in belief state  $b$ , then following the optimal policy for all time. This contradicts the fact that  $V_{\text{lat}}^*[b] > \frac{-1}{1-\gamma}$  and  $V_{\text{lat}}^{\pi_{\text{lat}}^*}[b] = \frac{-1}{1-\gamma}$ . Therefore,  $\pi_{\text{lat}}^*$  must not be the optimal policy. ■

We can strengthen our claim if (1) we guarantee that every lattice point reachable from the  $q_0$  has at least one feasible action and (2) it is possible to achieve the goal with non-zero probability. Under those assumptions we know that  $V_{\text{lat}}^*[b] > \frac{-1}{1-\gamma}$  and theorem 1 guarantees that  $\pi_{\text{lat}}^*$  will never take an infeasible action. One simple way to satisfy condition (1) is to require that all actions are reversible; i.e. make the lattice an undirected graph.

### V. ONLINE POMDP PLANNER

Lat-POMDP has a reward function that changes whenever obstacles are added to, removed from, or moved within the environment. An offline POMDP solver, like point-based value iteration [40] or SARSOP [28] would require re-computing the optimal policy for each problem instance.

Instead, we use DESPOT, a state-of-the-art online POMDP solver that uses regularization to balance the size of the policy against its quality [48]. DESPOT incrementally explores the action-observation tree rooted at  $b(s_0)$  by performing a series of trials. Each *trial* starts at the root node, descends the tree, and terminates by adding a new leaf node to the tree.

In each step, DESPOT chooses the action that maximizes the upper bound  $\bar{V}[b]$  and the observation that maximizes *weighted excess uncertainty*, a regularized version of the gap  $\bar{V}[b] - \underline{V}[b]$  between the upper and lower bounds. This search strategy heuristically focuses exploration on the optimally reachable belief space [28]. Finally, DESPOT backs up the upper and lower bounds of all nodes visited by the trial.

There are key two challenges in solving Lat-POMDP:

First, we must construct the configuration lattice. We use lazy evaluation to interleave lattice construction with planning. By doing so, we only evaluate the parts of the lattice that are visited by DESPOT (section V-A).

Second, we must provide DESPOT with upper and lower bounds on the optimal value function to guide its search. We derive these bounds from a relaxed version of the problem that ignores obstacles by only considering the pose of the movable object relative to the hand (sections V-B to V-D).

#### A. Lattice Construction

DESPOT uses upper and lower bounds to focus its search on belief states that are likely to be visited by the optimal policy. We exploit this fact to avoid constructing the entire lattice. Instead, we interleave lattice construction with planning and

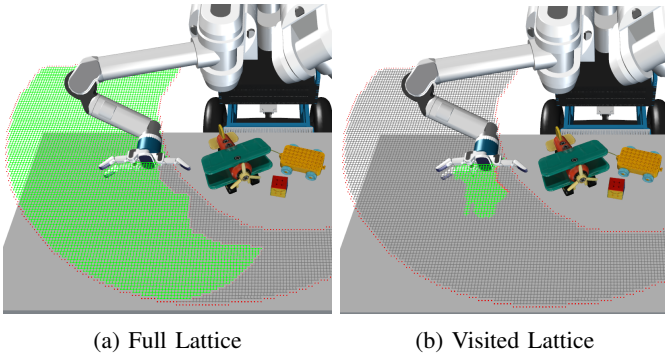


Fig. 3: Comparison of the (a) full lattice to the (b) subset visited by DESPOT, with feasible (—), infeasible (—) and unevaluated (—) edges. We interleave planning with lattice construction to construct only the subset of the lattice visited by the search. This figure is best viewed in color.

only instantiate the lattice edges visited by the search, similar to the concept *lazy evaluation* used in motion planning [3, 15].

We begin with no pre-computation and run DESPOT until it queries the transition model  $T_{\text{lat}}$ , observation model  $\Omega_{\text{lat}}$ , or reward function  $R_{\text{lat}}$  for a state-action pair  $(x_r, a_{\text{lat}})$  that has not yet been evaluated. When this occurs, we pause the search and check the feasibility of the action by running the local planner  $\phi(x_r, a_{\text{lat}})$ . We use the outcome of the local planner to update the Lat-POMDP model and resume the search.

Figure 3 shows the (a) full lattice and (b) subset evaluated by DESPOT. Feasible action templates are shown as green, infeasible action templates are shown as red, and unevaluated action templates are shown as gray. Note that DESPOT evaluates only a small number of state-action pairs, obviating the need for performing computationally expensive collision checks on much of the lattice.

Note that it is also possible to use a hybrid approach by evaluating some parts of the lattice offline and deferring others to be computed online. For example, we may compute inverse kinematics solutions, kinematic feasibility checks, and self-collision checks in an offline pre-computation step. These values are fixed for a given support surface and, thus, can be used across multiple problem instances. Other feasibility tests, e.g. collision with obstacles in the environment, are lazily evaluated online.

### B. Hand-Relative POMDP

Recall from section IV that the motion of the end effector and the object is independent of the pose of the end effector  $x_r$  or the robot configuration  $q$ . We use this insight to define a *hand-relative POMDP* Rel-POMDP  $(S_{\text{rel}}, A_{\text{lat}}, O, T_{\text{rel}}, \Omega_{\text{rel}}, R_{\text{rel}})$  with a state space that only includes the pose  $x_{o,\text{rel}} = x_r^{-1}x_o \in S_{\text{rel}}$  of the movable object relative to the hand. The hand-relative transition model  $T_{\text{rel}}$ , observation model  $\Omega_{\text{rel}}$ , and reward function  $R_{\text{rel}}$  are identical to the original model when no un-modelled contact occurs. Rel-POMDP is identical to the hand-relative POMDP models used in prior work [17, 26, 27] and is equivalent to assuming

that environment is empty and the robot is a lone end effector actuated by an incorporeal planar joint.

### C. Lat-POMDP Upper Bound

Rel-POMDP is a relaxation of Lat-POMDP that treats all actions as feasible. As such:

**Theorem 2.** *The optimal value function  $V_{\text{rel}}^*$  of Rel-POMDP is an upper bound on the optimal value function  $V_{\text{lat}}^*$  of Lat-POMDP:  $V_{\text{rel}}^*[b] \geq V_{\text{lat}}^*[b]$  for all  $b \in \Delta$ .*

*Proof:* We define a *scenario*  $\psi = (s_0, \psi_1, \psi_2, \dots)$  as an abstract simulation trajectory that captures all uncertainty in our POMDP model [35, 48]. A scenario is generated by drawing the initial state  $s_o \sim b(s_0)$  from the initial belief state and each random number  $\psi_1 \sim \text{uniform}[0, 1]$  uniformly from the unit interval. Given a scenario  $\psi$ , we assume that the outcome of executing a sequence of actions is deterministic; i.e. all stochasticity is captured in the initial state  $s_0$  and the sequence of random numbers  $\psi_1, \psi_2, \dots$ .

Suppose we have a policy  $\pi$  for Rel-POMDP that executes the sequence of actions  $a_{\text{lat},1}, a_{\text{lat},2}, \dots$  in scenario  $\psi$ . The policy visits the sequence of states  $s_{\text{rel},1}, s_{\text{rel},2}, \dots$  and receives the sequence of rewards  $R_1, R_2, \dots$ .

Now consider executing  $\pi$  in the same scenario  $\psi$  on Lat-POMDP. Without loss of generality, assume that  $\pi$  first takes an infeasible action or makes un-modelled contact with the environment at timestep  $H$ . The policy receives the same sequence of rewards  $R_1, R_2, \dots, R_2, \dots, R_{H-1}, -1, -1, \dots$  as it did on Rel-POMDP until timestep  $H$ . Then, it receives  $-1$  reward for taking an infeasible action, transitions to the absorbing state  $s_{\text{invalid}}$ , and receives  $-1$  reward for all time.

Policy  $\pi$  achieves value  $V_{\text{rel},\psi}^\pi = \sum_{t=0}^{\infty} \gamma^t R_t$  on Rel-POMDP and  $V_{\text{lat},\psi}^\pi = \sum_{t=0}^{H-1} \gamma^t R_t - \frac{\gamma^H}{1-\gamma}$  on Lat-POMDP in scenario  $\psi$ . Since  $R_t \geq -1$ , we know that  $V_{\text{rel},\psi}^\pi \geq V_{\text{lat},\psi}^\pi$ . The value of a policy  $V^\pi = E_\psi[V_\psi^\pi]$  is the expected value of  $\pi$  over all scenarios.

Consider the optimal policy  $\pi_{\text{lat}}^*$  of Lat-POMDP. There exists some Rel-POMDP policy  $\pi_{\text{mimic}}$  that executes the same sequence of actions as  $\pi_{\text{lat}}^*$  in all scenarios. From the reasoning above, we know that  $V_{\text{rel}}^{\pi_{\text{mimic}}} \geq V_{\text{lat}}^*$ . We also know that  $V_{\text{rel}}^* \geq V_{\text{rel}}^{\pi_{\text{mimic}}}$  because the value of any policy is a lower bound on the optimal value function.

Therefore,  $V_{\text{rel}}^* \geq V_{\text{rel}}^{\pi_{\text{mimic}}} \geq V_{\text{lat}}^*$ ; i.e. the optimal value function  $V_{\text{rel}}^*$  of Rel-POMDP is an upper bound on the optimal value function  $V_{\text{lat}}^*$  of Lat-POMDP. ■

This result implies that *any* upper bound  $\bar{V}_{\text{rel}}$  is an upper bound on the value of the optimal value function  $\bar{V}_{\text{rel}} \geq V_{\text{rel}}^* \geq V_{\text{lat}}^*$ . Therefore, we may also use  $\bar{V}_{\text{rel}}$  as an upper bound on Lat-POMDP. The key advantage of doing so is that the  $\bar{V}_{\text{rel}}$  may be pre-computed once per hand-object pair. In contrast, the same upper bound on  $\bar{V}_{\text{lat}}$  must be re-computed for each problem instance.

### D. Lower Bound

We also use Rel-POMDP as a convenient method of constructing a lower bound  $\underline{V}_{\text{lat}}$  for Lat-POMDP. As explained

above, the value of any policy is a lower bound on the optimal value function. We use offline pre-computation to compute a rollout policy on  $\pi_{\text{rollout}}$  for Rel-POMDP once per hand-object pair. For example, we could use MDP value iteration to compute a QMDP policy [32] or a point-based method [28, 40] to find a near-optimal policy.

Given an arbitrary policy  $\pi_{\text{rollout}}$  computed in this way, we construct an approximate lower bound  $\underline{V}_{\text{lat}}$  for Lat-POMDP by estimating the value  $V_{\text{lat}}^{\pi_{\text{rollout}}}$  of executing  $\pi_{\text{rollout}}$  on Lat-POMDP using rollouts. Approximating a lower bound with a *rollout policy* is commonly used in POMCP [46], DESPOT [48], and other online POMDP solvers.

## VI. EXPERIMENTAL RESULTS

We validated the efficacy of the proposed algorithm by running simulation experiments on a robot equipped with a 7-DOF Barrett WAM arm [45] and the BarrettHand [51] end-effector. The robot attempts to push a bottle into the center of its palm on a table littered with obstacles.

### A. Problem Definition

The state space of the problem consists of the configuration space  $Q = \mathbb{R}^7$  of the robot and the pose of the object  $X_o$  relative to the end effector. The robot begins in a known initial configuration  $q_0$  and the initial pose of the bottle  $x_o$  is drawn from a Gaussian distribution centered in front of the palm with a covariance matrix of  $\Sigma^{1/2} = \text{diag}[5 \text{ mm}, 10 \text{ cm}]$ . At each timestep, the robot chooses an action  $a_{\text{lat}}$  that moves 1 cm at a constant Cartesian velocity in the  $xy$ -plane, receives an observation  $o$  from its  $n_o = 2$  fingertip contact sensors, and receives a reward  $R$ .

The goal is to push the object into the  $4 \text{ cm} \times 6 \text{ cm}$  goal region  $X_{\text{goal}}$  centered in front of the palm. We evaluate the performance of the policy by: (1) computing the probability that the object is in the goal region at each timestep and (2) computing the discounted sum of reward with  $\gamma = 0.99$ . If the robot takes an infeasible action, the simulation immediately terminates and the robot receives  $-1$  reward for all remaining timesteps

1) *Transition Model*: The motion of the object is assumed to be quasistatic [33] and is simulated using the Box2D physics simulator [7]. We simulate uncertainty in the model by sampling the hand-object friction coefficient and center of the object-table pressure distribution from Gaussian distributions at each timestep [12, 26].

2) *Observation Model*: After each timestep the robot receives a binary observation from a contact sensor on each of its fingertips. We assume that the sensors perfectly discriminate between contact and no-contact [26, 27], but provide no additional information about where contact occurred on the sensor. The robot must take information-gathering actions, by moving side-to-side, to localize the object relative to the hand.

3) *Rel-POMDP Discretization*: We discretize  $S_{\text{rel}}$  with a 1 cm resolution over a region of size  $20 \text{ cm} \times 44 \text{ cm}$  centered around the palm. We compute a transition model, observation model, and reward function for Rel-POMDP by taking the

expectation over the continuous models by assuming that each discrete state represents a uniform distribution over the corresponding area of the continuous state space. We use these discrete models in both our Rel-POMDP and Lat-POMDP experiments.

States that leave the discretized region around the hand are treated as un-modelled contact. As in prior work [17, 27], discretization is necessary to speed up evaluation of the model and to enable calculation of the QMDP and SARSOP baseline policies described below.

### B. State Lattice Construction

The robot's actions tessellate  $X_r$  into a lattice centered at  $T_{\text{ee}}(q_0)$  with a resolution of  $\Delta x_r = \Delta y_r = 1 \text{ cm}$ . First, we select a configuration  $q_{\text{lat}}(x_r)$  using an iterative inverse kinematics solver initialized with the solution of an adjacent lattice point. Then, we use a Cartesian motion planner to find a trajectory that connects adjacent points while satisfying the constraints imposed by the action templates. In most configurations, the 1 cm lattice is sufficiently dense to simply connect adjacent lattice points with a straight-line trajectory in configuration space. Forward kinematics, inverse kinematics, and collision detection is provided by the Dynamic Animation and Robotics Toolkit (DART) [1].

As described in section V-A, the kinematic structure of the lattice is computed offline for the height of the support surface, but no collision checking is performed. Collision checks are deferred to runtime when the planner queries the feasibility of an action template. At that point, the edge is checked for collision against the environment using the Flexible Collision Library (FCL) [37].

### C. Policies

We compare the quality of the policy produced by the proposed algorithm (Lat-DESPOT) against several policies:

1) *Rel-QMDP*: Choose the action at each timestep that greedily optimizes the single-step  $Q$ -value of the MDP value function associated from Rel-POMDP [32]. QMDP is optimal if all uncertainty disappears after executing one action and, as a result, does not perform multi-step information-gathering actions. QMDP is commonly used in robotic applications due to its speed, simplicity, and good performance in domains where information is easily gathered [13, 22].

2) *Rel-SARSOP*: Compute a near-optimal policy for Rel-POMDP using SARSOP, an offline point-based method [28]. Rel-SARSOP serves as a baseline to demonstrate that information-gathering is beneficial and to verify that DESPOT is capable of finding a near-optimal solution. SARSOP has been shown to perform well on Rel-POMDP in prior work [17, 27]. As in that work, we used the implementation of SARSOP provided by the APPL toolkit and allowed it to run for 10 minutes.

3) *Rel-DESPOT*: Use DESPOT to find a near-optimal policy for Rel-POMDP [48]. The solver uses Rel-QMDP as an upper bound and rollouts of Rel-QMDP as a lower bound (see section V-D for an explanation of a rollout policy). We



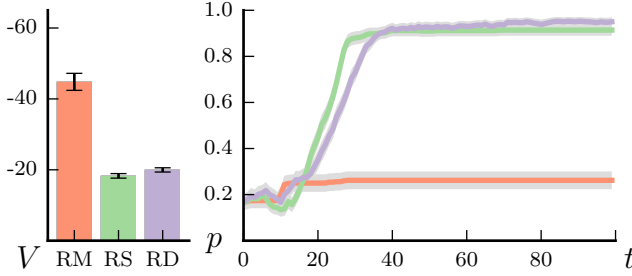


Fig. 4: Performance of the Rel-QMDP (RM ■ —), Rel-SARSOP (RS ■ —), and Rel-DESPOT (RD ■ —) policies on Rel-POMDP. (Left) Value  $V_{\text{rel}}$  achieved by each policy computed over 100 timesteps. Note that the  $y$ -axis is inverted; lower (less negative) is better. (Right) Probability that the movable object is in  $X_{\text{goal}}$  as a function of timestep. Error bars denote a 95% confidence interval. Best viewed in color.

use the implementation of DESPOT provided by the APPL toolkit and tuned the number of trials, number of scenarios, regularization constant, and gap constant on a set of training problem instances that are distinct from the results presented in this section.

4) *Lift-QMDP and Lift-SARSOP*: Use the state lattice to evaluate the feasibility of the action returned by Rel-QMDP and Rel-SARSOP, respectively, before executing it. If the desired action is infeasible, instead execute the action with the highest estimated  $Q$ -value. This represents a heuristic solution for modifying a Rel-POMDP policy to avoid taking infeasible actions.

5) *Lat-DESPOT*: The proposed algorithm described in section V. We run DESPOT [48] on Lat-POMDP using Rel-QMDP as the upper bound and rollouts of Lift-QMDP as the lower bound (see section V-D for an explanation of a rollout policy). Just as with Rel-DESPOT, we used the implementation of DESPOT provided by the APPL toolkit and tuned all parameters on a set of training problem instances.

#### D. Rel-POMDP Experiments

We begin with a preliminary experiment on Rel-POMDP. Our goal is to: (1) validate that discretizing Rel-POMDP does not harm the performance of the optimal policy on the underlying continuous problem, (2) confirm that information-gathering is necessary, (3) verify that Rel-DESPOT does not perform worse than Rel-SARSOP, and (4) estimate  $V_{\text{rel}}^*[b_0]$ .

1) *Discretization Validation*: Figure 4 shows simulation results averaged over 500 instances of the experiment described above. Figure 4-Left shows the value of each policy achieved in 100 timesteps on simulated on the discretized Rel-POMDP problem described in section VI-A3. Figure 4-Right shows the probability that the movable object is in  $X_{\text{goal}}$  at each timestep when simulated using the continuous model. The close agreement between the results suggests that discretizing the state space does not harm the performance of the policy on the underlying continuous problem.

2) *Necessity of Information-Gathering*: Rel-QMDP (■ —) performs poorly on this problem, achieving  $< 30\%$  success probability, because it pushes straight and does not attempt to localize the object. Rel-SARSOP (■ —) and Rel-DESPOT (■ —) execute information-gathering action by moving the hand laterally to drive the movable object into one of the fingertip contact sensors. Once the object has been localized, the policy successfully pushes it into the goal region. These results are consistent with prior work [17, 27] and confirm that information-gathering is necessary to perform well on this problem.

3) *Rel-DESPOT Policy*: Our intuition is that it is more difficult to solve Lat-POMDP than Rel-POMDP. Therefore, it is important that we verify that DESPOT can successfully solve Rel-POMDP before applying it to Lat-POMDP. Our results confirm this is true: Rel-DESPOT (■ —) achieves comparable value and success probability to Rel-SARSOP (■ —).

#### E. Lat-POMDP Experiments

We evaluate the proposed approach (Lat-DESPOT) on Lat-POMDP in four different environments: (a) an empty table, (b) obstacles on the right, (c) obstacles on the left, and (d) more complex obstacles on the right. Note that kinematic constraints are present in all four environments, even the empty table, in the form of reachability limits, self-collision, and collision between the arm and the table. Scenes (b), (c), and (d) are constructed out of objects selected from the YCB dataset [6].

Figure 5 shows results for each scene averaged over 500 instances of the problem. Just as in the Rel-POMDP experiments, Figure 5-Top shows the value  $V_{\text{lat}}$  achieved by each policy when evaluated on the discretized version of Lat-POMDP. Figure 5-Middle shows the probability that the movable object is in  $X_{\text{goal}}$  at each timestep, treating instances that have terminated as zero probability. Figure 5-Bottom shows the proportion of Rel-QMDP, Rel-SARSOP, and Rel-DESPOT policies that are active at each timestep; i.e. have not yet terminated by taking an infeasible action. Note that all six policies—even Rel-MDP, Rel-SARSOP, and Rel-DESPOT—are subject to kinematic constraints during execution.

1) *Baseline Policies*: Rel-QMDP (RM ■ —) and Lift-QMDP (LM ■ —) to perform poorly across all environments, achieving  $< 30\%$  success probability, because they do not take multi-step information-gathering actions. Figure 5 confirms this: both QMDP policies perform poorly on all four environments.

Rel-SARSOP (RS ■ —) and Rel-DESPOT (RD ■ —) perform well on environments (a) and (b) because they hit obstacles late in execution. The converse is true on environments (c) and (d): both policies hit obstacles so quickly that they perform worse than Rel-QMDP!

Lift-SARSOP (LS ■ —) performs near-optimally on environments (a) and (b) because it: (1) does not take infeasible actions and (2) gathers information. However, it performs no better than Rel-QMDP on problem (d). This occurs because Lift-SARSOP myopically considers obstacles in a one-step

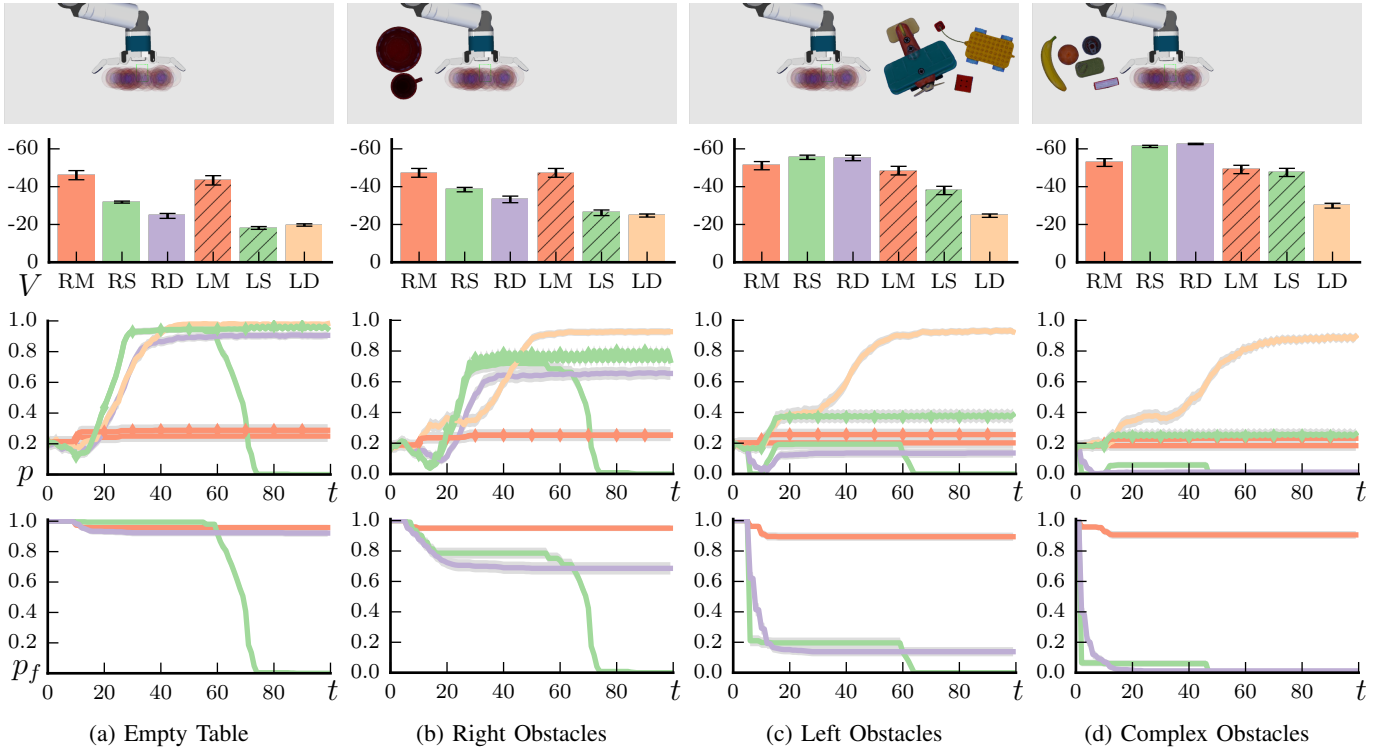


Fig. 5: Performance of the Rel-QMDP (RM  $\text{---}$ ), Rel-SARSOP (RS  $\text{---}$ ), Rel-DESPOT (RD  $\text{---}$ ), Lift-QMDP (LM  $\text{---}$ ), Lift-SARSOP (LS  $\text{---}$ ), and Lat-DESPOT (LD  $\text{---}$ ) policies on Lat-POMDP on four environments. (Top) Value achieved by each policy after 100 timesteps. Note that the  $y$ -axis is inverted; lower (less negative) is better. (Middle) Probability  $p = \Pr(s_t \in X_{\text{goal}})$  that the movable object is in the goal region at each timestep. (Bottom) Probability that the execution is feasible as a function of time. Lift-QMDP, Lift-SARSOP, and Lat-DESPOT are omitted because they do not take infeasible actions. In all cases, error bars denote a 95% confidence interval. Best viewed in color.

lookahead and may oscillate when blocked. Small changes in the environment are sufficient to induce this behavior: the key difference between environments (b) and (d) is the introduction of a red box that creates a cul-de-sac in the lattice.

2) *Lat-DESPOT Policy*: Our proposed approach, Lat-DESPOT ( $\text{---}$ ), avoids myopic behavior by considering action feasibility during planning. Lat-DESPOT performs no worse than Lift-SARSOP on environments (a) and (b) and outperforms it on environments (c) and (d). Unlike Rel-SARSOP, Lat-DESPOT identifies the cul-de-sac in (d) during planning and avoids becoming trapped in it. In summary, *Lat-DESPOT is the only policy that performs near-optimally on all four environments*.

Our unoptimized implementation of Lat-DESPOT took between 200  $\mu\text{s}$  and 2.4 s to select an action (a 1 cm motion of the end-effector) on a single core of a 4 GHz Intel Core i7 CPU. The policy was slowest to evaluate early in execution, when information-gathering is necessary, and fastest once the movable object is localized; i.e. once the upper and lower bounds become tight. The QMDP and SARSOP policies, which perform no planning online, took an average of 1.6  $\mu\text{s}$  and 218  $\mu\text{s}$  respectively. We are optimistic about achieving real-time performance from Lat-DESPOT by parallelizing scenario rollouts, optimizing our implementation of the algorithm,

and leveraging the loose coupling between the optimal pre- and post-contact policies [27].

3) *Upper Bound Validation*: Finally, we combine the data in fig. 4-Left and fig. 5-Top to empirically verify the bound we proved in theorem 2; i.e. the optimal value function  $V_{\text{rel}}^*$  of Rel-POMDP is an upper bound on the optimal value function  $V_{\text{lat}}^*$  of Lat-POMDP. Note that the value of Rel-SARSOP ( $\text{---}$ ) and Rel-DESPOT ( $\text{---}$ ) on Rel-POMDP (fig. 4-Left) are greater (i.e. less negative) than the value of all policies we evaluated on Lat-POMDP (fig. 5-Top). The data supports our theory: the optimal value achieved on Rel-POMDP is no worse than the highest value achieved on Lat-POMDP in environment (a) and greater than the highest value achieved in environments (b), (c), and (d).

## VII. DISCUSSION

In this paper, we formulated the problem of planar contact manipulation under uncertainty as a POMDP in the joint space of robot configurations and poses of the movable object (section I). We simplify the problem by constructing a lattice in the robot's configuration space and prove that, under mild assumptions, the optimal policy of Lat-POMDP will never take an infeasible action (section IV). We find a near-optimal policy for Lat-POMDP using DESPOT [48], a state-of-the-art online



POMDP solver, guided by upper and lower bounds derived from Rel-POMDP (section V).

Our simulation results show that Lat-DESPOT outperforms five baseline algorithms on cluttered environments: it achieves a  $> 90\%$  success rate on all environments, compared to the best baseline (Lift-SARSOP) that achieves only a  $20\%$  success rate on difficult problems. They also highlight the importance of reasoning about both object pose uncertainty and kinematic constraints during planning. Lat-DESPOT is a promising first step towards using recent advances in online POMDP solvers, like DESPOT [48], to achieve that goal. However, Lat-DESPOT has several limitations that we plan to address in future work:

First, our approach assumes that the robot has perfect proprioception and operates in an environment with known obstacles. In practice, robots often have imperfect proprioception [4, 24] and uncertainty about the pose of *all* objects in the environment. We hope to relax both of these assumptions by replacing the deterministic transition model for robot configuration with a stochastic model that considers the probability of hitting an obstacle.

Second, we are excited to scale our approach up a larger repertoire of action templates (including non-planar motion), solving more complex tasks, and planning in environments that contain multiple movable objects. Solving these more complex problems will require more informative heuristics. We are optimistic that more sophisticated Rel-POMDP policies, e.g. those computed by a point-based method [43] or Monte Carlo Value Iteration [2], could be used to guide the search.

Third, our approach commits to a single inverse kinematics solution  $q_{\text{lat}}(x_r)$  for each lattice point. This prevents robots like HERB, which has a seven degree-of-freedom manipulator, from using redundancy to avoid kinematic constraints. We plan to relax this assumption in future work by generating multiple inverse kinematic solutions for each lattice point and instantiating an action template for each. Our intuition is that many solutions share the same connectivity and, thus, may be treated identically during planning.

Finally, we plan to implement Lat-DESPOT on a real robotic manipulator and evaluate the performance of our approach on real-world manipulation tasks.

#### ACKNOWLEDGEMENTS

This work was supported by a NASA Space Technology Research Fellowship (award NNX13AL62H), the National Science Foundation (awards IIS-1218182 and IIS-1409003), the U.S. Office of Naval Research, and the Toyota Motor Corporation. We would like to thank Rachel Holladay, Shervin Javdani, Jennifer King, Stefanos Nikolaidis, and the members of the Personal Robotics Lab for their helpful input. We would also like to thank Nan Ye for assistance with the APPL toolkit.

#### REFERENCES

- [1] Dynamic Animation and Robotics Toolkit. <http://dartsim.github.io>, 2013.
- [2] H. Bai, D. Hsu, W.S. Lee, and V.A. Ngo. Monte Carlo value iteration for continuous-state POMDPs. In *Workshop on the Algorithmic Foundations of Robotics*, 2011.
- [3] R. Bohlin and L.E. Kavraki. Path planning using lazy PRM. In *IEEE International Conference on Robotics and Automation*, pages 521–528, 2000.
- [4] B. Boots, A. Byravan, and D. Fox. Learning predictive models of a depth camera & manipulator from raw execution traces. In *IEEE International Conference on Robotics and Automation*, 2014.
- [5] M. Brokowski, M. Peshkin, and K. Goldberg. Curved fences for part alignment. In *IEEE International Conference on Robotics and Automation*, 1993. doi: 10.1109/ROBOT.1993.292216.
- [6] B. Calli, A. Singh, A. Walsman, S.S. Srinivasa, P. Abbeel, and A.M. Dollar. The YCB object and model set: Towards common benchmarks for manipulation research. In *International Conference on Advanced Robotics*, 2015.
- [7] E. Catto. Box2D. <http://box2d.org>, 2010.
- [8] B. Cohen, S. Chitta, and M. Likhachev. Single- and dual-arm motion planning with heuristic search. *International Journal of Robotics Research*, 2013.
- [9] M. Dogar, K. Hsiao, M. Ciocarlie, and S.S. Srinivasa. Physics-based grasp planning through clutter. In *Robotics: Science and Systems*, 2012.
- [10] M.R. Dogar and S.S. Srinivasa. Push-grasping with dexterous hands: Mechanics and a method. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2010. doi: 10.1109/IROS.2010.5652970.
- [11] M.R. Dogar and S.S. Srinivasa. A planning framework for non-prehensile manipulation under clutter and uncertainty. *Autonomous Robots*, 33(3):217–236, 2012. doi: 10.1007/s10514-012-9306-z.
- [12] D.J. Duff, J. Wyatt, and R. Stolkin. Motion estimation using physical simulation. In *IEEE International Conference on Robotics and Automation*, 2010. doi: 10.1109/ROBOT.2010.5509590.
- [13] R. Emery-Montemerlo, G. Gordon, J. Schneider, and S. Thrun. Approximate solutions for partially observable stochastic games with common payoffs. In *Autonomous Agents and Multiagent Systems*, 2004.
- [14] M.A. Erdmann and M.T. Mason. An exploration of sensorless manipulation. *IEEE Journal of Robotics and Automation*, 1988. doi: 10.1109/56.800.
- [15] K. Hauser. Lazy collision checking in asymptotically-optimal motion planning. In *IEEE International Conference on Robotics and Automation*, 2015.
- [16] P. Hebert, T. Howard, N. Hudson, J. Ma, and J.W. Burdick. The next best touch for model-based localization. In *IEEE International Conference on Robotics and Automation*, 2013. doi: 10.1109/ICRA.2013.6630562.

- [17] M. Horowitz and J. Burdick. Interactive non-prehensile manipulation for grasping via POMDPs. In *IEEE International Conference on Robotics and Automation*, 2013. doi: 10.1109/ICRA.2013.6631031.
- [18] K. Hsiao. *Relatively robust grasping*. PhD thesis, Massachusetts Institute of Technology, 2009.
- [19] K. Hsiao, L.P. Kaelbling, and T. Lozano-Pérez. Grasping POMDPs. In *IEEE International Conference on Robotics and Automation*, 2007. doi: 10.1109/ROBOT.2007.364201.
- [20] K. Hsiao, T. Lozano-Pérez, and L.P. Kaelbling. Robust belief-based execution of manipulation programs. In *Workshop on the Algorithmic Foundations of Robotics*, 2008.
- [21] S. Javdani, M. Klingensmith, J.A. Bagnell, N.S. Pollard, and S.S. Srinivasa. Efficient touch based localization through submodularity. In *IEEE International Conference on Robotics and Automation*, 2013. doi: 10.1109/ICRA.2013.6630818.
- [22] S. Javdani, J.A. Bagnell, and S.S. Srinivasa. Shared autonomy via hindsight optimization. In *Robotics: Science and Systems*, 2015.
- [23] L.P. Kaelbling, M.L. Littman, and A.R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 1998. doi: 10.1016/S0004-3702(98)00023-X.
- [24] M. Klingensmith, T. Galluzzo, C. Dellin, M. Kazemi, J.A. Bagnell, and N. Pollard. Closed-loop servoing using real-time markerless arm tracking. In *IEEE International Conference on Robotics and Automation/Humanoids Workshop*, 2013.
- [25] M.C. Koval, J.E. King, N.S. Pollard, and S.S. Srinivasa. Robust trajectory selection for rearrangement planning as a multi-armed bandit problem. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2015.
- [26] M.C. Koval, N.S. Pollard, and S.S. Srinivasa. Pose estimation for planar contact manipulation with manifold particle filters. *International Journal of Robotics Research*, 34(7):922–945, 2015. doi: 10.1177/0278364915571007.
- [27] M.C. Koval, N.S. Pollard, and S.S. Srinivasa. Pre- and post-contact policy decomposition for planar contact manipulation under uncertainty. *International Journal of Robotics Research*, 2015. doi: 10.1177/0278364915594474. In press.
- [28] H. Kurniawati, D. Hsu, and W.S. Lee. SARSOP: Efficient point-based POMDP planning by approximating optimally reachable belief spaces. In *Robotics: Science and Systems*, 2008.
- [29] S.M. LaValle and S.A. Hutchinson. An objective-based framework for motion planning under sensing and control uncertainties. *International Journal of Robotics Research*, 1998. doi: 10.1177/027836499801700104.
- [30] Q. Li, C. Schürmann, R. Haschke, and H. Ritter. A control framework for tactile servoing. In *Robotics: Science and Systems*, 2013.
- [31] M. Likhachev and D. Ferguson. Planning long dynamically feasible maneuvers for autonomous vehicles. *International Journal of Robotics Research*, 28(8):933–945, 2009.
- [32] M.L. Littman, A.R. Cassandra, and L.P. Kaelbling. Learning policies for partially observable environments: Scaling up. *International Conference on Machine Learning*, 1995.
- [33] K.M. Lynch, H. Maekawa, and K. Tanie. Manipulation and active sensing by pushing using tactile feedback. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 1992. doi: 10.1109/IROS.1992.587370.
- [34] M.T. Mason. Mechanics and planning of manipulator pushing operations. *International Journal of Robotics Research*, 5(3):53–71, 1986. doi: 10.1177/027836498600500303.
- [35] A.Y. Ng and M. Jordan. PEGASUS: A policy search method for large MDPs and POMDPs. In *Conference on Uncertainty in Artificial Intelligence*, 2000.
- [36] S.C.W. Ong, S.W. Png, D. Hsu, and W.S. Lee. POMDPs for robotic tasks with mixed observability. In *Robotics: Science and Systems*, June 2009.
- [37] Jia Pan, Sachin Chitta, and Dinesh Manocha. FCL: A general purpose library for collision and proximity queries. In *IEEE International Conference on Robotics and Automation*, pages 3859–3866, 2012.
- [38] P. Pastor, L. Righetti, M. Kalakrishnan, and S. Schaal. Online movement adaptation based on previous sensor experiences. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2011. doi: 10.1109/IROS.2011.6095059.
- [39] A. Petrovskaya and O. Khatib. Global localization of objects via touch. *IEEE Transactions on Robotics*, 27(3):569–585, 2011. doi: 10.1109/TRO.2011.2138450.
- [40] J. Pineau, G. Gordon, and S. Thrun. Point-based value iteration: An anytime algorithm for POMDPs. In *International Joint Conference on Artificial Intelligence*, 2003.
- [41] M. Pivtoraiko and A. Kelly. Efficient constrained path planning via search in state lattices. In *International Symposium on Artificial Intelligence, Robotics, and Automation in Space*, 2005.
- [42] R. Platt, A.H. Fagg, and R.A. Grupen. Nullspace grasp control: theory and experiments. *IEEE Transactions on Robotics*, 26(2):282–295, 2010. doi: 10.1109/TRO.2010.2042754.
- [43] J.M. Porta, N. Vlassis, M.T.J. Spaan, and P. Poupart. Point-based value iteration for continuous POMDPs. *Journal of Machine Learning Research*, 7:2329–2367, 2006.
- [44] S. Ross, J. Pineau, S. Paquet, and B. Chaib-Draa. Online planning algorithms for POMDPs. *Journal of Artificial Intelligence Research*, 2008. doi: 10.1613/jair.2567.
- [45] K. Salisbury, W. Townsend, B. Eberman, and D. DiPietro. Preliminary design of a whole-arm manipulation system (WAMS). In *IEEE International Conference on Robotics and Automation*, 1988. doi: 10.1109/ROBOT.

1988.12057.

- [46] D. Silver and J. Veness. Monte-Carlo planning in large POMDPs. In *Advances in Neural Information Processing Systems*, 2010.
- [47] R.D. Smallwood and E.J. Sondik. The optimal control of partially observable Markov processes over a finite horizon. *Operations Research*, 21(5):1071–1088, 1973. doi: 10.1287/opre.21.5.1071.
- [48] A. Somani, N. Ye, D. Hsu, and W.S. Lee. DESPOT: On-line POMDP planning with regularization. In *Advances in Neural Information Processing Systems*, 2013.
- [49] S.S. Srinivasa, D. Berenson, M. Cakmak, A. Collet, M.R. Dogar, A.D. Dragan, R.A. Knepper, T. Niemueller, K. Strabala, and M. Vande Weghe. HERB 2.0: Lessons learned from developing a mobile manipulator for the home. *Proceedings of the IEEE*, 100(8):1–19, 2012.
- [50] F. Stulp, E. Theodorou, J. Buchli, and S. Schaal. Learning to grasp under uncertainty. In *IEEE International Conference on Robotics and Automation*, pages 5703–5708, 2011. doi: 10.1109/ICRA.2011.5979644.
- [51] W. Townsend. The BarrettHand grasper—programmably flexible part handling and assembly. *Industrial Robot: An International Journal*, 27(3):181–188, 2000. doi: 10.1108/01439910010371597.
- [52] H. Zhang and N.N. Chen. Control of contact via tactile sensing. *IEEE Transactions on Robotics and Automation*, 16(5):482–495, 2000. doi: 10.1109/70.880799.