

Team Harp: Amazon Picking Challenge Stowage Task

16-662 Robot Autonomy, Carnegie Mellon University

Abhishek Bhatia, Alex Brinkman, Feroze Naina, Ihsane Debbache, and Rick Shanor

Project Video Link: <https://youtu.be/SN3n3u0Wb2s>

Introduction

The stowage task for the Amazon Picking Challenge requires autonomous decision making to control arm motions, end effector actuation, perception, and classification. The goal is to pick 12 items out of an unstructured bin and place them onto a 12 bin shelf and return the contents of every shelf bin. To accomplish this task, a UR5 6 DOF robotic arm and suction end effector system is used to manipulate items. A head mounted Kinect 2 for Windows RGB-D sensor is used to identify grasp points out of the unstructured bin and a base mounted Kinect used for item identification after grasping.

System Setup – Hardware

The Arm is a UR5 manipulator from Universal Robots and comes with its only controller and power supply. The Base and end effector are custom designed and fabricated. The grasping system consists of a suction cup and vacuum equipped with a pressure sensor that is controlled by a desktop. The Kiva pot shelf and red stowage bin are provided by Amazon and will be the same used in the competition. There is a computer-controlled actuator underneath the stowage bin to tilt the stowage bin as needed.

System Setup – Software

The UR5 comes equipped with a touch screen interface for typical manufacturing assembly. We connect to the controller over LAN socket to router and control the robot from a networked computer. ROS is used to manage the sensors, world model, controllers, and task-level executives. The TF package is in charge of managing the robot state, collision objects, and target item grasp points, and drop off set points. Caffe is used to train and get predictions from the Convolution Neural Network used for identification. PCL and OpenCV are used for the various pointcloud segmentation operations and image filtering operations as needed. Once the competition is completed, our code will be available on Github for general use.

Arm Planning

Arm planning is implemented through ROS and Moveit! using Open Motion Planning Library planners. A URDF of the robotic arm is provided by the manufacturer and modified to include our custom end effector. The default planning group starts at the base mounting point and extends through all 6 joints to the tool attachment mount. Our planning group is extended to include the end effector to enable target poses defined in the world coordinates for suction cup locations. Collision models for the mounting frame, order bin, and kiva pod shelf are loaded in using the MoveIt! planning scene interface. Plan details are specified using Moveit! like maximum planning time and which OMPL planner to use.

Each single query plan took approximately 5 seconds to find acceptable paths. Given the competition time constraint of 15 minutes, planning between set points would take too long and inhibit the effectiveness of the robot. A feature to precompute the path plans was developed and implemented in our system. The trajectory replay feature learns paths when learning mode is enabled. Whenever the trajectory-playback execution order is specified, the feature looks through the learned database for similar start configurations and goal poses and execute the stored trajectory. Otherwise, a Cartesian or single-query plan must be computed.

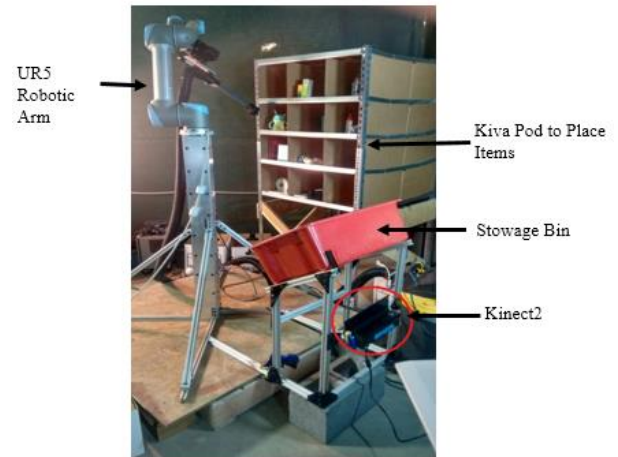


Figure 1: System Depiction

Our implementation of arm control implements a ROS blocking server that can accept planning details and execution orders. The motivation was to abstract planning as another step in our state controller so that infinite planning loops are not possible. Multiple plans and planning details can be tuned for the needs of each motion plan request. The available execution orders are trajectory-playback, Cartesian plan, strict Cartesian plan, fast single-query, and slow single-query. Trajectory-playback is a feature that replays precomputed trajectories. The Cartesian plan is a fast planner that uses FastIK to quickly solve the inverse kinematics of the arm at short intervals along the path. Straight-line paths are computed from start and goal poses. Additional waypoints can be specified and the Cartesian planner will attempt to find collision-free, straight-line plans between each waypoint. The Strict version of the Cartesian execution order forces the arm to be able to travel to the goal state whereas the normal Cartesian execution order will result and execute a partial path. Finally, the single query execution order uses the OMPL planner to perform a traditional motion planning query. The only difference between fast and slow planners is the allowable planning time. The OMPL planner could be specified for each request. RRT* was our default planner for its ability to find valid plans and improve them as time allows. RRT plans were fast but caused large swinging motions that are not desired for our picking application.

Grasp Point Selection

To determine grasp surfaces and valid grasp points, we first determine clusters corresponding to a set of points that are part of the same smooth surface. This is done using the region growing segmentation and clustering algorithm. We used the base code provided as part of the ‘Region Growing Segmentation’ pointcloud tutorial [1] to generate valid clusters. Once the clusters are identified, the clusters are scored and sorted to find the best grasp point. The score is computed for each clusters based on the maximum number of points, height of each cluster, area of the horizontal surface of each cluster (x-y axis), and direction of normal. The cluster with the maximum score is selected as the best cluster. The grasp point is then determined as the centroid of the best cluster, with the height of this point updated as the maximum height of the cluster.

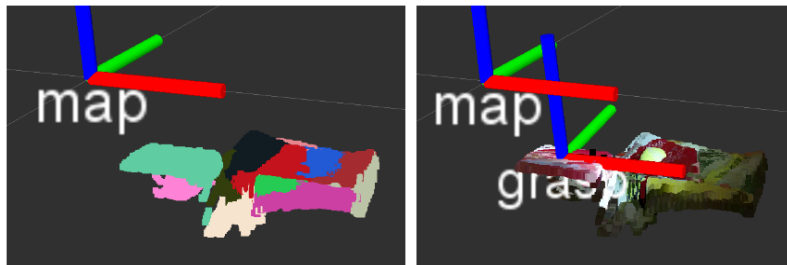


Figure 2: RGB Segmentation and Grasp Selection

Grasping Execution

Once grasp points are identified, the arm is moved to a setpoint to prepare for the pickup. The suction system is turned on and the variation in pressure is recorded. The grasp detection pressure is set as the measured average pressure minus 2.5 times the standard deviation. Using this dynamic setting ensures we are robust to small changes in pressure, hose configuration, height above sea level, and other factors that may influence unrestricted vacuum pressure. Using the 2.5 sigma boundary ensures we are as sensitive as possible to grasp detection without risking false positives. With the suction system on, the arm is commanded over the item and iterated down onto the item until pressure feedback is detected or the goal has been reached. If the grasp was successful, the arm moves the item over the identification camera to perform item identification. Otherwise, the arm is commanded over the bin where the stowage bin is reimaged and another attempt is made. First pass picking success using this approach is approximately 90%.

Item Identification - Database Generation

Any approach for item identification requires a large amount of training data to create a classifier. To generate this data for 40 items, a turntable was created to aid in database generation. The turntable consists of an actuated turntable, capable of 360 degrees of rotation and an actuator to vary the view angle of the Kinect v2 RGB-D sensor. The turntable is shown in Figure 3.



Figure 3: Turntable for Dataset Collection

Manually masking each input image would have required a lot of time so HSV thresholding and convex hull filters were applied to the images to

automatically mask the input images. Each HSV filter had to be tuned by hand for each item. Items that contained green elements on the edges required a different background color. Specular items against the green background would reflect green hues so a cardboard background was chosen for specular and green items.

Once the HSV threshold was tuned per item, a convex hull routine was applied to the image to create cleaner filters. The convex hull operation finds the smallest polygon that contains all of the pixels that pass the HSV color filter. This ensures that holes in the masks do not appear in the final mask which could cause the identification classifier to learn fictitious features. Figure 3 shows the output of each operation

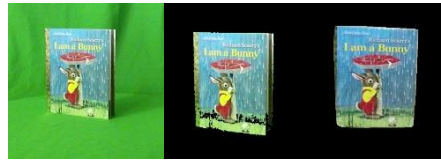


Figure 4: Original image (left), the HSV color filter image (center), and the final mask after Convex Hull (right)

In total, approximately 100 masked images were captured for each of the 40 items in the item dictionary. These images were rotated, distorted, mirrored, lightened, and darkened to create approximately 400,000 images for future classifier training.

Item Identification – Online Segmentation

Online segmentation is required to produce the mask similar to the images created for the database. Since the robotic system geometry is known, a cubic region of interest is created and pointcloud points that fall outside this region of interest are removed from the image. Care was taken to ensure the region of interest was large enough to capture full images of the item without including bad pixel values Figure 4 shows the raw and masked images.



Figure 5: Raw (left) and Masked (right)

Item Identification - SLIC Superpixels

Superpixel generation is a color based segmentation technique. A superpixel is defined as a group of pixels with similar characteristics. To solve our problem, we are using the SLIC (Simple Linear Iterative Clustering) superpixel algorithm to generate superpixels from an input image, primarily because of its low computational overhead [2]. Considering the scope of this project, we did not modify much within the superpixel generation pipeline and directly used the off the shelf package provided by scikit-image ‘skimage.segmentation.slic’ that takes in a RGB image (3D array) and outputs an integer mask indicating segment masks [4]. Once we generate the valid superpixels for each input (masked) image, we input these superpixels to the CNN described in sections below for training and testing/identification purposes. Figure 6 shows the output of the superpixel operation.

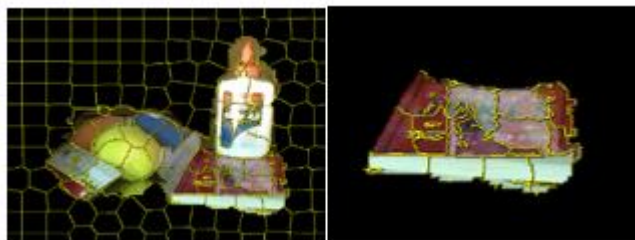


Figure 6: SLIC Superpixels

Item Identification - CNN for Superpixel Identification

To capture the large amount of appearance variation from item to item, a convolution neural network (CNN) was chosen to be the classifier. The popular deep net AlexNet was selected since its architecture and training made it well suited to the Amazon picking application. The CNN takes in a 255x255 RGB pixel array and passes the image through a series of five convolutional layers that filter and preprocess the image for distinctive features. Next, three fully-connected layers are trained to detect the identifying features for each item.

Our implementation of AlexNet requires an additional output layer to map the 1000 class predictions to the 40 class predictions needed for our application. This final layer consists of an SVM classifier with no initial training. The CNN takes in a 255x255 RGB image of each superpixel and outputs the probability that the superpixel belongs to each class. This is repeated for each superpixel in the segmented image to create an Nx40 prediction matrix for N superpixels.

The CNN was trained on images collected by our turntable apparatus and segmented into superpixels. The initial parameters were used for the net and our 400,000 training images were used to adjust the parameter and output SVM later using a learning rate of .001 and batch size of 250 training images. The model was loaded into an Nvidia Titan GPU for training and requires approximately 6 hours of computation to retrain the model.

Item Identification - Local Item Prediction

The goal of local item prediction is to identify the grasped item out of the 12 possible items in the request list to guide the robot executive decision-making. The probability of correct identification is low for a 12 choose 1 decision, but still useful to the operation of the system.

The output of the CNN produces an Nx40 matrix of predictions for N superpixels and the 40 items in the item dictionary. From this output, the items not contained in the 12 item request list can be ignored and the probabilities for the remaining 12 items renormalized. Finally, the average prediction is taken over all the superpixels for each item. This results in a 1x12 matrix of confidence for each of the 12 possible items. The local prediction is made as the item with the greatest confidence.

Item Identification - Global Item Prediction

The goal of global item prediction is to identify the entire set of items out of the 12 possible items in the request list. During autonomous operation, the robot will grasp and pickup each of the items and perform local prediction on each. After all images are captured, segmented, and passed through the CNN, global prediction compares the full set of prediction probabilities to make the best overall prediction.

The global item prediction begins with the Nx40 matrix of predictions for each of the 12 items. As with local prediction, impossible items are removed from the superpixel matrix. As before, each matrix is averaged across superpixels and renormalized to get the average item confidence.

Next, the global prediction finds the most confident item in all of the image confidence tables and assigns that image to the item label. The predictions for the selected item label is removed from the remaining images and each image confidence table is renormalized. Now the process repeats. With the set of remaining items, the highest confidence prediction is taken and assigned to the correct image. The predictions for the best item keep getting removed until all items have been assigned to images.

Item Identification – Results

The item identification process was tested by acquiring several images for each item, serving as ground truth, and choosing 12 images at random to do item identification. The segmentation and CNN operations were applied to the images and then global prediction was performed on the output prediction tables. The study was performed on 10,000 random permutations of images and compared to the ground truth. Overall, 75% accuracy was achieved by this method. A

confusion matrix was created comparing the global prediction results against the ground truth. Figure 6 shows several easily confused items from the confusion matrix.



Figure 7: Examples of Several Confusion Sets

Place Operation

The place operation occurs once the item has been grasped and imaged for identification. The arm is commanded to one of the 12 shelf bins and then translates into the center of the bin. The suction system is disengaged and the item deposited in the bin. The arm withdraws and moves to look down into the stowage bin, ready to pick up the next item.

Results

The current system performs first-pass picking with 90% accuracy. Local prediction is approximately 58% accurate and global item prediction is 75% accurate. Overall, we can pick and place approximately 84% of items and according to the competition rules score full points on approximately 63% of attempts. These results are promising but the system will require refinements before the competition in the end of June.

Future Work

Plans to improve the performance of the robot are focused on the identification and stowage operations. Identification will be improved by revisiting the segmentation to reduce data loss in the pointcloud that could be reducing the confidence of the predictions. The CNN will be trained on more data as we acquire more testing data to better reflect the images it will be actually see at runtime. Finally, multiple viewpoints may be acquired to ensure there is at least one good image of the item.

The other main area of improvement is the stowage operation. Currently, there is no pose estimation of the grasped item and the placement on the shelf goes to the center of the shelf. This works well for most of the items but larger items will need to be handled better. We want to implement a simple pose estimate by approximating a bounding box from the identification pointcloud. From that, a simple offset may be added to the shelf placement for improved placing.

References

- [1] PointCloud Library http://pointclouds.org/documentation/tutorials/region_growing_segmentation.php
- [2] Achanta, Radhakrishna, et al. Slic superpixels. No. EPFL-REPORT-149300. 2010.
- [3] Wikipedia. 2016 https://en.wikipedia.org/wiki/Lab_color_space
- [4] Scikit-Image. 2016 <http://scikit-image.org/docs/dev/api/skimimage.segmentation.html#skimimage.segmentation.slic>
- [5] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." Advances in neural information processing systems. 2012.
- [6] The Quadratic-Chi Histogram Distance Family . Werman Ofirpele. <http://www.ariel.ac.il/sites/ofirpele/publications/ECCV2010.pdf>

