

Supplementary Material for “Mechanical Search: Multi-Step Retrieval of a Target Object Occluded by Clutter”

Michael Danielczuk^{*1}, Andrey Kurenkov^{*2}, Ashwin Balakrishna¹, Matthew Matl¹, David Wang¹, Roberto Martín-Martín², Animesh Garg², Silvio Savarese², Ken Goldberg¹

This document describes supplementary experiments and details for the ICRA 2019 submission “Mechanical Search: Multi-Step Retrieval of a Target Object Occluded by Clutter.”

I. EXTENDED RESULTS

Tables I, II, III, and IV give a detailed breakdown of each policies selected actions and success rate over the 1000 simulated trials and 50 trials on the physical robot for each policy on 15 object heaps. In simulation, pushing actions result in higher success rates. On the physical system, the human policy selects pushing actions much more frequently to clear multiple occluding objects from the target object. We will explore this discrepancy further in future work.

| <i>Simulation Policy</i> | <i>Success Rate</i> | <i>Mean Actions</i> |
|----------------------------|---------------------|---------------------|
| Random | 88.8% | 11.26 ± 0.15 |
| Preempted Random | 89.7% | 8.55 ± 0.16 |
| Preempted Random + Pushing | 94.3% | 8.87 ± 0.15 |
| Largest-First | 90.3% | 6.35 ± 0.14 |
| Largest-First + Pushing | 93.3% | 6.51 ± 0.14 |

TABLE I: Success rate and mean number of actions (with standard error of the mean) for extraction for 1000 trials of each policy tested in simulation. The largest-first policies extract the target most efficiently, and pushing shows ability to increase overall success rate.

| <i>Simulation Policy</i> | <i>Suction</i> | <i>Parallel-Jaw</i> | <i>Push</i> |
|----------------------------|----------------|---------------------|-------------|
| Random | 7015 | 4467 | 0 |
| Preempted Random | 6168 | 2870 | 0 |
| Preempted Random + Pushing | 6180 | 2825 | 274 |
| Largest-First | 5266 | 1718 | 0 |
| Largest-First + Pushing | 5116 | 1706 | 259 |

TABLE II: Breakdown of action selection for each policy in simulation over 1000 trials. Policies typically attempt many more suction grasps due to better accessibility in clutter, and only attempt pushes a small fraction of the time.

II. SIAMESE NETWORK IMPLEMENTATION DETAILS

The Siamese network architecture involves first passing each input 512×512 RGB image through a ResNet-50 architecture pretrained on ImageNet. During training of the Siamese network, these weights remained fixed. The featurizations of the input images are then concatenated and passed through two dense, fully connected layers: the first with 1024 neurons and ReLU activations, and the second with a single output neuron and sigmoid activation, whose output

| <i>Physical Policy</i> | <i>Success Rate</i> | <i>Mean Actions</i> |
|----------------------------|---------------------|---------------------|
| Random | 92% | 10.87 ± 0.66 |
| Preempted Random | 90% | 6.71 ± 0.61 |
| Preempted Random + Pushing | 98% | 6.31 ± 0.63 |
| Largest-First | 94% | 4.85 ± 0.51 |
| Largest-First + Pushing | 96% | 6.00 ± 0.63 |
| Human | 98% | 3.06 ± 0.32 |

TABLE III: Success rate and mean number of actions (with standard error of the mean) for extraction for 50 trials of each policy tested on the physical system. All policies achieve success rates of over 90% due to effective low-level grasping policies, but the human outperforms the best policy by 37% in terms of mean number of actions, suggesting that there is considerable room for action selection policy improvement.

| <i>Physical Policy</i> | <i>Suction</i> | <i>Parallel-Jaw</i> | <i>Push</i> |
|----------------------------|----------------|---------------------|-------------|
| Random | 275 | 300 | 0 |
| Preempted Random | 282 | 76 | 0 |
| Preempted Random + Pushing | 250 | 58 | 19 |
| Largest-First | 222 | 47 | 0 |
| Largest-First + Pushing | 244 | 59 | 18 |
| Human | 99 | 27 | 44 |

TABLE IV: Breakdown of action selection for each policy in 50 physical trials. The human pushes much more frequently than the other policies, especially to clear multiple occluding objects at the beginning of the trial.

represents the probability that the two input images are of the same object. The motivation for this architecture is to allow the Siamese network to learn a distance metric over the ResNet-50 featurizations. The training dataset for the Siamese network consists of 5 views of each of the objects used in physical experiments. For each view, we generated a total of 10 additional images: 5 randomly rotated versions of the original image as well as 5 rotated versions that are partially occluded. To simulate occlusions, we took randomly rotated and scaled binary masks of a dataset of synthetic objects, and overlay the masks on the original object. We only used occlusions that covered at least 20% and at most 80% of the original pixels of the object. For training, we sampled 10,000 positive and 10,000 negative image pairs, where a positive image pair consists of an original image of an object and one of the 10 augmented images and a negative image pair consists of an original image of an object and one of the 10 augmented images of an entirely different object. The network is then trained with a contrastive loss function for 10 epochs using a batch size of 64 and the Adam optimizer with a learning rate of 0.0001. For physical experiments, a recognition confidence threshold $t_r = 0.9$ was used.

III. SIMULATED HEAP GENERATION

Simulated heaps are generated by sampling: 1) N objects from a dataset of over 1600 3D models, 2) a heap center

* Authors have contributed equally and names are in alphabetical order.
¹University of California, Berkeley ²Stanford University

around the center of the bin, and 3) planar pose offsets for each object around the heap center. Then, using the Bullet Physics Engine, sampled objects are dropped one by one into the bin from a fixed height at their pose offset from the heap center, and all objects are allowed to come to rest (i.e. all velocities of all objects go to zero). Once all objects have been added to the heap, the modal and amodal segmentation masks for each object are rendered from the camera’s perspective. The modal segmask of an object is a segmask of the portion of the object visible from the perspective of the camera (accounting for occlusions), while the amodal segmask of an object is a segmask of the object’s exact position in the scene given ground truth information from the simulation environment. Using these masks, the target object is chosen to be the object with the smallest ratio between modal and amodal segmask area (i.e., the least visible object in the bin). This metric is used as a proxy for finding the most occluded object.

IV. POLICY PARAMETERS

Simulation Policy Parameters: Grasp confidence thresholds of $t_{\text{thresh}} = 0.15$ and $t_{\text{high}} = 0.3$ are used in simulation

for the high-level action selector to determine whether to execute grasp actions from the low level grasp policies. In experimental trials, these values were found to provide a balance between avoiding grasp failures and quickly clearing objects from the bin as soon as sufficiently good grasps become available.

Physical Policy Parameters: In physical experiments, grasp confidence thresholds of $t_{\text{thresh}} = 0.1$ and $t_{\text{high}} = 0.3$ were used for action selection to determine whether to execute grasp action plans from the low-level grasp action policies. These values are similar to those used in simulation, but t_{thresh} is made slightly lower for physical experiments since it we observed that low confidence grasps succeeded more often in physical experiments than in simulation, which was designed to be conservative to encourage good transfer to reality. Additionally, $t_{\text{high}} = 0.5$ was used for policies that included low-level pushing action policies, so that pushing would be further encouraged over low-quality grasp actions.