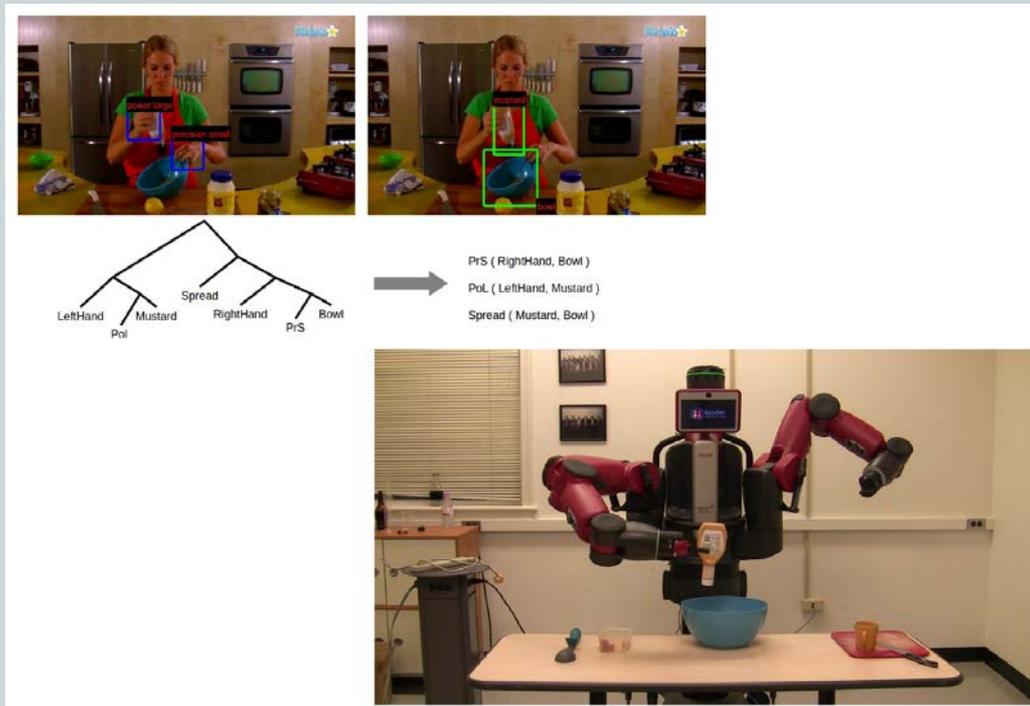


# What Happened to the Representations of Perception ?

Cornelia Fermüller  
Computer Vision Laboratory  
University of Maryland

# Why we study manipulation actions?

## 2. Learning from humans to teach robots



# The robot learns to mix a drink

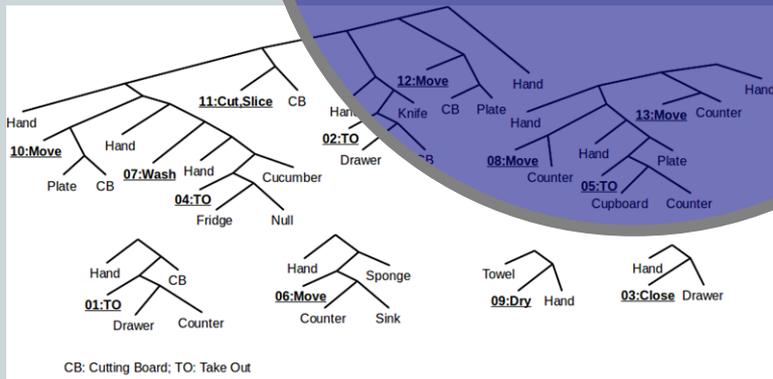
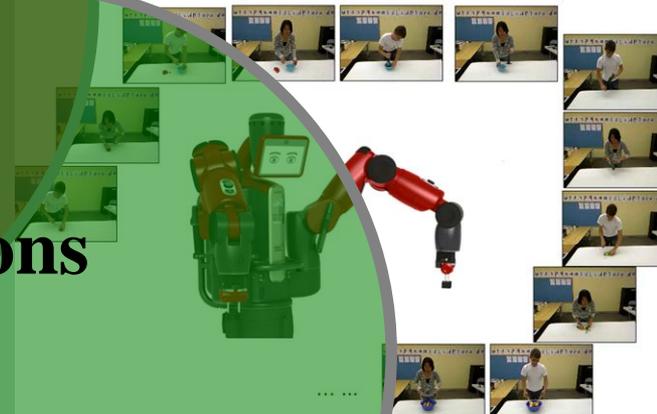
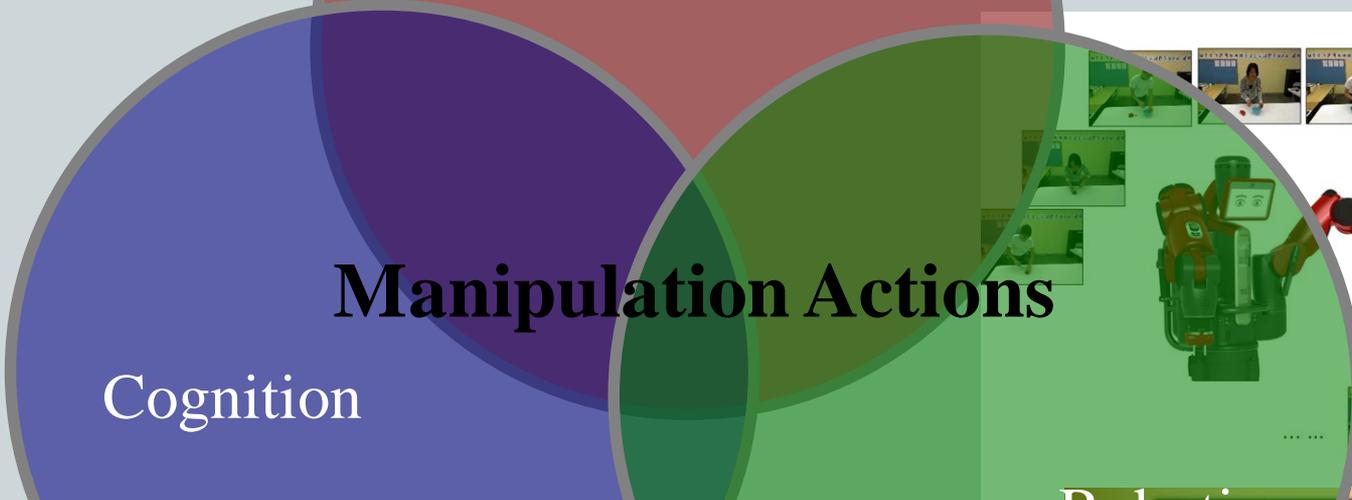
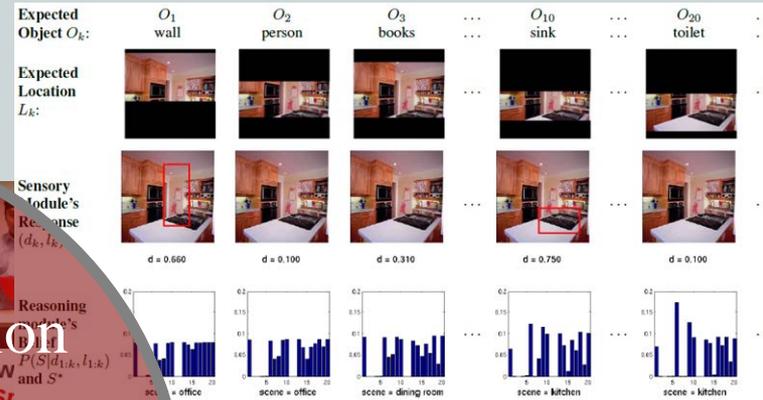




**Computer Vision**

Reasoning about the scene's  $P(S|d_k, l_k, l_1, k)$  and  $S^*$

<Tool: Knife> <Action: Cut>  
 <Tool: Peeler> <Action: Peel>  
 <Tool: Cup> <Action: Drink>  
 <Tool: Bow> <Action: Spoon>

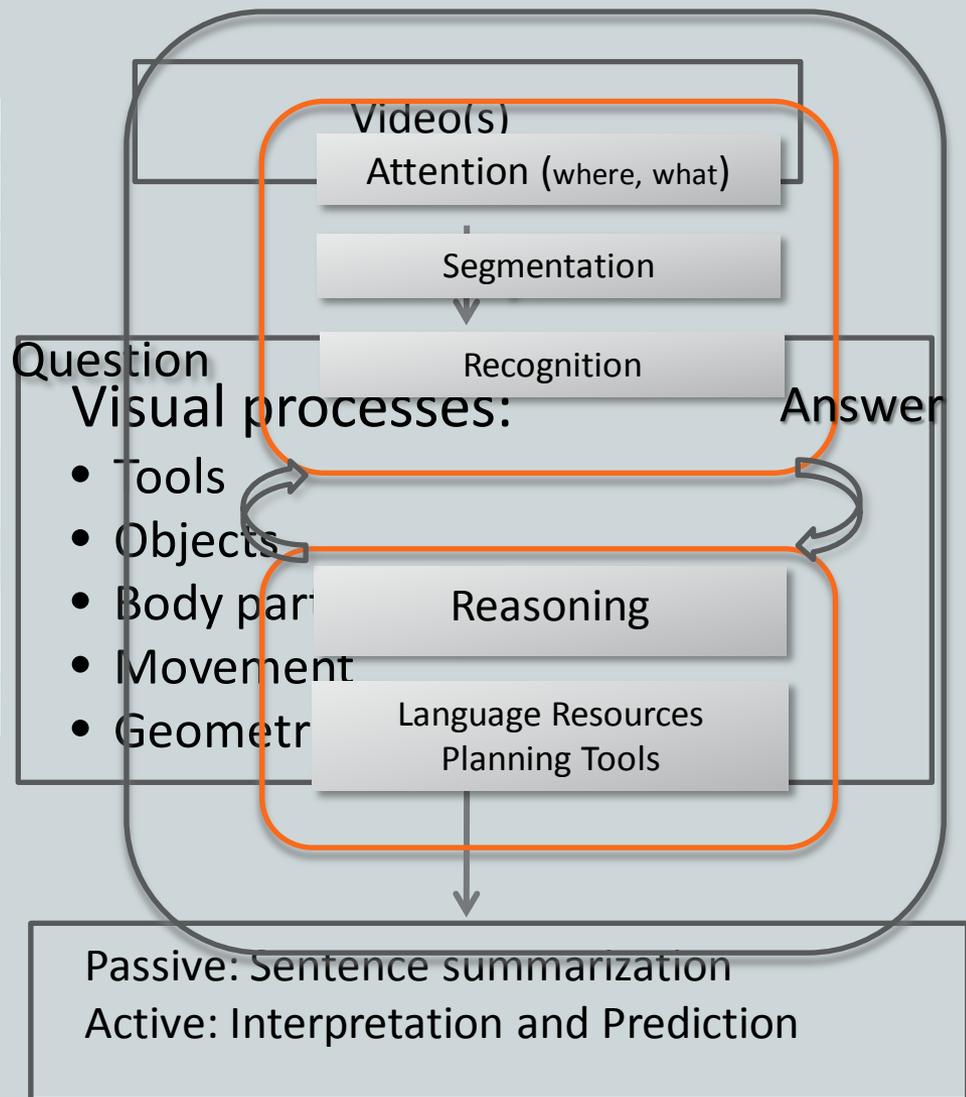


#Event	Action	Subject	Agent	From/At	To
01	Take Out	Hand	CB	Drawer	Counter
02	Take Out	Knife	CB	Drawer	CB
03	Close	Hand	Drawer	Null	Null
04	Take Out	Hand	Cucumber	Fridge	Null
05	Take Out	Hand	Plate	Cupboard	Counter
06	Move	Hand	Sponge	Sink	Counter
07	Wash	Hand	Cucumber	Null	Null
08	Move	Hand	Plate	Counter	Null
09	Dry	Towel	hand	Null	Null
10	Move	Hand	Cucumber	Plate	CB
11	Cut Slice	Knife	Cucumber	CB	Null
12	Move	Hand	Cucumber	CB	Plate
13	Move	Hand	Plate	Counter	Null

**Robotics**

"Stir the coffee..."

# The approach: A dialogue

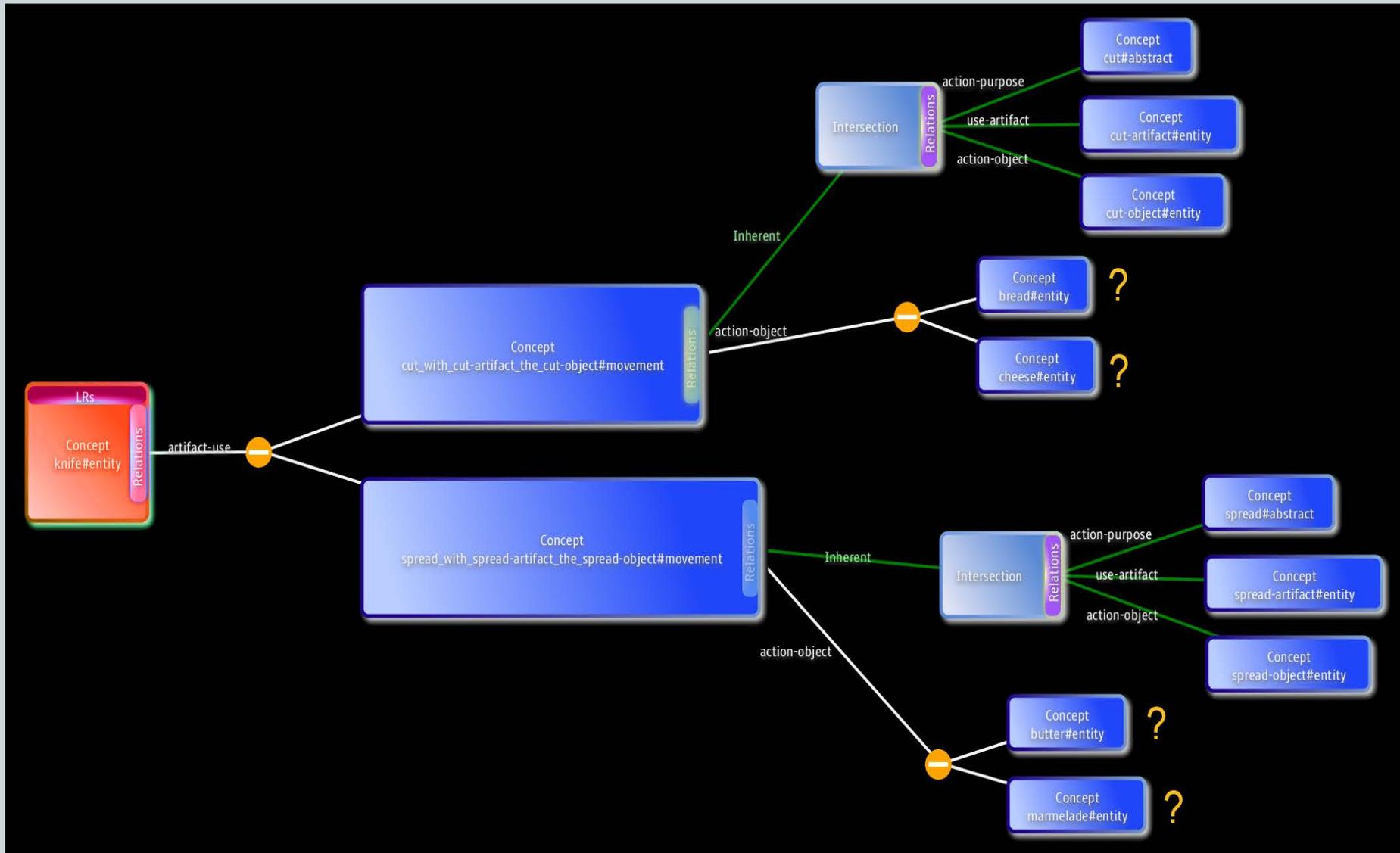


# Hand segmentation and tool detection



Hand segmentation  
Object next to hand

<Tool: Knife>

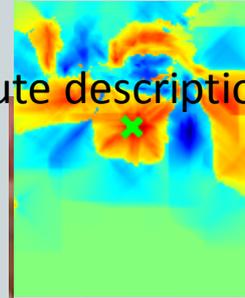


# Object recognition



Attention operator

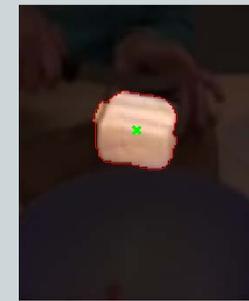
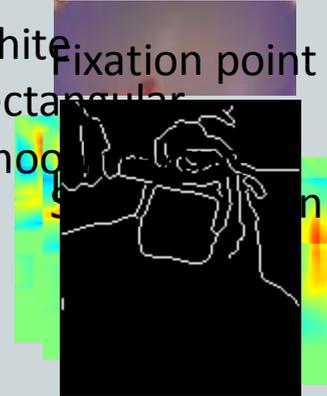
Attribute description



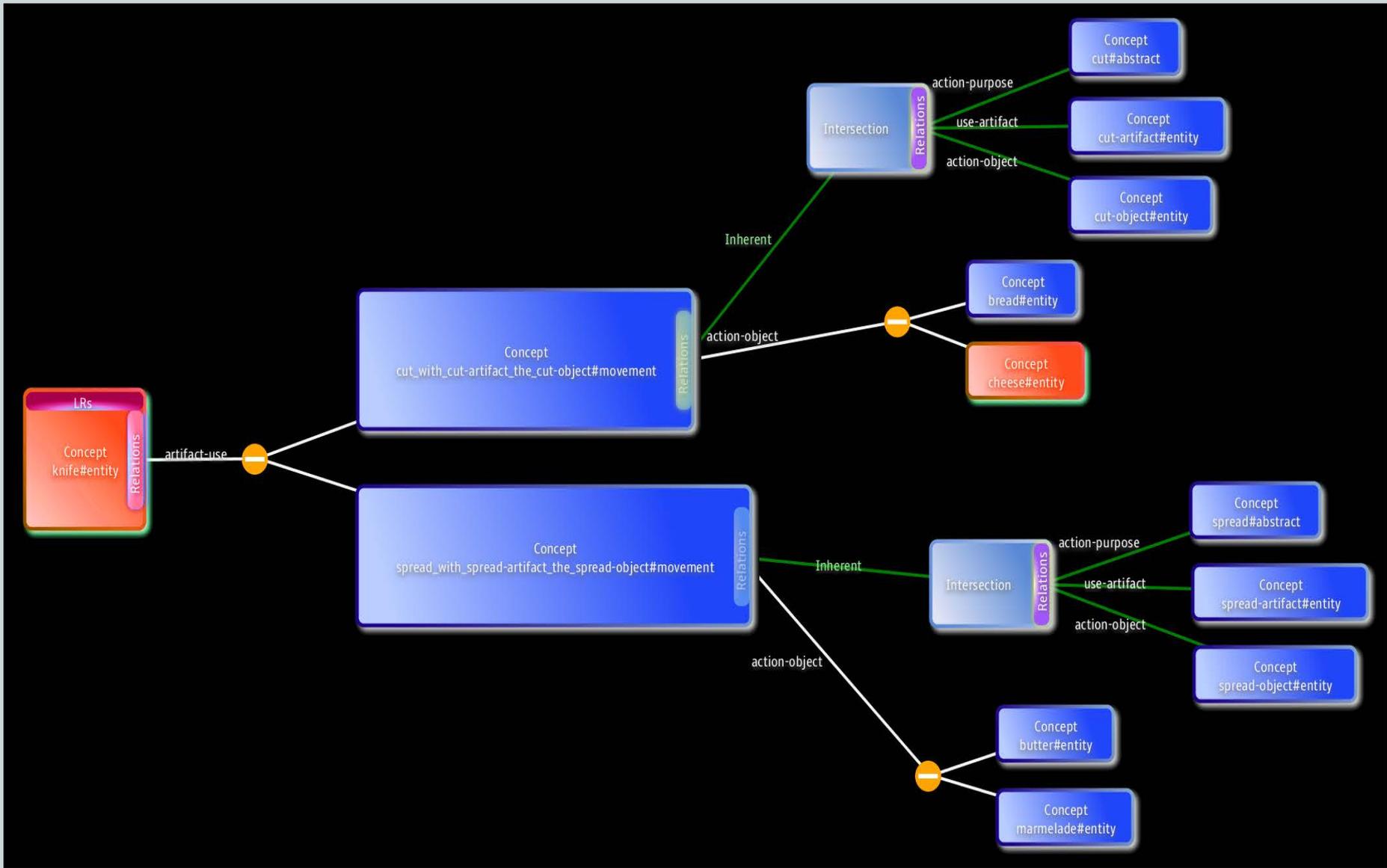
White: Fixation point

Rectangular

Smooth



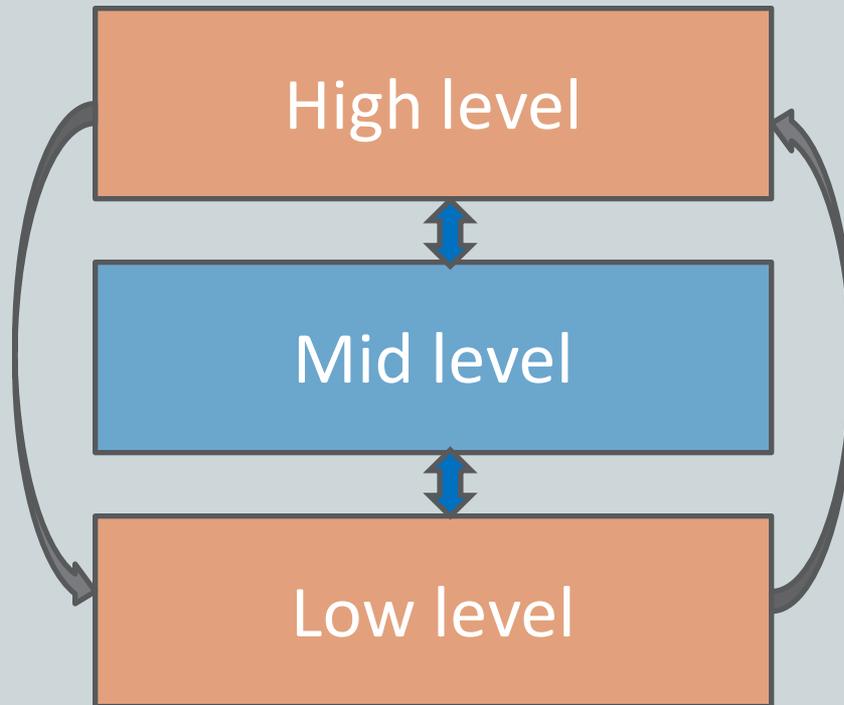
Segmentation



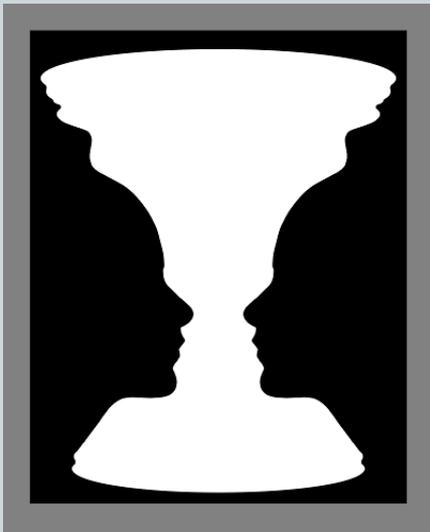
↓  
**Cheese**

# Vision Processes

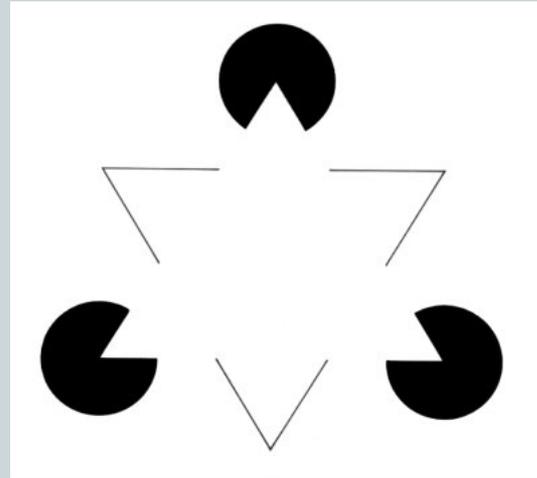
- Mid-level process (bottom-up and top-down) for object recognition



# Visual illusions demonstrating Gestalt principles



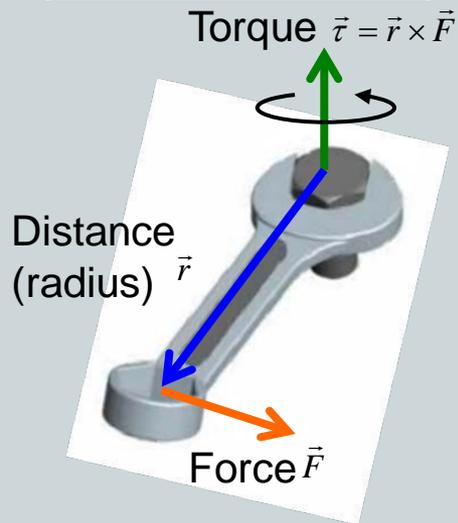
*Rubin, 1915*



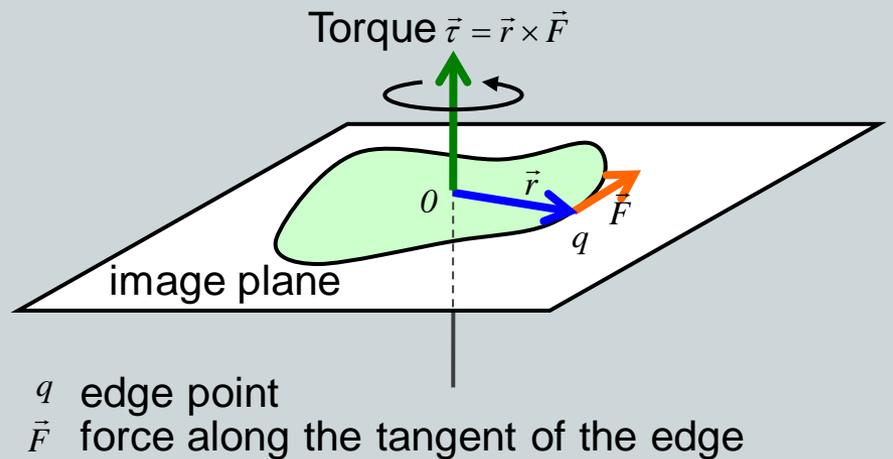
*Kanizsa, 1976*

# Torque in Images

## Torque in Physics

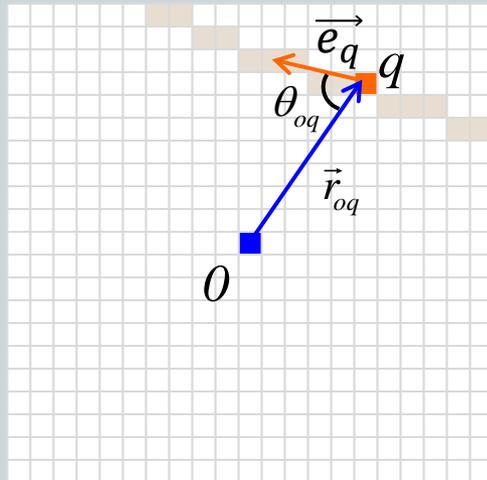


## Torque in Images



# Definition of Torque in Images

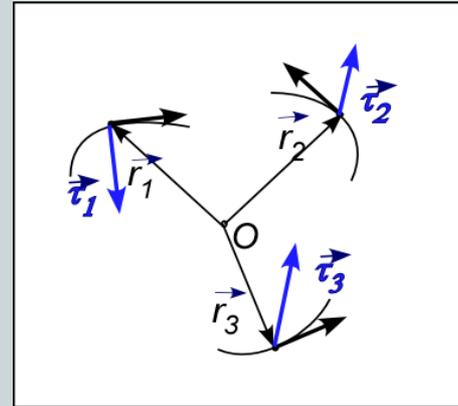
## Discrete Edge Points



Torque at point  $q$ :  $\vec{\tau}_{oq} = \vec{r}_{oq} \times \vec{e}_q$

### Value of the torque

$$\begin{aligned} \tau_{oq} &= |\vec{r}_{oq}| \cdot |\vec{e}_q| \cdot \sin \theta_{oq} \\ &= |\vec{r}_{oq}| \sin \theta_{oq} \end{aligned}$$

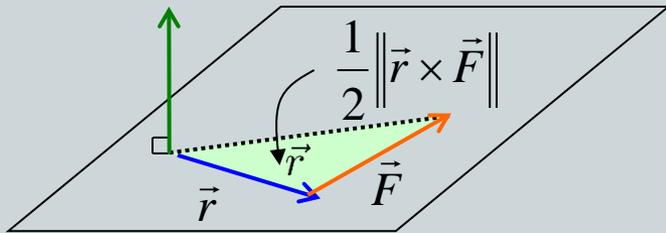


### Torque of a patch

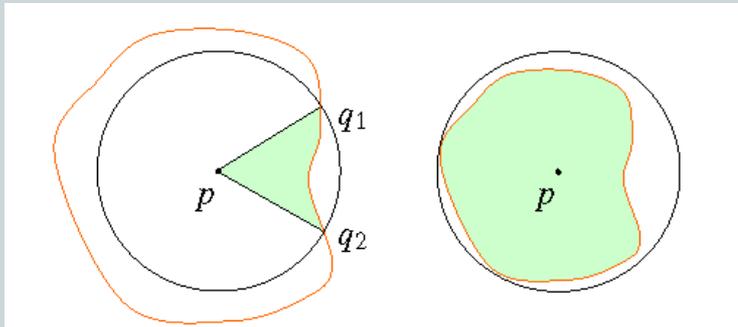
$$\tau_P = \frac{1}{2|P|} \sum_{q \in E(P)} \tau_{pq}$$

$|P|$  area of the image patch  $P$   
 $E(P)$  set of edge points within the patch  $P$

# Torque of an Image Patch



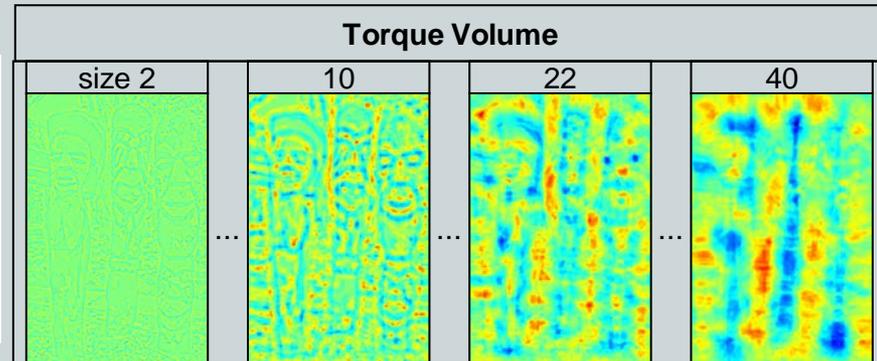
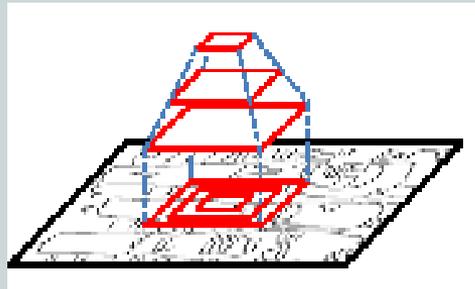
The triangle enclosed by vector  $\vec{r}$  and  $\vec{F}$  is equivalent to  $\|\vec{r} \times \vec{F}\|/2$



Torque of image patch is related to the area enclosed by a contour

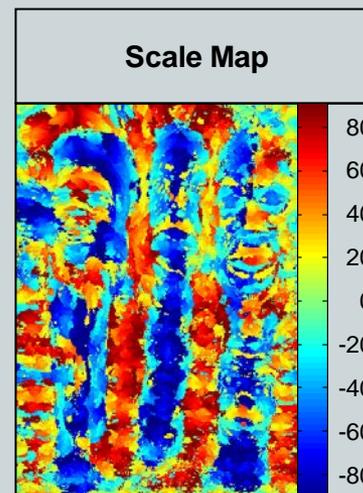
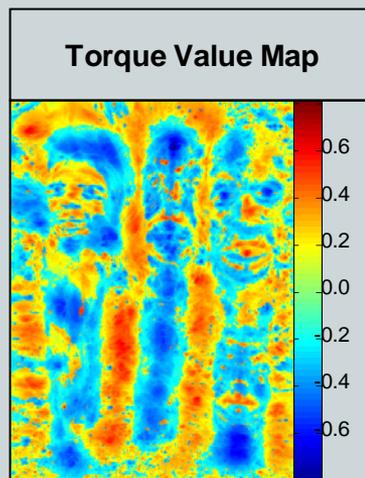
\*Disk or rectangle patches are used in our experiments

# Using the Torque

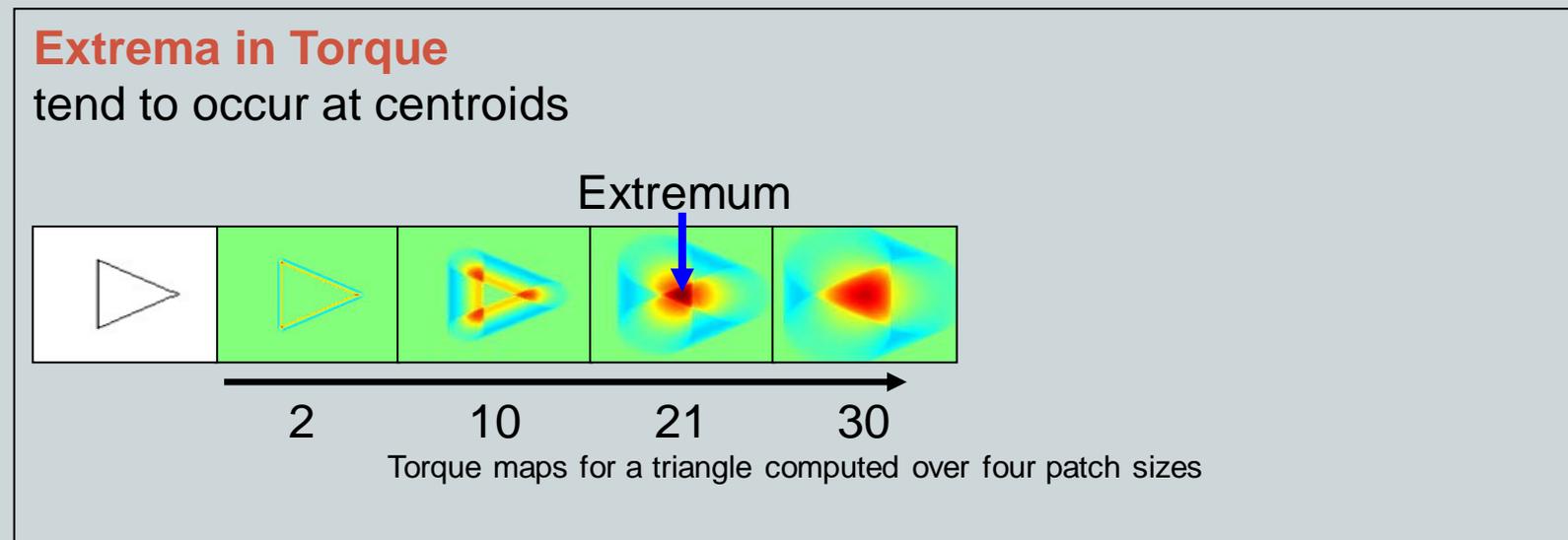
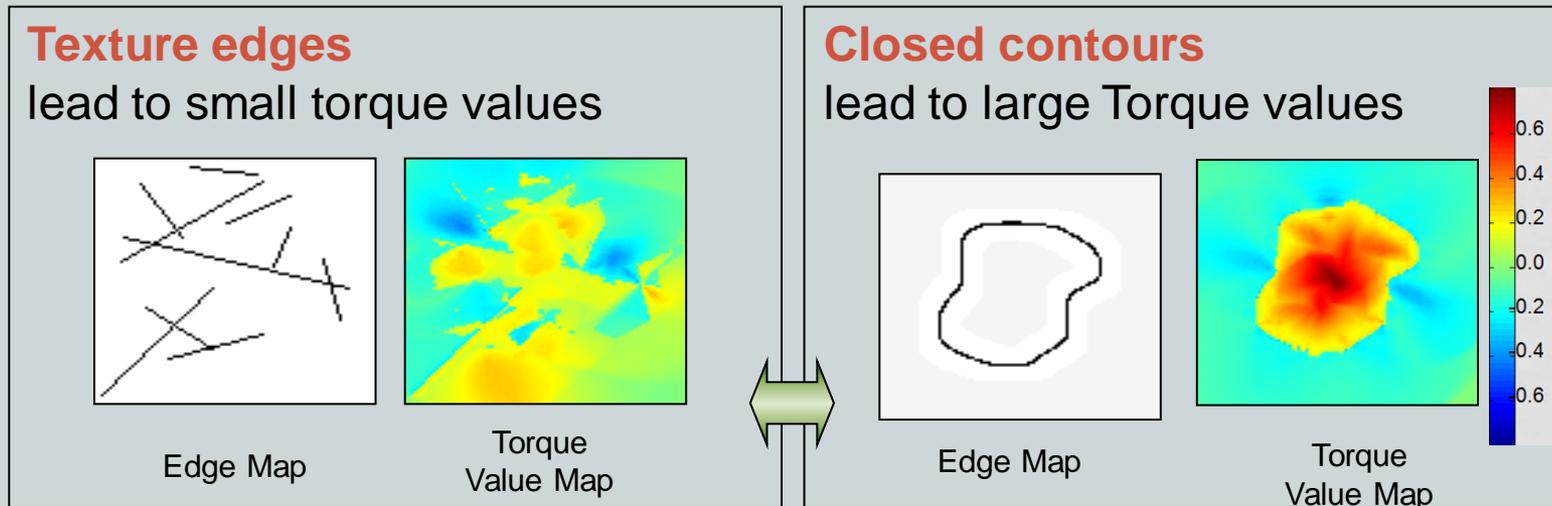


Torque values for different patch sizes

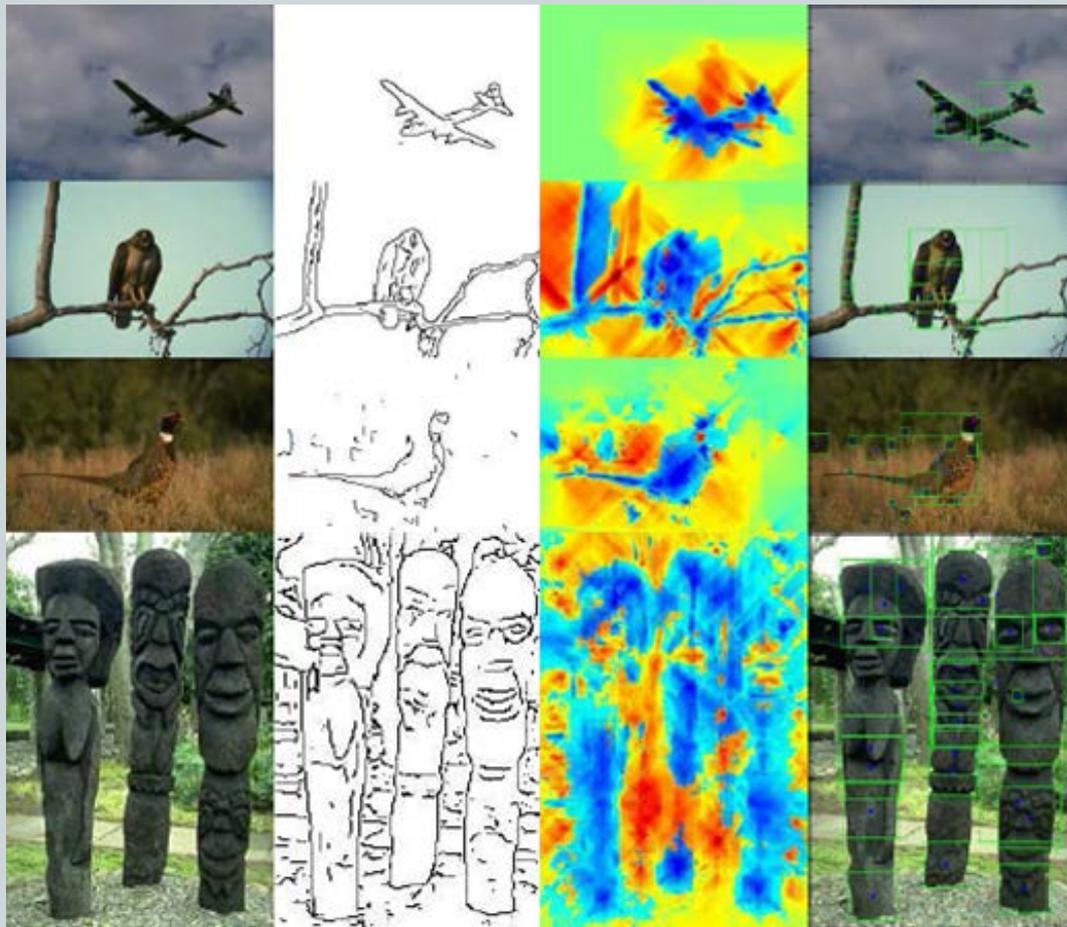
Combination of the torque values from all patch sizes



# Key properties of the Torque



# Torque Extrema



- a. Image
- b. Pb edges
- c. Torque value map
- d. Minima in Torque volume

a

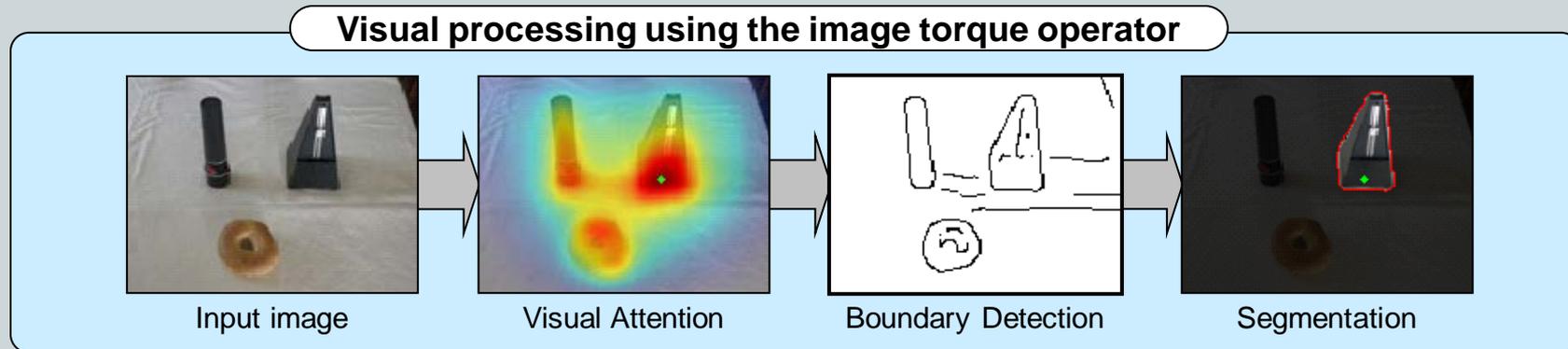
b

c

d

# Application

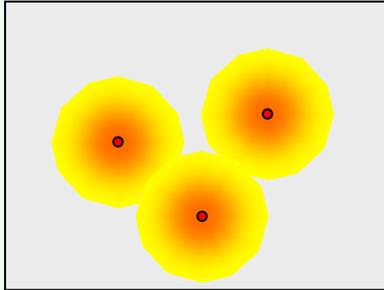
An active approach to finding an object in the scene consists of three modules: visual attention, boundary detection, and foreground segmentation.



We showed that by adding the torque we can improve state of the art methods

# Visual Attention

saliency map using torque

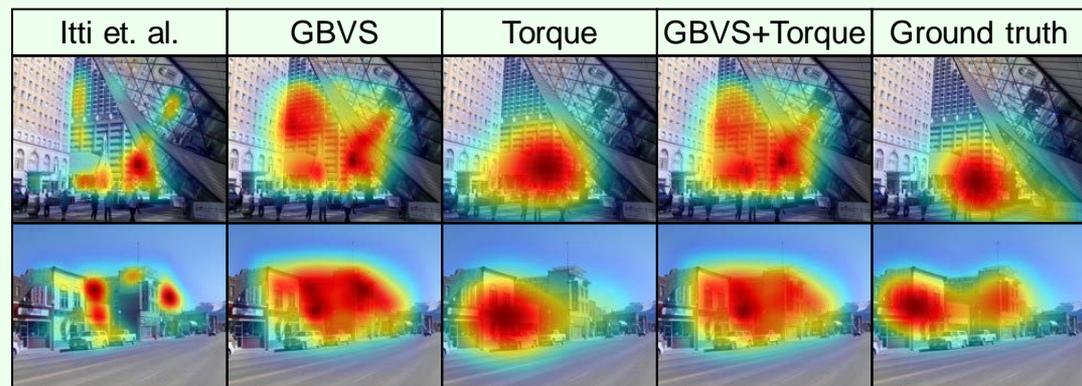
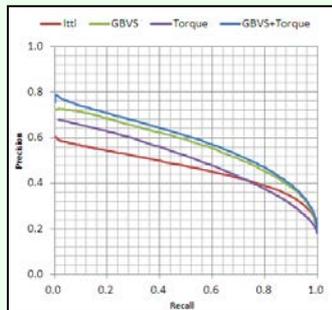


Gaussian distributions centered at torque extrema

Method	F-measure
Itti et. al.	0.53
GBVS(Harel et al.,2009)	0.59
Torque	0.54
<b>GBVS+Torque</b>	<b>0.60</b>

Improved by Torque

Evaluation on dataset by Judd et. al. (2009): F-measure and precision-recall curve.  
GBVS+Torque is with weights 0.7 and 0.3.



Examples of visual attention for two test images

# 3D volumetric video segmentation



# Detecting object specific contours

## Feedback to Mid-level Vision

1. Training: Obtain “prototypical contours” from annotated ground truth contours



2. Run time.

- **Match** partial contour fragments to models
- **Reweigh** torque based on matching scores

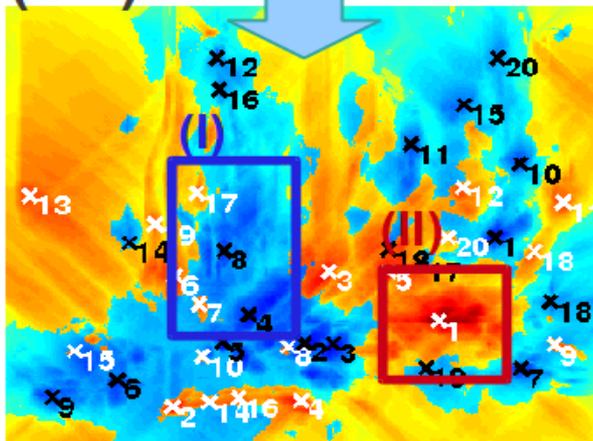
- C. L Teo, C.Fermüller, Y. Aloimonos. “A Gestaltist approach to contour-based object recognition: Combining bottom-up and top-down cues,” Intern.Journal of Robotics Research, 2015.

- M. Maynard, A. Guha (Y. Aloimonos, C. Fermüller) “Feedback from Vision,” Qualcomm Innovation Fellowship Award 2016.

(b-1)

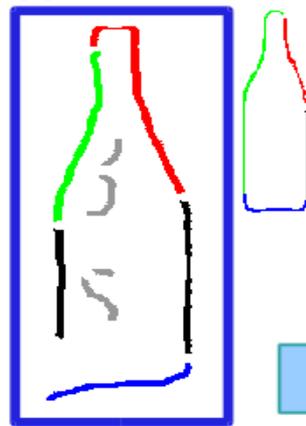


(b-2)

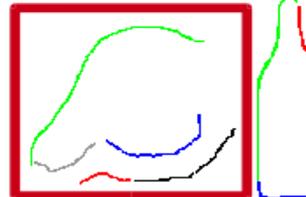


(c)

(I)



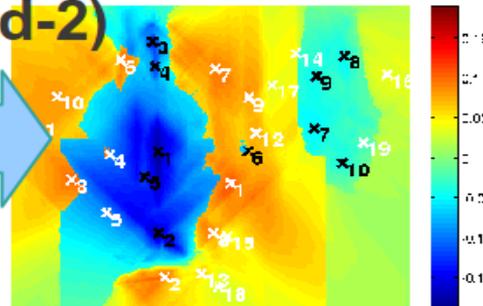
(II)



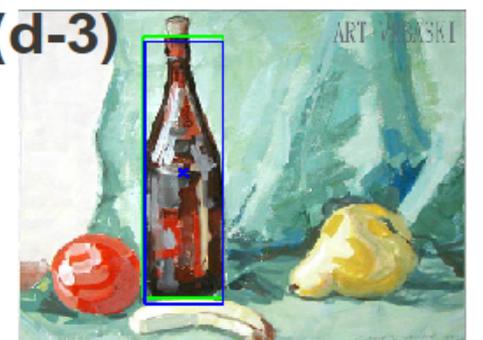
(d-1)



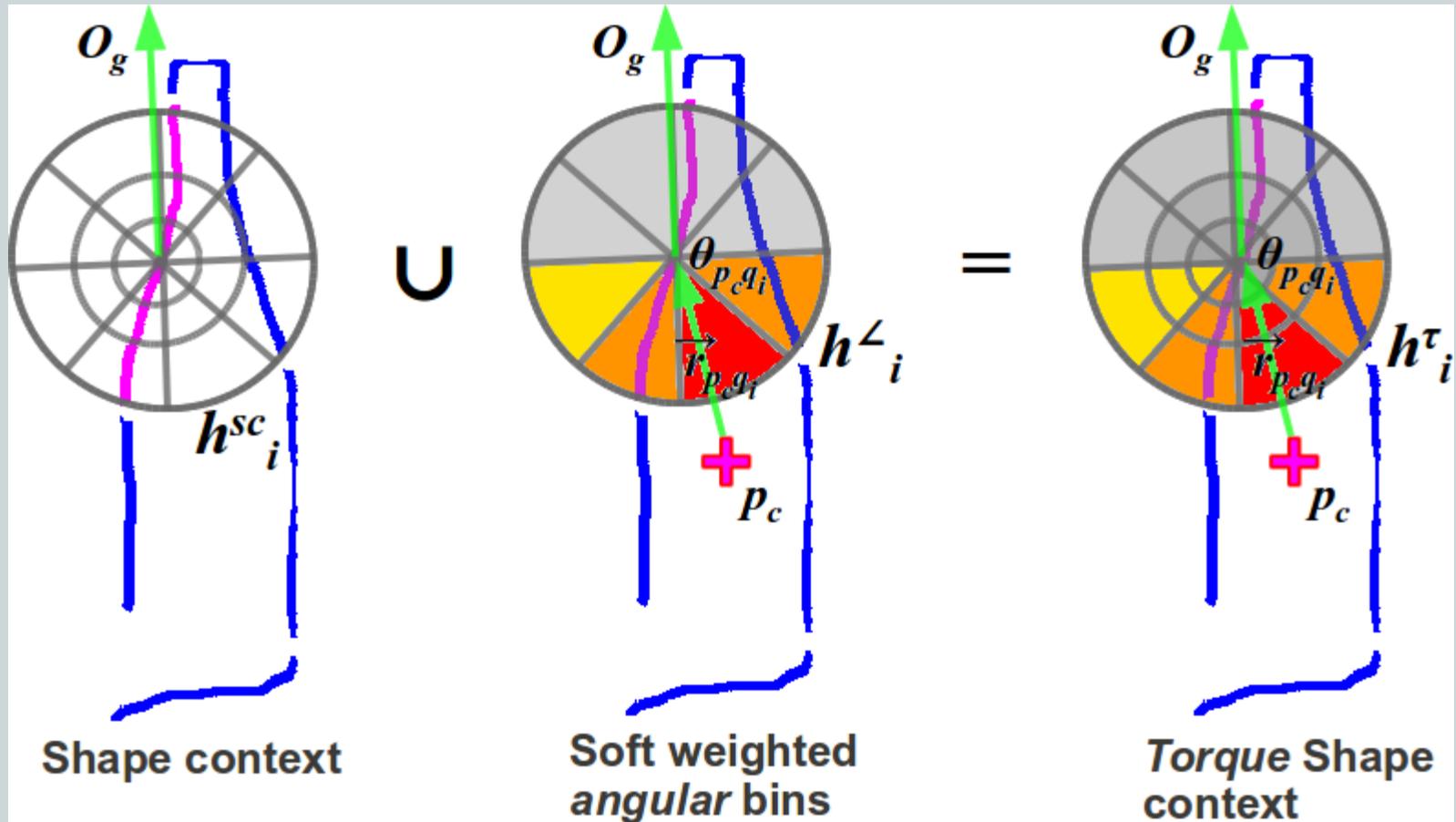
(d-2)



(d-3)



# Partial Contour Matching: Torque Shape-Context

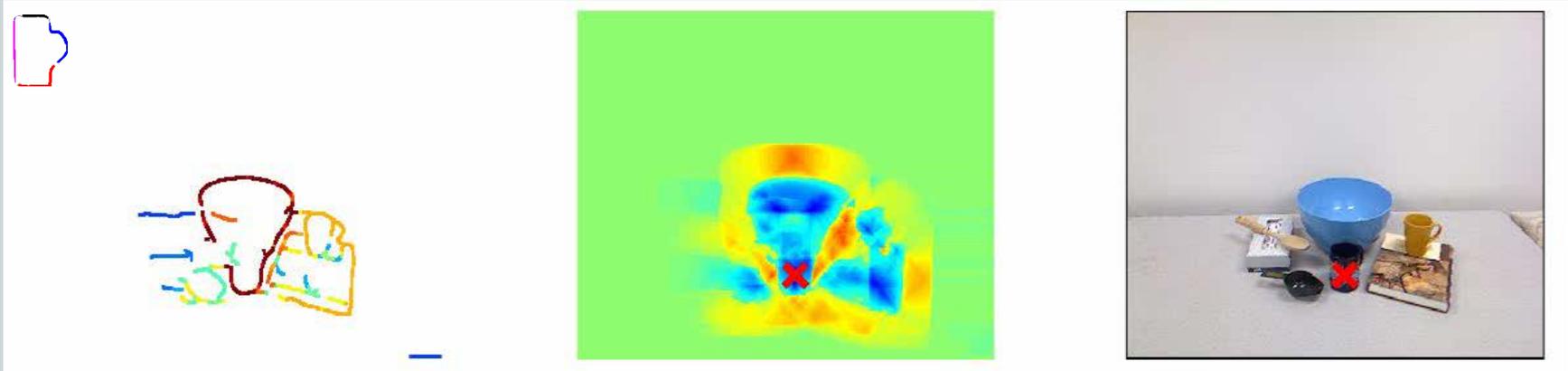


# Example results on Robot

Edges

$$\tau_P^m$$

RGB Input



# What is Border Ownership?

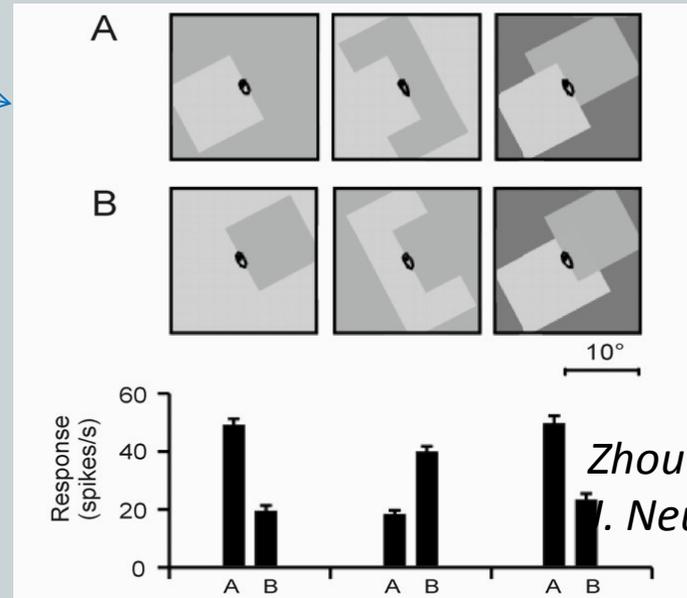


# Motivations: Psychological & Biological

V2 cell that prefers bright Figure on the left

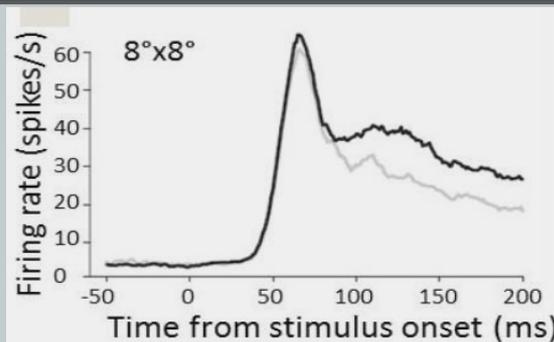


Selective nature of border ownership neurons (V2 & V4):



Zhou et al.,  
*J. Neuroscience* 2000

Fast response time,  
<75ms from  
stimulus onset.

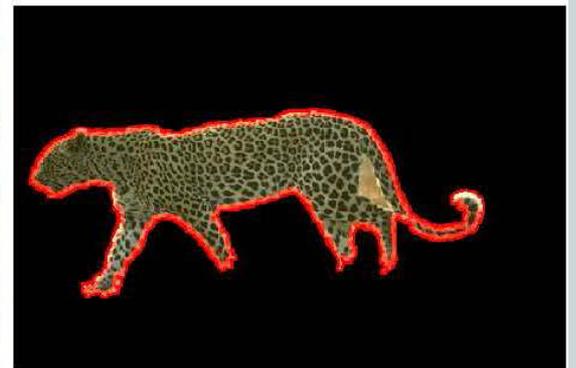
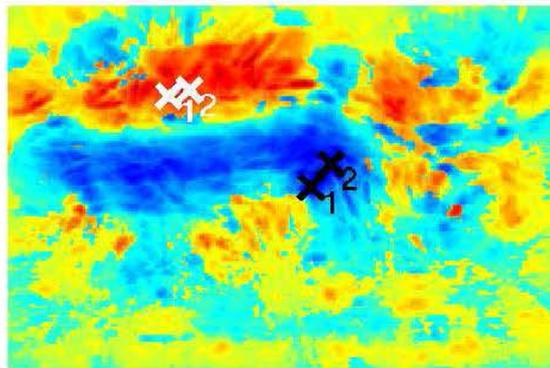
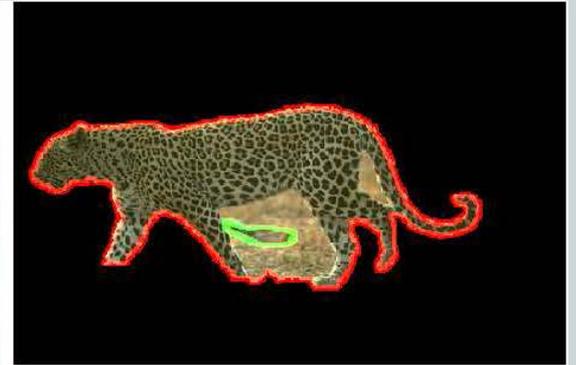
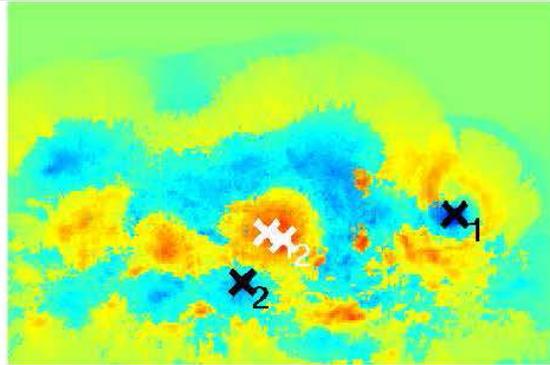


Sugihara et al.,  
*J. Neurophysiology* 2011

Figure-ground organization and  
attention are closely related

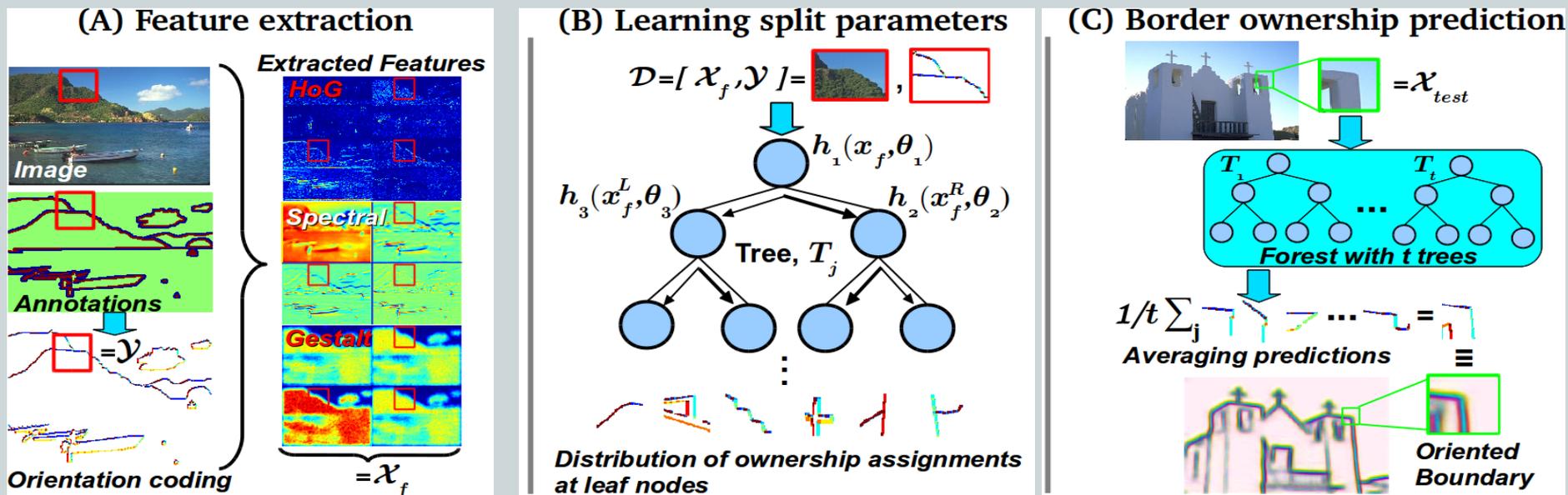
Craft et al.,  
*J. Neurophysiology* 2007

# Segmentation



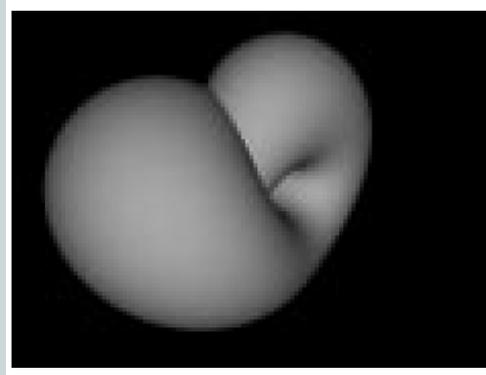
# Approach Overview

1. Extract patch-based features sensitive to ownership
2. Train a Structured Random Forest (SRF) that saves ownership structure at leaf nodes
3. Fast inference using SRF by averaging responses over all decision trees



# Feature Extraction: Local Ownership Cues

**Extremal edges** or *image folds* are characteristic changes in intensity along boundaries.

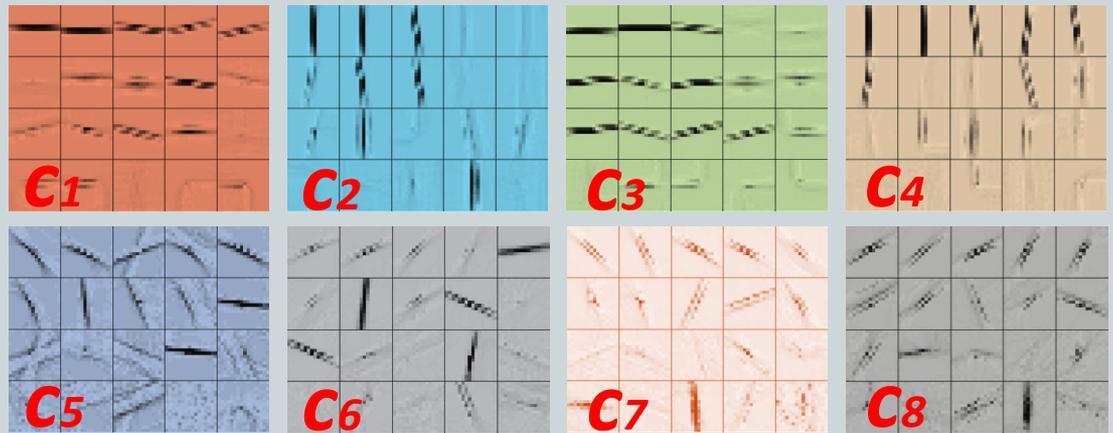
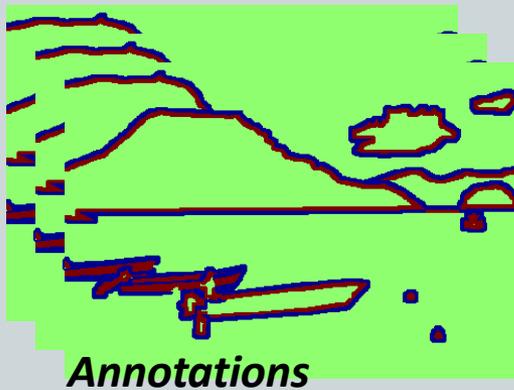


*Huggins & Zucker, ICCV 2001*

Psychophysical experiments have shown them to be one of the **strongest cues** for ownership.

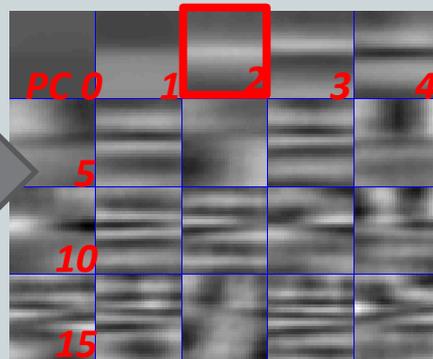
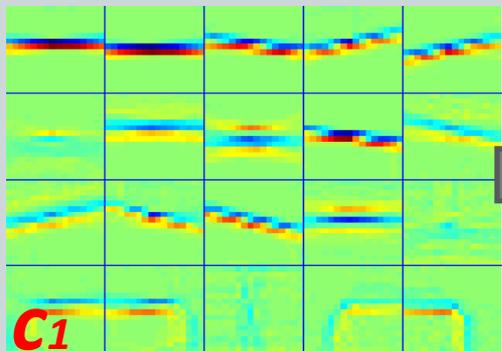
*Ghose & Palmer, J. Vision 2010*

# Local Ownership Cues



“Sketch token” clusters of 8 ownership directions

*Lim et al., CVPR'13*



**PC2** displays grayscale variations indicative of extremal edges.

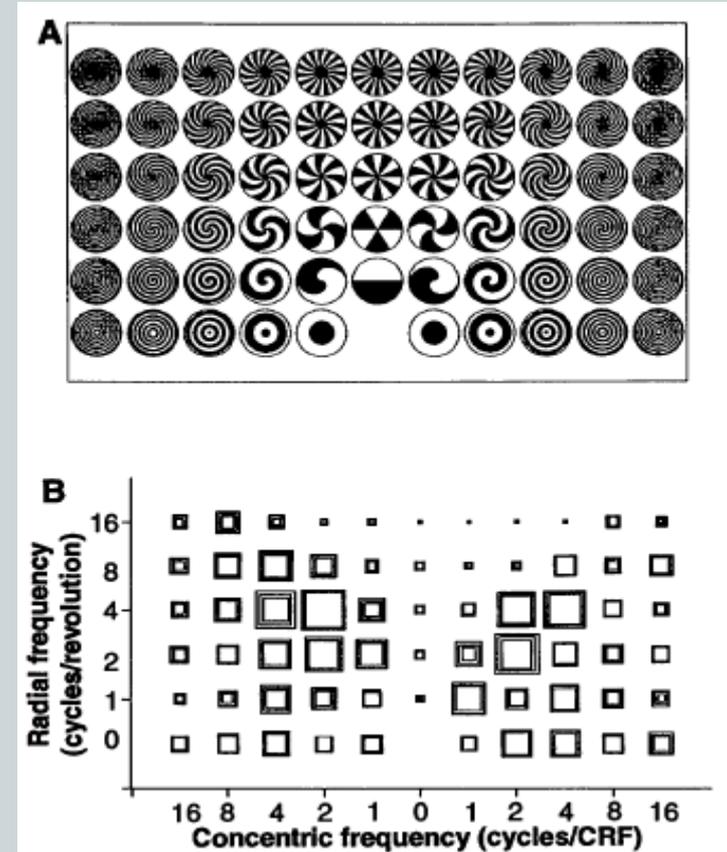
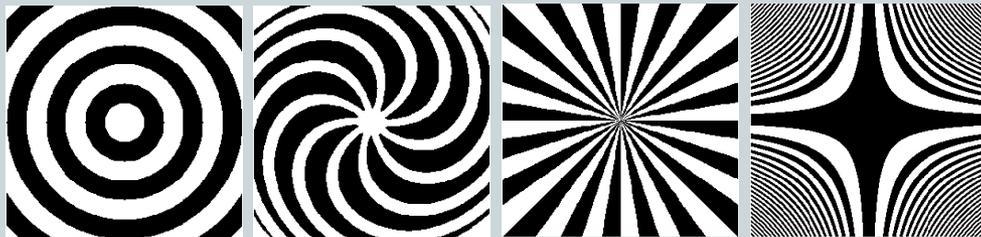
*Ramenahalli et al., CISS'11*

# Feature Extraction: Global Ownership Cues

Border ownership is also determined by longer range (global) contextual cues.

*Craft et al., J. Neurophysiology 2007*

Implementation through visual operators that capture four grouping or “Gestalt” patterns:



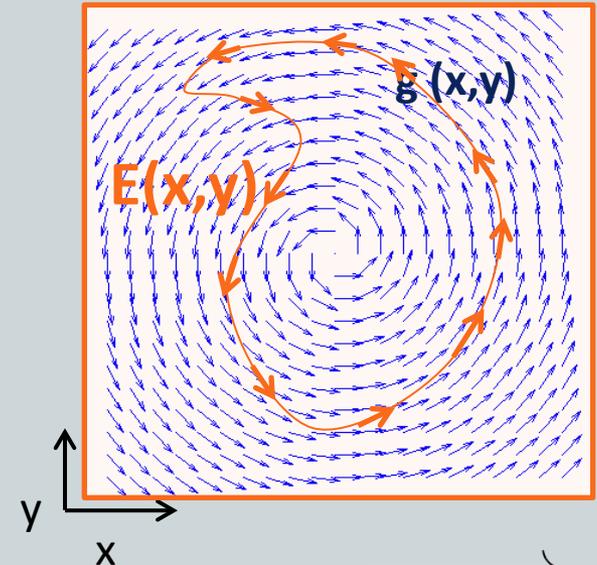
A: Gratings. B. Responses of a V4 cell

Cells tuned to these patterns have been observed area V4 of macaques :

*Gallant et al., Science 1993*

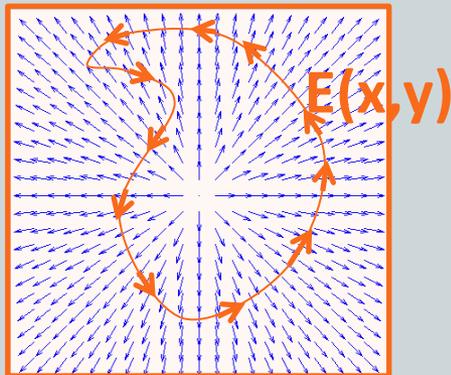
# Global Ownership Cues

Rewriting the image torque as a *scalar product of the edges* (tangent vectors)  $E(x,y)$  and a circular gradient field  $g(x,y)$ .



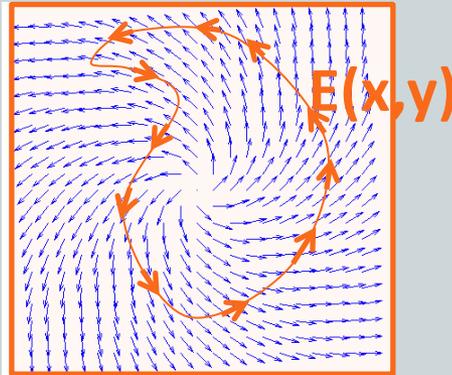
$$\tau(x, y) = \vec{E}(x, y) \times \vec{r} \quad \equiv \quad \tau(x, y) = \vec{E}(x, y) \cdot \overrightarrow{(-y, x)}$$

Radial



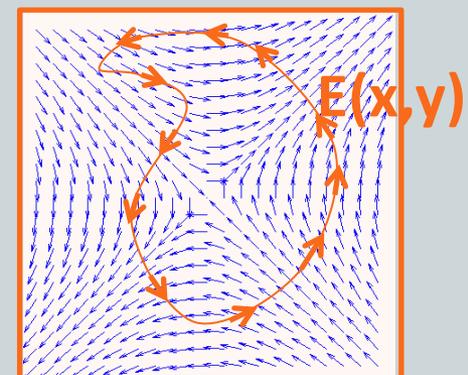
$$g(x, y) = (x, y)$$

Spiral



$$g(x, y) = (ax - y, y + ax)$$

Hyperbolic gradient field



$$g(x, y) = (ay, x)$$

# Results

Predicted **boundaries (red)** and ownership (**FG: green, BG:blue**)



**BSDS** (100 training/100 testing)

*Martin et al., PAMI 2004*

**NYU-Depth** (795 training/ 654 testing)

*Silberman et al., ECCV 2012*

Ownership prediction accuracy:

*Ren et al., ECCV 2006*  
*Leichter & Lindenbaum, ICCV 2009*

Feature set	BSDS	NYU-Depth
HoG	72.0%	66.0%
+ Spectral (no contour tokens)	73.1% (72.0%)	67.0% (65.6%)
+ Spectral (contour tokens)	74.0% (72.3%)	68.1% (66.7%)
+ Gestalt patterns	74.4% (72.7%)	<b>68.4%</b> (66.7%)
All features + Spectral (NYU)	<b>74.7%</b> (72.8%)	-
Global-CRF	69.1%	-
2.1D-CRF	68.9%	-

( ) denotes use of single features

Boundary prediction accuracy:

*Arbelaez et al., PAMI 2011*  
*Dollar et al., PAMI 2015*

Method	BSDS-500	NYU-Depth
Our approach	0.73,0.74,0.76	0.63,0.64,0.60
gPb-owt-ucm	0.73, <b>0.76</b> ,0.73	0.63,0.66,0.56
SE	0.73,0.75, <b>0.77</b> (SE-SS)	<b>0.65,0.67,0.65</b> (SE-RGB)

# Results



**Red: Boundaries, Green: Foreground, Blue: Background**

# Symmetry in 2D

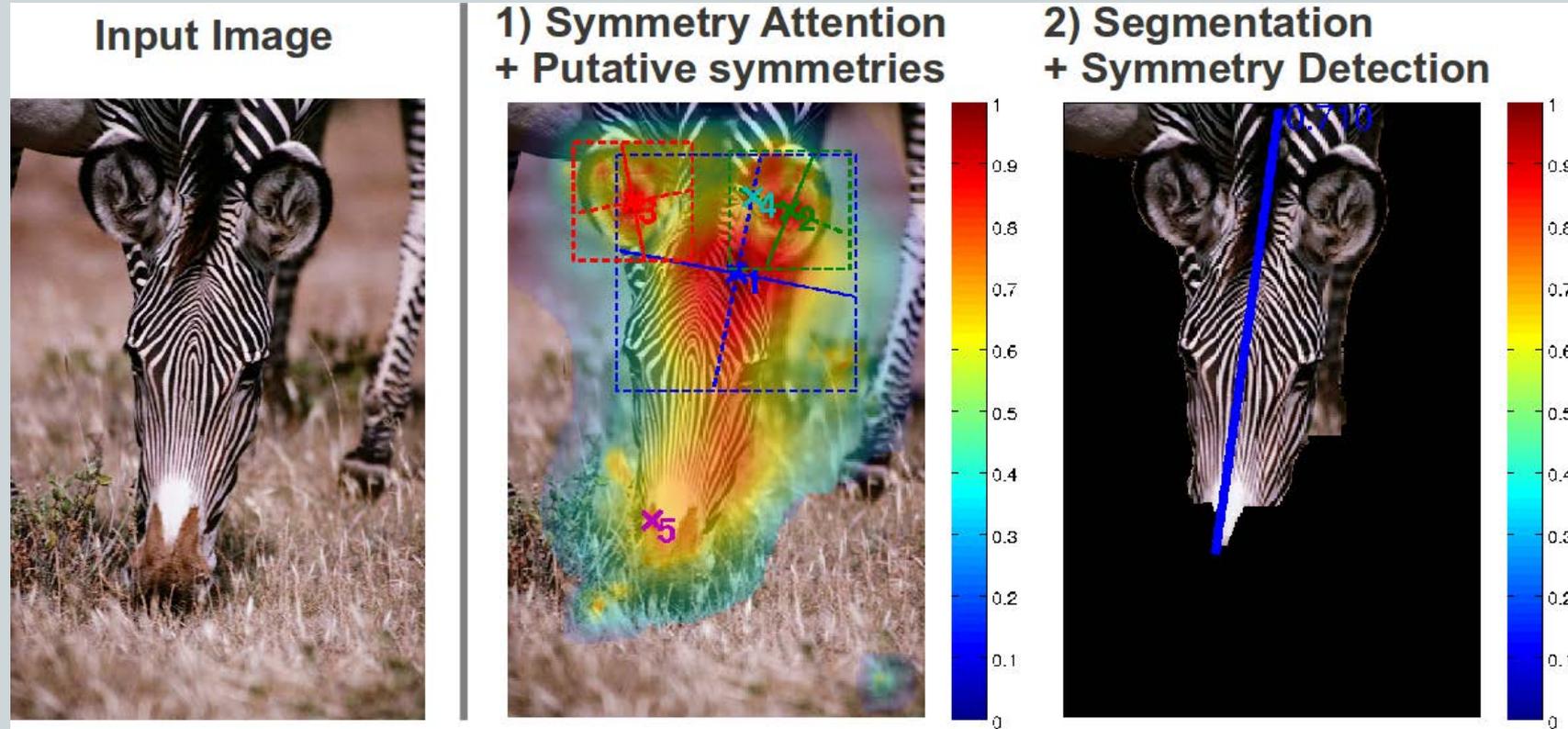
- Goal: Detect symmetries in complex environments containing *clutter*:



- Key challenges:
  - *Where* to compute the symmetries? → **Attention**
  - *How* to compute the symmetries reliably? → **Statistics**

# Proposed Solution

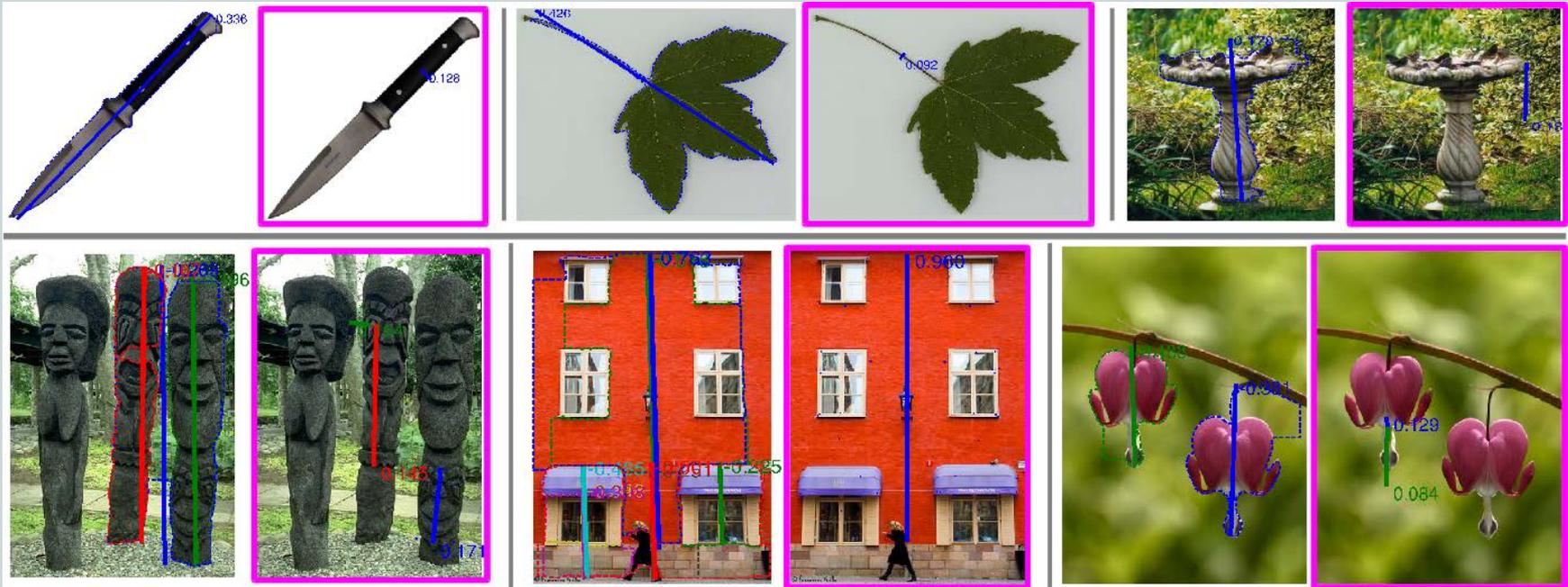
(2 steps):



Segmentation applied per fixation point: resolves the scale issue since the object region is selected

# Results

Top: singles



Bottom: multiples

Left: Our approach,

Right: Loy & Eklundh, ECCV 2006

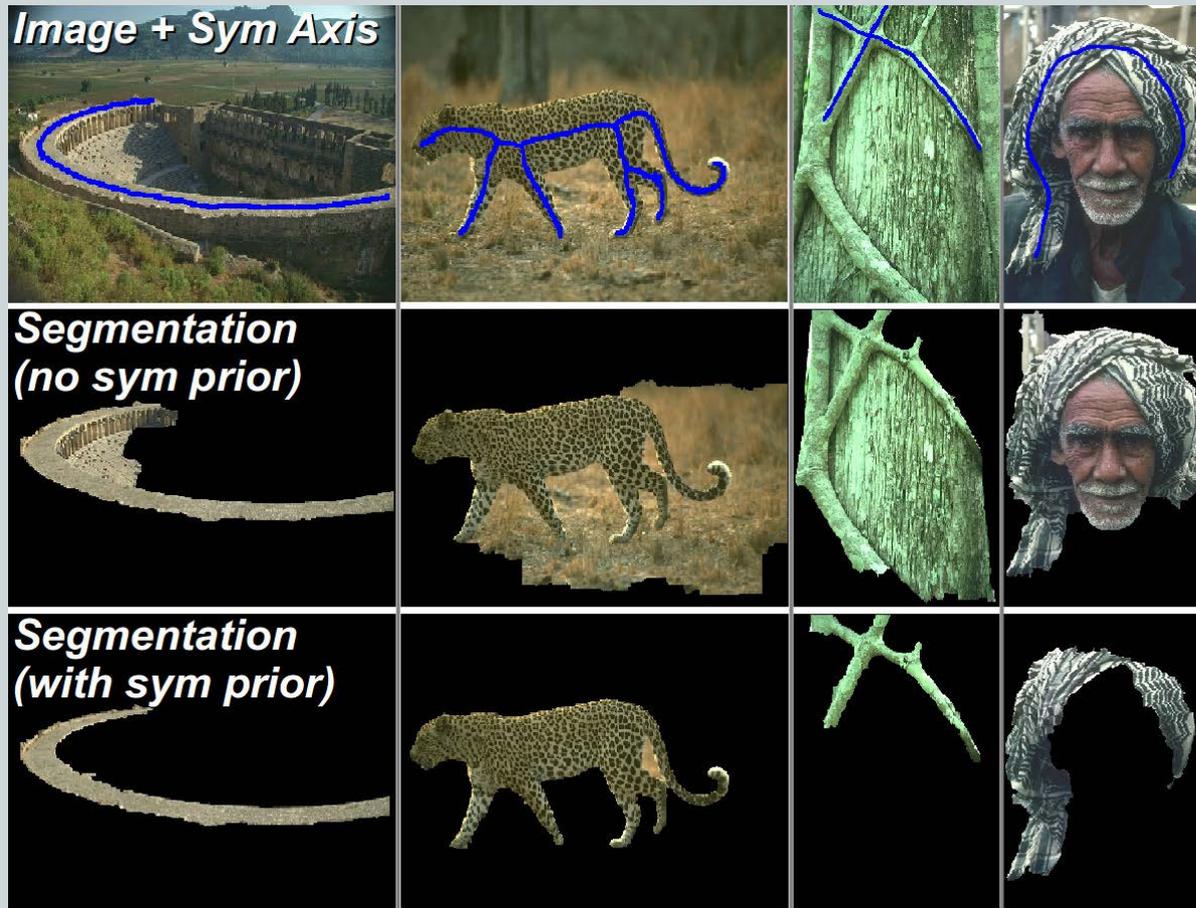
# Segmentation with Symmetry

$$E(f) = \sum_{p \in \mathcal{P}} U_p(f_p) + \sum_{(p,q) \in \mathcal{N}} V_{pq}(f_p, f_q) + \sum_{(p,p') \in \mathcal{N}_{sym}} S_{pp'}(f_p, f_{p'}) + \sum_{(p,q) \in \mathcal{N}} B_{pq}(f_p, f_q)$$

Unary + Pairwise (edges)

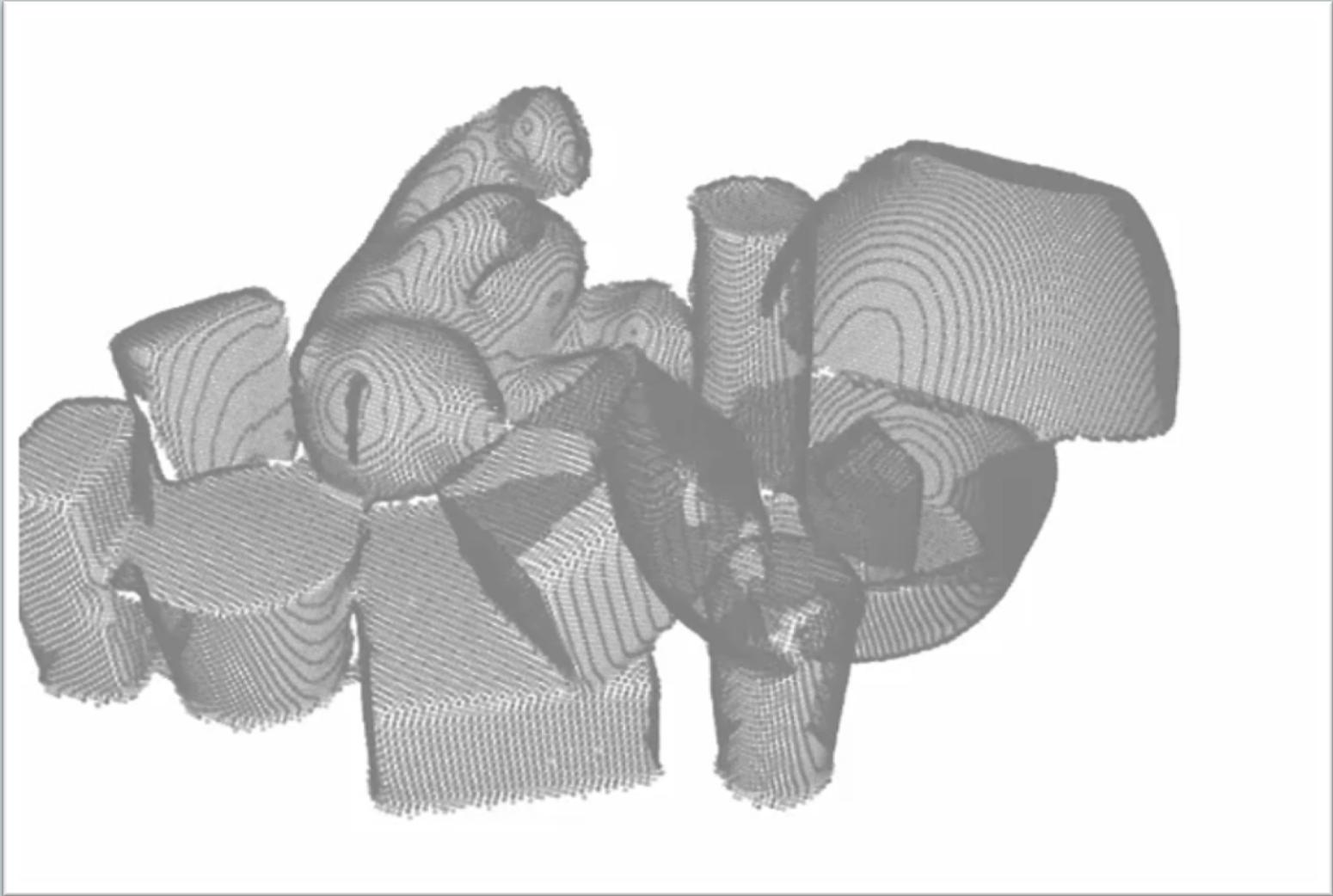
**Cross-symmetry + 'Ballooning'**

# Segmentation with symmetry axes





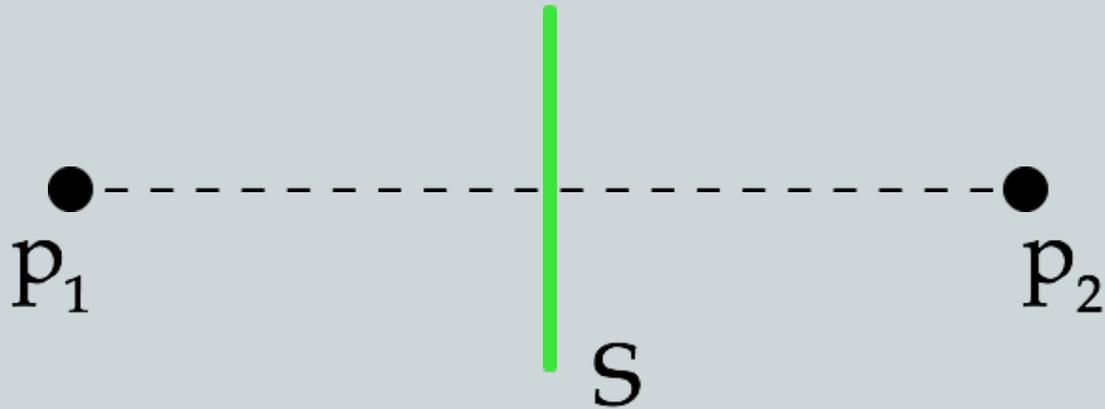
# Data



Ecins, C. Fermüller, and Y. Aloimonos, "Cluttered Scene Segmentation Using the Symmetry Constraint," ICRA, 2016.

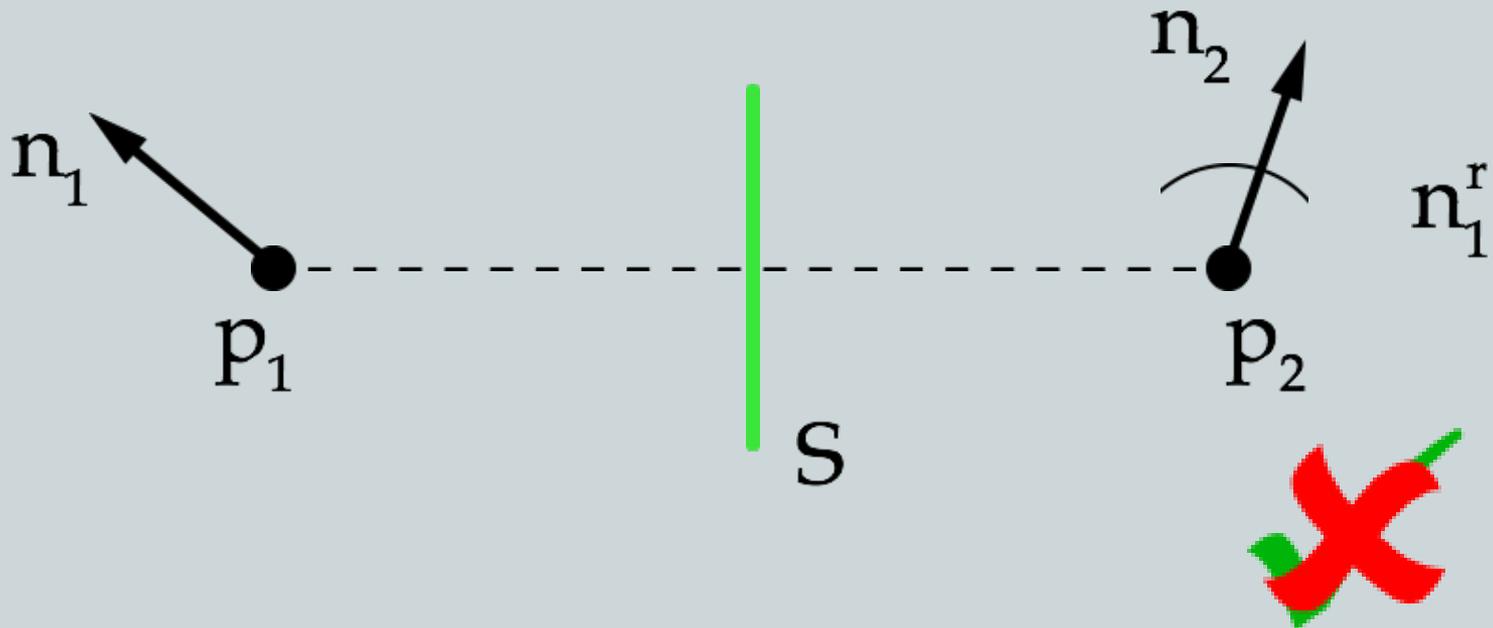
# Symmetric point correspondence

Two points in space uniquely define a reflectional symmetry plane.

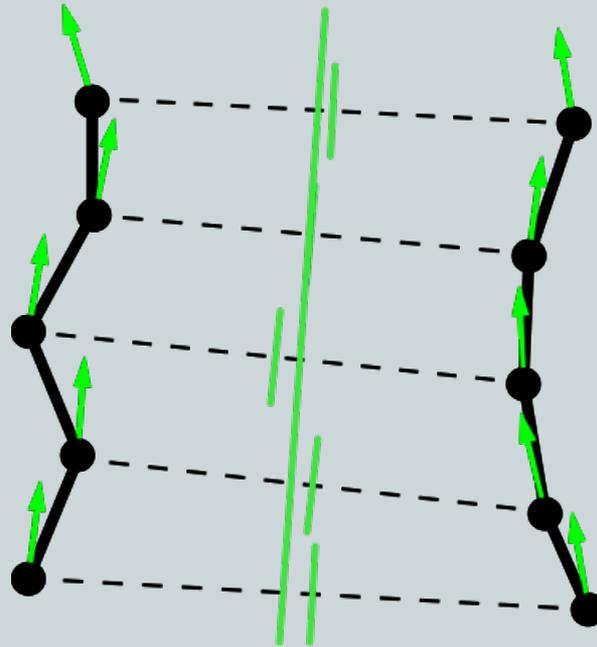


# Symmetric point correspondence

Two points in space uniquely define a reflectional symmetry plane  
Two oriented points in space form a symmetric match if their reflected normals align



# Symmetry detection

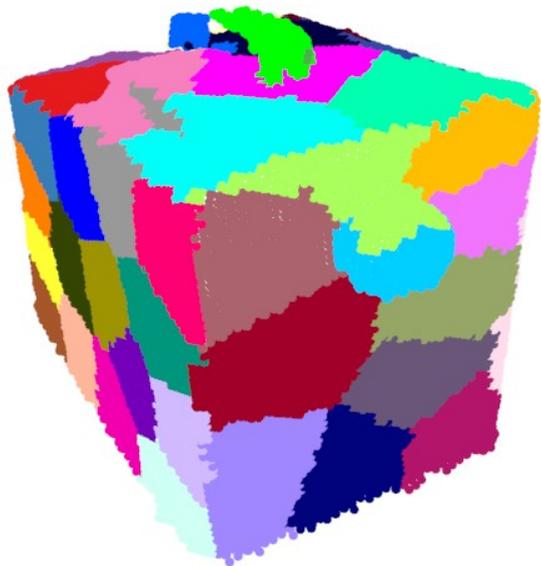


- Find symmetric correspondences between points
- Get a symmetry hypothesis for each correspondence
- Filter hypotheses using mean shift clustering

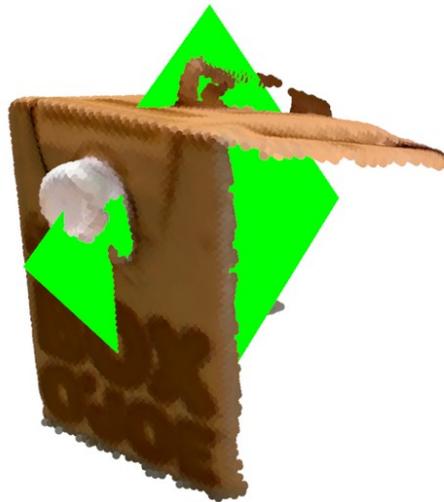
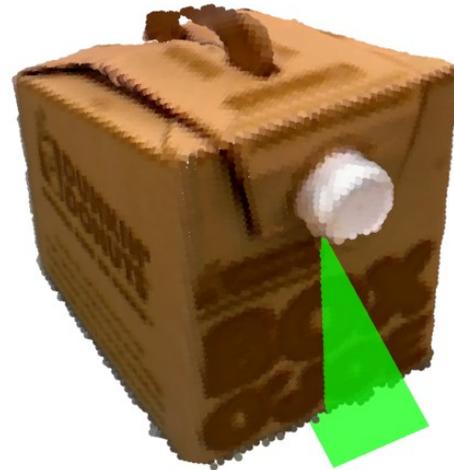
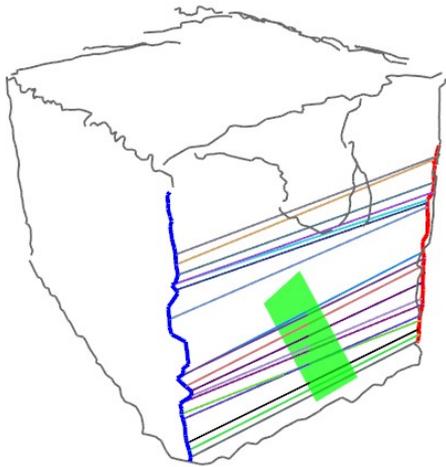
# Input Point clouds



# Features: Edges



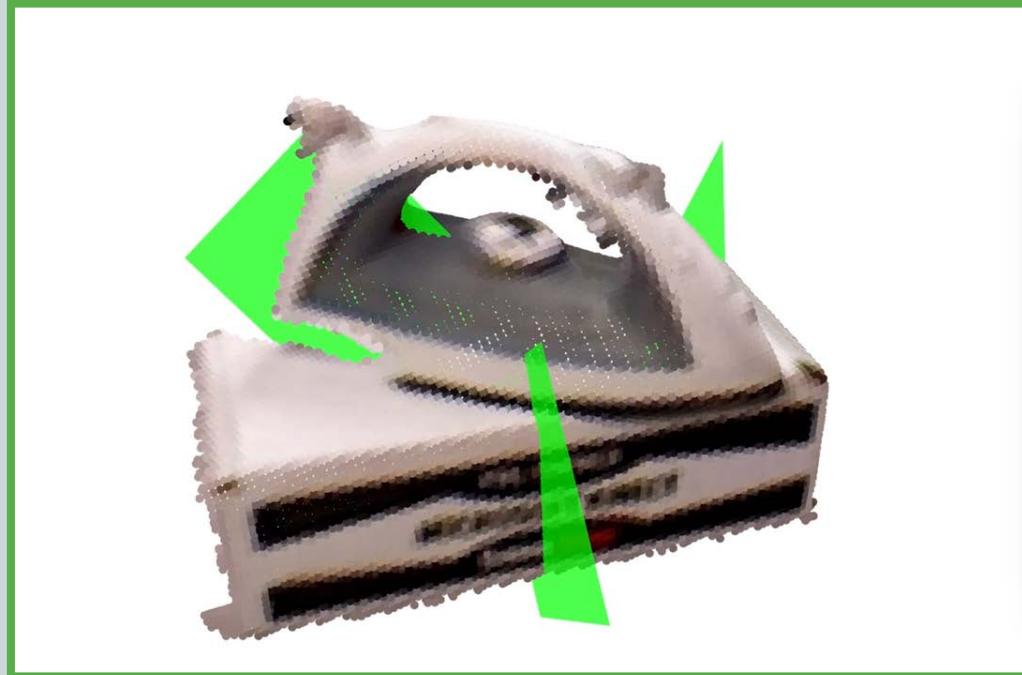
# Symmetry Detection



# Segmentation

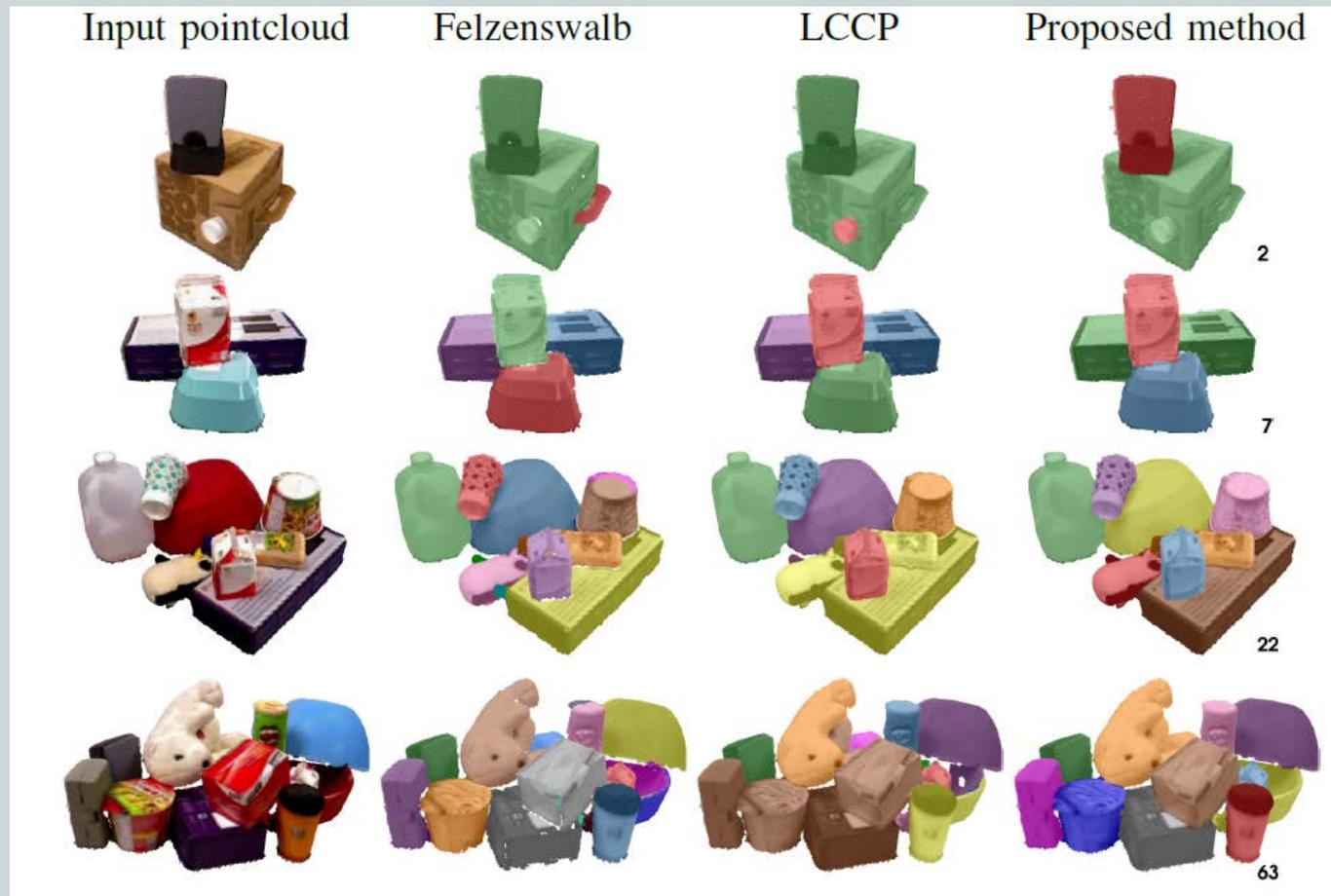
Grouping principles used:

- Convexity (old)
- Symmetry consistency (new)



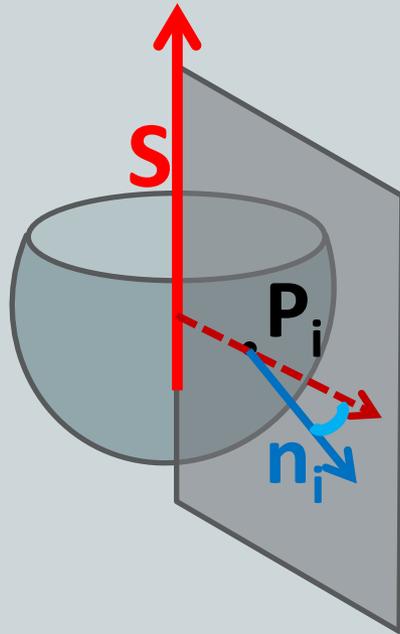


# Segmentation with 3D Symmetry

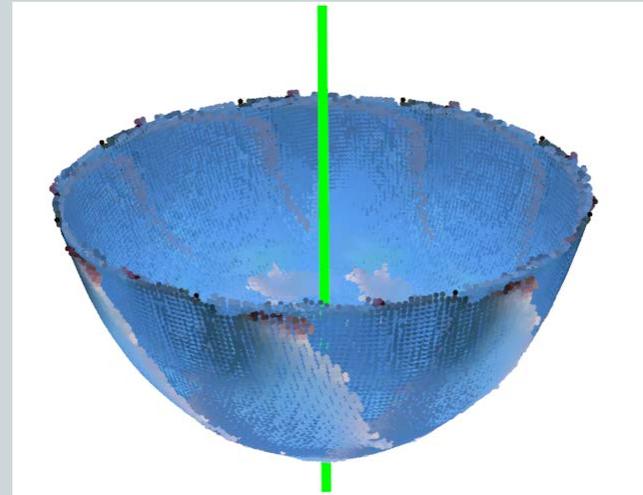
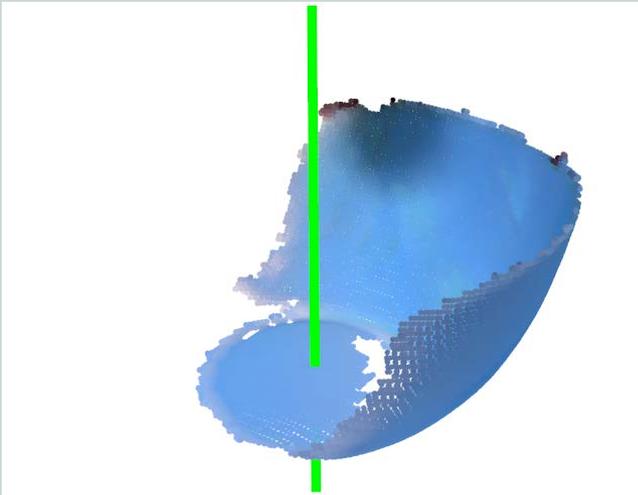


- A. Ecins, C. Fermüller, and Y. Aloimonos, "Cluttered Scene Segmentation Using the Symmetry Constraint," ICRA, 2016.
- LCCP: S. C. Stein, M. Schoeler, J. Papon, and F. Worgötter, "Object partitioning using local convexity," CVPR, 2014
- Felzenswalb adaptation: A. Karpathy, S. Miller, and L. Fei-Fei, "Object discovery in 3d scenes via shape analysis," ICRA, 2013.

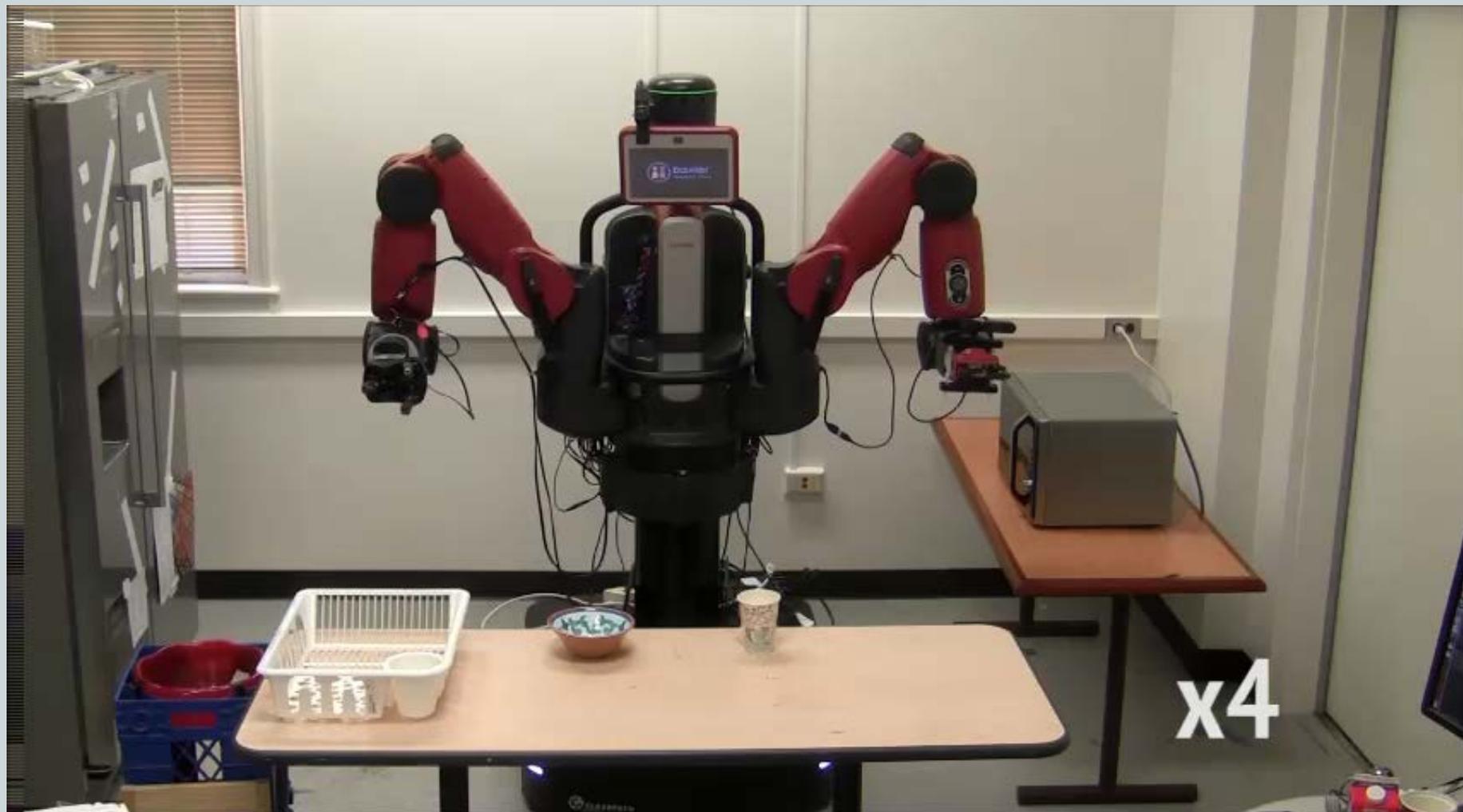
# Rotational Symmetry



$$\min_S \sum_i \angle(\text{Plane}(S, P_i), n_i)$$



# Heating a dish in the microwave



# Summary

- Mid-level concepts implemented as image operators
- Bottom-up principles of closure, symmetry, border-ownership
- Top-down task driven modulation of mid-level features
- Symmetry in 3D for object detection and segmentation