



2D-3D Object Categorization for Task-based Grasping

Marianna Madry

Dan Song

Danica Kragic



Computer Vision and Active Perception Lab, KTH, Stockholm, Sweden

{madry, dsong, danik}@csc.kth.se

Motivation

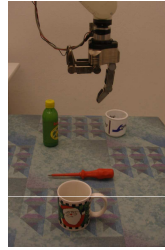
- How to grasp an object?



- Humans classify objects according to their functionality, i.e. depending on tasks they afford [Greene'94]
- Grasp knowledge can be transferred between objects that belong to the same category
- Fusion of 2D and 3D may provide more robust system

Goal

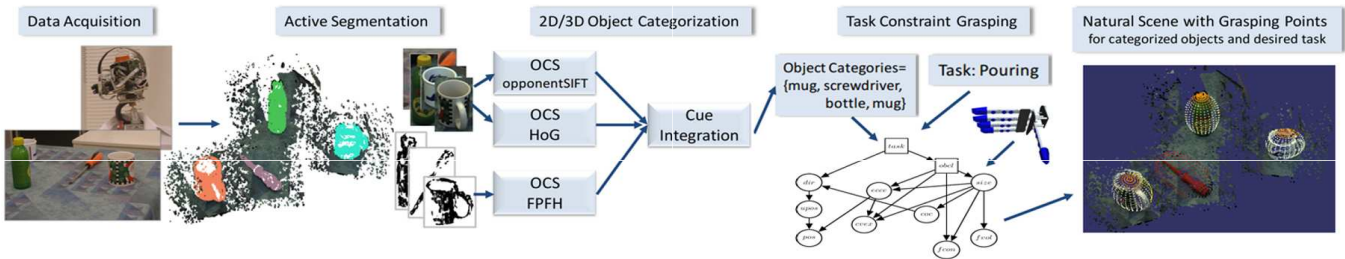
- Where to grasp to pour a liquid or hand over an object?
- Goal: Finding grasping points for a desired task in a natural scene
- Task-based grasp depends on:
 - Embodiment
 - Scene content
- Mid-goal: Finding object categories



Contributions

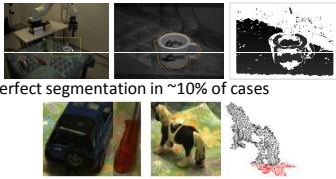
- Object Categorization:
 - Evaluation of several 2D and 3D appearance, color and shape descriptors on real stereo data
 - Fusion of 2D and 3D object categorization with high categorization rate (up to 92% for 11 object categories)
- Task-based Grasping:
 - Robot can choose objects that afford a desired task
 - Robot can plan the grasp that satisfies the constraints posed by the task
- Integration of the 2D-3D Object Categorization System with the active segmentation module [Björkman '10] and the probabilistic grasp reasoning system [Song '10]

System Overview



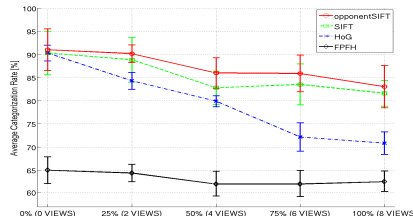
Active Segmentation

- Attention mechanisms in the peripheral view direct the foveal camera towards region of interest
- Imperfect segmentation in ~10% of cases



Object Representation

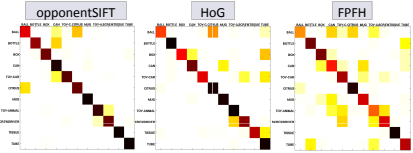
- Evaluation of several 2D and 3D descriptors encoding different object properties: appearance (SIFT), color (opponentSIFT), 2D shape (HoG), 3D shape (FPFH)
- We build a single-cue OCS for each descriptor:
 - Spatial pyramid for 2D; Bag-of-words for 3D
 - Classification: SVMs with a χ^2 kernel
- Results:
 - Performance under varying viewpoint condition



Setup-50:

Descriptor	Av. Categ. Rate	σ
SIFT	82.8%	3.6%
opponentSIFT	86.0%	3.3%
HoG	79.9%	1.2%
FPFH	62.0%	2.8%

Confusion matrices – complementarity for cue integration:

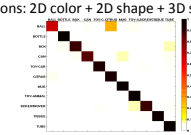


2D-3D Object Categorization

- Confidence Measure:
 - normalized distance of a sample to the hyperplane for OaA
- Cue Integration:
 - Fusion of evidences from the single-cue OCSs at the high level
 - Evaluation of the linear and nonlinear integration methods:

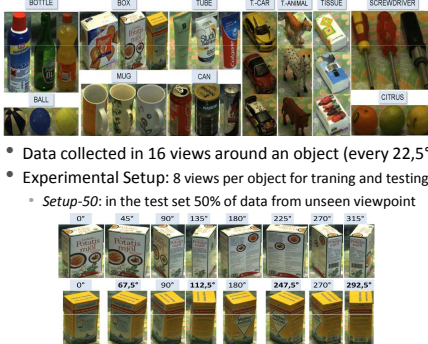
Descriptor (D1+D2)	Average		Max Rate		Product Rule		Sum Rule	
	Categ. Rate	σ	Avg. (D1,D2)	σ	Avg. (D1,D2)	σ	Avg. (D1,D2)	σ
opponentSIFT+FPFH	84.5%	3.8%	11.54, 22.4%	89.9%	4.9%	15.93, 27.9%	82.8%	4.81, 23.2%
opponentSIFT+HoG	80.6%	3.7%	10.6, 24.5%	87.8%	4.3%	1.9, 25.9%	80.9%	3.2, 4.9, 23.9%
opponentSIFT+HoG	81.6%	3.4%	4.1, 2.9%	80.0%	0.8%	91, 6.1%	87.4%	0.6%
HoG+FPFH	78.9%	3.5%	0.6, 17.9%	83.1%	6.5%	3.1, 21.6%	83.4%	4.6%

- Results:
 - Performance under varying viewpoint condition
 - Best combinations: 2D color + 2D shape + 3D shape descriptors
 - 2D-3D OCS significantly outperforms the best single-cue OCSs
 - Linear weighted summation is better than the complex methods



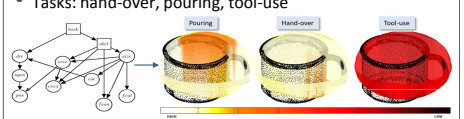
110 Object Stereo Database

- 11 object categories x 10 object instances per category
- Data: 2D (RGB image) and 3D (point cloud)
- Data collected in 16 views around an object (every 22.5°)
- Experimental Setup: 8 views per object for training and testing
- Setup-50: in the test set 50% of data from unseen viewpoint



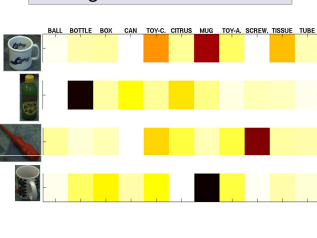
Task-constrained Bayesian Network

- Training on the synthetic 3D object models from the Princeton Shape Benchmark
- Tasks: hand-over, pouring, tool-use



Where to grasp to perform the desired task in the real scene?

Categorization Confidence



Hand-over

Pouring

Tool-use

