

Contents

| | | |
|----------|---|----------|
| 1 | Introduction / <i>Malthe Høj-Sunesen</i> | 2 |
| 2 | Grasping | 3 |
| 2.1 | Motivation / <i>Malthe Høj-Sunesen</i> | 3 |
| 2.2 | Simplifying objects to primitive models / <i>Malthe Høj-Sunesen</i> . | 3 |
| 2.3 | Learning to grasp through attempts / <i>Malthe Høj-Sunesen</i> . . . | 4 |
| 2.4 | Elementary grasping actions / <i>Hsin-Yu Lee</i> | 6 |
| 3 | The vision system / <i>Hsin-Yu Lee</i> | 8 |
| 3.1 | ECV system / <i>Hsin-Yu Lee</i> | 8 |

1 Introduction *Malthe Høj-Sunesen*

According to ISO 8373 [ISO, 2012], at least two different types robots exist: Industrial robots and service robots. An industrial robot is defined as aa “automatically controlled, reprogrammable, multipurpose manipulator programmable in three or more axes”, while a service robot is defined as a “robot that performs useful tasks for humans or equipment excluding industrial automation applications”. The classical application of an industrial robot is to have the robot do a predefined behavior repeatedly, while service robots are still very much under development. Due to hardware and software concerns, robots in the industry have previously not seen adaptive behavior, so elements must be aligned in a specific way. Humans, and indeed most animals, are able to look at objects and grasp accordingly. A lot of research is going into making the robot able to understand what it is “looking” at much like humans can, and how to grasp it. This research into grasping¹ objects using only visual cues is the focus point for this report.

¹For the purposes of this report, *grasping* is to pick up an object.

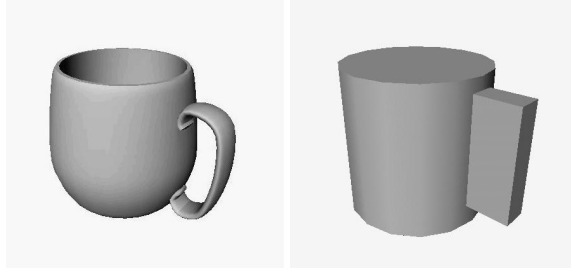


Figure 1: “A mug and its primitive representation”; a cylinder and a box. From [Miller et al., 2003].

2 Grasping

2.1 Motivation *Malthe Høj-Sunesen*

Humans spend years learning how to grasp objects. Babies have a hard time figuring out how to grasp even the most simple objects, and parents solve that problem by giving babies and children plastic cutlery, bouncy, soft toys and always walking around with an eye on each finger. We come to expect of a child to drop toys, knock over glasses, and the like.

A robot is not allowed to fail in the same way. When a robot’s hand grasps something we expect it to not let it go — or worse, drop it — before it is supposed to. In a tightly controlled production line that is not a problem. Using embodied AI the parts can be aligned perfectly for the robot and the robot can assemble the parts correctly.

In a not so tightly controlled environment among people it is a bigger problem. If a service robot is supposed to clean up mess left after a human, it is almost guaranteed that the parts are not aligned as a robot could predict. If an industrial robot can figure out the best grasp autonomously for an object it would decrease operator dependency, leading to faster setup and lower costs for the company.

2.2 Simplifying objects to primitive models *Malthe Høj-Sunesen*

Biederman suggested that elements can be broken down to geons, basic elements describing one feature of an object.

In [Miller et al., 2003] the idea behind geons is used to help a robot simulator find good grasps. The robot knows how to grasp each shape primitive (equivalent to geon). Any object is then reduced to its shape primitives where applicable. This allows a simulator to know which points are good to grasp, resulting in simpler calculations. An example of this reduction can be seen in Figure 1 with shape primitive building bricks shown in Figure 2.

Reducing the visual information in this way will give the simulator a



Figure 2: “Examples for grasp generation on single primitives. The balls represent starting positions for the center of the palm. A long arrow shows the grasp approach direction, and a short arrow shows the thumb direction. In most grasp locations, two or more grasp possibilities are shown, each with a different thumb direction.” From [Miller et al., 2003].

simpler task, as it does not have to simulate thousands of possible grasps but only the grasps based on the preshape grasps per primitive representation. An example of the found grasps can be seen in Figure 3.

2.3 Learning to grasp through attempts *Malthe Høj-Sunesen*

Much like [Miller et al., 2003] in Section 2.2 tried to emulate how the human vision works according to Biederman, so do [Detry et al., 2011] try to emulate how a child learns to grasp objects. Any parent will tell you that their child did not quite know how to actively grasp² toys from the beginning. Where to hold is one of the problems.

The approach in [Detry et al., 2011] is to let a robotic platform learn how to grasp a single object. Using stereo vision, the 3D features of an object can be calculated. The system will then try to calculate where that object can be picked up. An example of where the system calculates a toy pan can be grasped can be seen in Figure 4 on the following page. s

After calculating where the object can be grasped, the robot will start to grasp the object, time and time again. In [Detry et al., 2011], the robot performed more than 2000 grasps. During the grasp trials the robot and

²Let alone letting it go again!

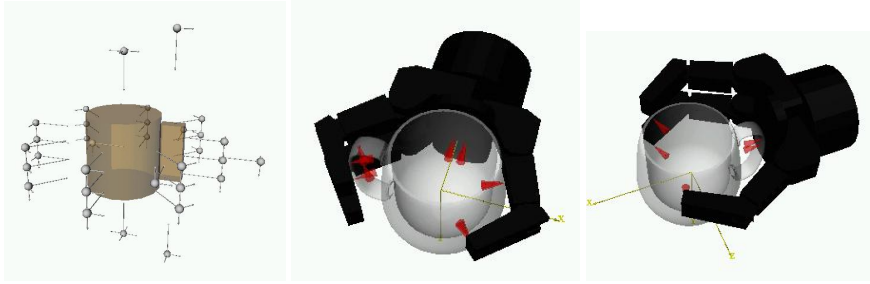


Figure 3: The primitive mug representation and the two best grasps. The red cones indicate point-of-contact. From [Miller et al., 2003].

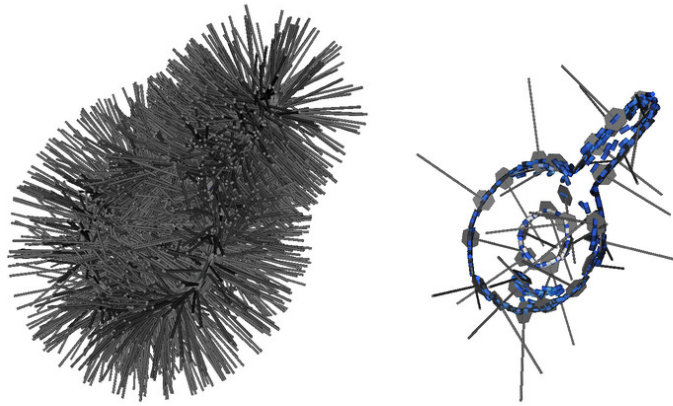


Figure 4: Left: A full graph of the possible grasp positions. Right: Clearer showing of what the sticks mean; the stick is the robotic hand's translation while the paddle at the end is the point where the fingers close. From [Detry et al., 2011]

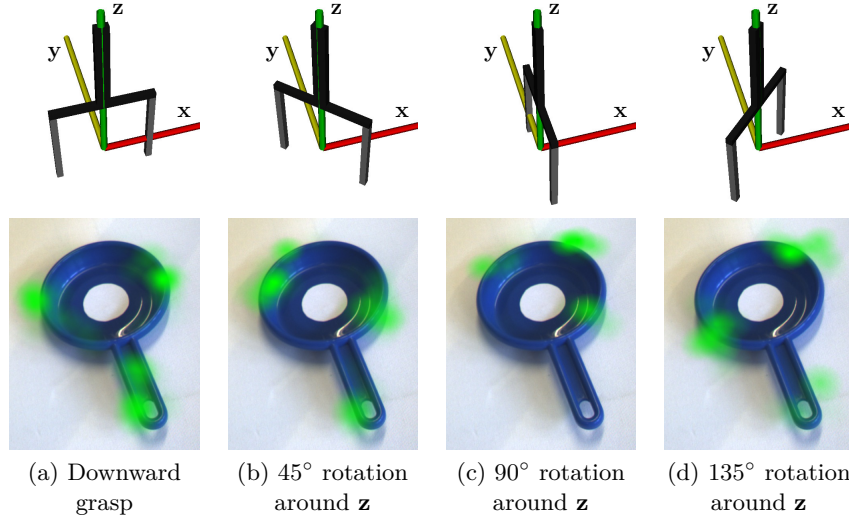


Figure 5: “Various projections of a grasp density generated by a robot”. The greener the pixel, the higher the density and thus probability of a good grasp. From [Detry et al., 2011].

vision system will see if the grasp is stable, ie. if the object is not moving. Using all the trials the robot system can build a model of the possibility that a grasp will be successful. The result is visualized in Figure 5.

In the end, the robot has learned how to pick up an object from any angle, and is able to choose the best possible grasp in any situation. This is indeed how children and grown-ups know how to pick up objects as well.

2.4 Elementary grasping actions *Hsin-Yu Lee*

The Elementary Grasping Actions, which also called EGA, is the specific grasping gestures used in the paper [Kootstra et al., 2012] that was implemented once after we find the object by the ECV system. The EGA was also mentioned and defined in the paper [Pugeault et al., 2010]. The Figure 6 on the following page shows the three different ways of grasping object base on edge and surface information separately. Basically, the directions, approaching ways and positions are pre-defined, only slightly changed according to the size of the surface or the distance between two selected contours.

We can simply choose one method to grasp something; we also can use the simulator to help us select the best grasping action through these methods. There are several experiments presented in the paper [Kootstra et al., 2012] shows the performance of each different actions and in the circumstance of using multiple method at the same time.

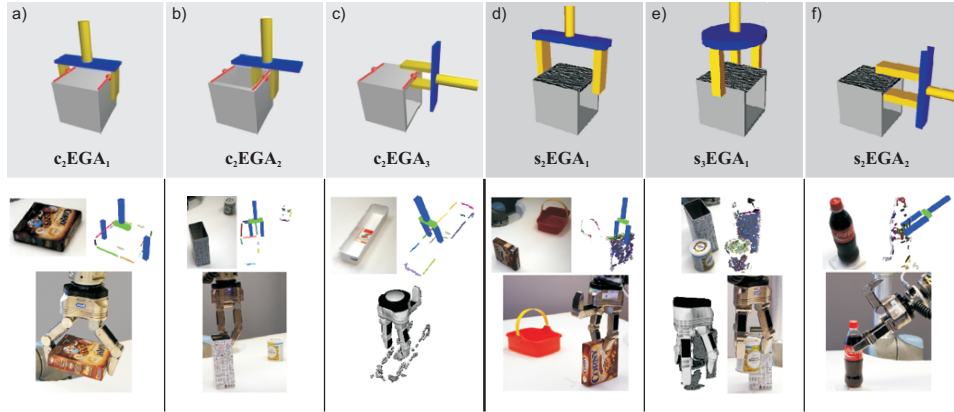


Figure 6: The elementary grasping actions (EGA) are illustrated on the top with some examples of real grasps on the bottom. (a-c) The three contour-based EGAs. The red lines indicate the selected contours. (d-f) The surface-based EGAs. The dark faces show the selected surface. The first letter in the naming scheme marks the type of features used to generate a grasp. 'c' stands for contour and 's' for surface. The first subscript stands for two or three fingers. The last subscript marks the general type of grasp, where '1' is an encompassing grasp, '2' is a pinch grasp from the top and '3' is a pinch grasp from the side of the surface. For each type of grasping action, an example is shown consisting of an original image, a snapshot of the ECV representation along with the selected grasp, and the grasp execution in the simulation/real setup. From [Kootstra et al., 2012]

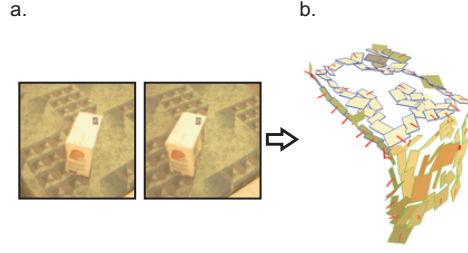


Figure 7: a. The images captured from the stereo cameras. b. The analyzing result of implementing the ECV on the stereo images. From [Kootstra et al., 2012].

3 The vision system *Hsin-Yu Lee*

If we don't want to use shape primitives as our way to generate grasping gestures, there is an alternative way when we want way to recognize the object just base on the stereo images from the cameras. The approach in [Kootstra et al., 2012] tries to imitate the way that human vision system try to recognize unknown things. The analyzing system is called "biological-motivated hierarchical vision system" [Pugeault et al., 2010], which also called "ECV", the "Early Cognitive Vision". Literally, the system was inspired by the primate's vision system. By using this system, the 3D features of edge and surface of the object are naturally aligned together. The basic analyzing process of the system can be seen in Figure 7.

3.1 ECV system *Hsin-Yu Lee*

Figure 8 on the following page shows the hierarchical vision system "ECV" implemented in the paper [Kootstra et al., 2012]. The system recognizes objects by 2D and 3D geometrical and appearance relations between visual entities at the different levels of the hierarchy in two major domains, which are edge and surface. The process of recognizing the edge is: First, use the algorithm of image processing to transfer the stereo images into the 2D line segment images. Second, use the mathematic way to find out the smallest line segments. Then, combine the line segments into a larger segment. Finally, the edge of the objects will come out. It's also similar to the process of forming the surfaces of the object.

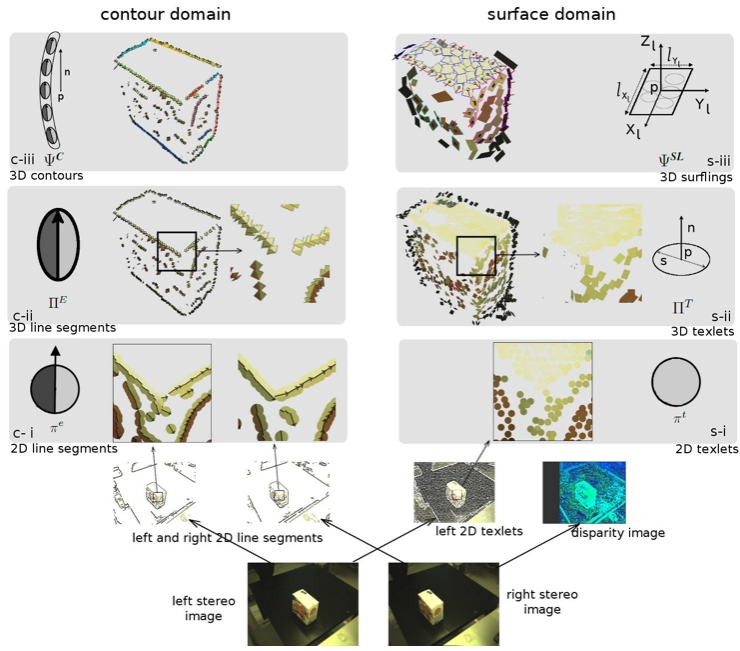


Figure 8: The hierarchical representation of contour and texture information in the ECV system. The stereo images at the bottom are the real world images captured by the camera, the others pictures show the process in the simulator try to find the object based on the edge and surface information. From [Kootstra et al., 2012].

References

- [Detry et al., 2011] Detry, R., Kraft, D., Kroemer, O., Bodenhagen, L., Peters, J., Krüger, N., and Piater, J. (2011). Learning grasp affordance densities. *Paladyn, Journal of Behavioral Robotics*, 2:1–17.
- [ISO, 2012] ISO (2012). Robots and robotic devices — vocabulary. ISO 8373:2008(en), International Organization for Standardization, Geneva, Switzerland.
- [Kootstra et al., 2012] Kootstra, G., Popović, M., Jørgensen, J. A., Kuklinski, K., Miatliuk, K., Kragic, D., and Krüger, N. (2012). Enabling grasping of unknown objects through a synergistic use of edge and surface information. *Int. J. Rob. Res.*, 31(10):1190–1213.
- [Miller et al., 2003] Miller, A., Knoop, S., Christensen, H., and Allen, P. (2003). Automatic grasp planning using shape primitives. In *Robotics and Automation, 2003. Proceedings. ICRA '03. IEEE International Conference on*, volume 2, pages 1824–1829 vol.2.
- [Pugeault et al., 2010] Pugeault, N., Wörgötter, F., and Krüger, N. (2010). Visual primitives: Local, condensed, semantically rich visual descriptors and their applications in robotics. *International Journal of Humanoid Robotics*, 07(03):379–405.