# Scene Representation for Robots:
# From Elementary Behaviors to Grasping
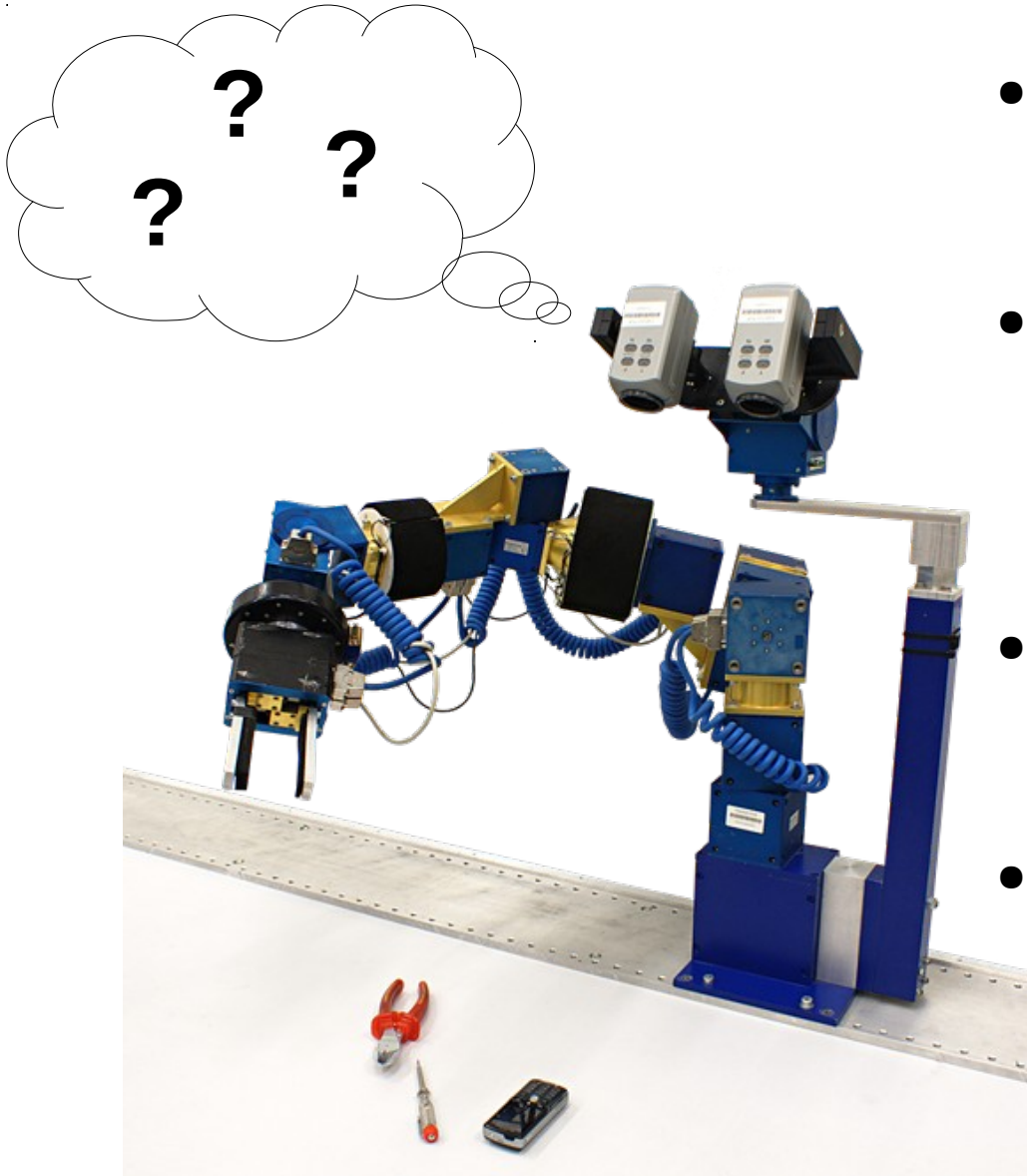
## Stephan Zibner

# Outline

- Representing Scenes
  - In Humans
  - Robotic Scenario
- Building Blocks
- Elementary Behaviors
  - Exploration
  - Maintenance
  - Query
- Reaching and Grasping

# Representing Scenes

# Scene Representation



- internal representation of environment

- foundation for every higher cognitive operation and action

- stable despite eye and body movements
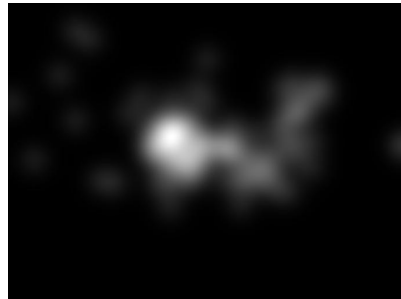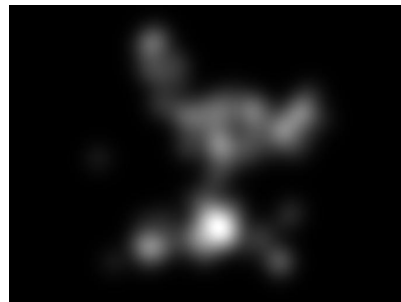
- limited capacity, link to long-term memory

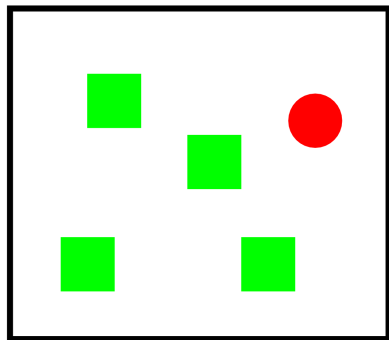# Scene Representation in Humans

image

mean eye placement



- human eye movement is not random (saliency)

Bruce, Tsotsos (2009)

# Scene Representation in Humans
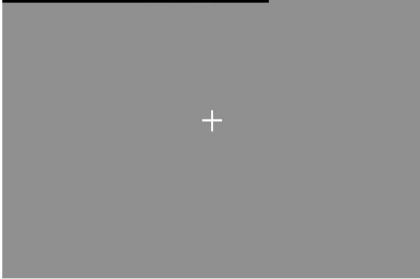

feature search


conjunctive search



Bottom: Treisman, Gelade (1980)

- human eye movement is not random (saliency)

- visual search highlights target objects

# Scene Representation in Humans
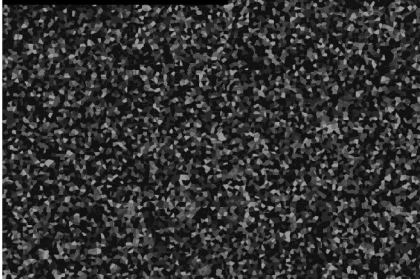


1. Fixation, 1000 ms
2. Initial Scene, 20 s
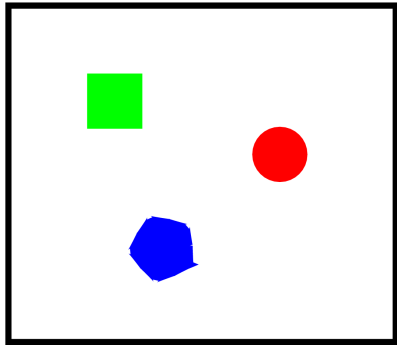3. Dot Onset, 150 ms
4. Initial Scene, 200 ms
5. Mask, 200 ms
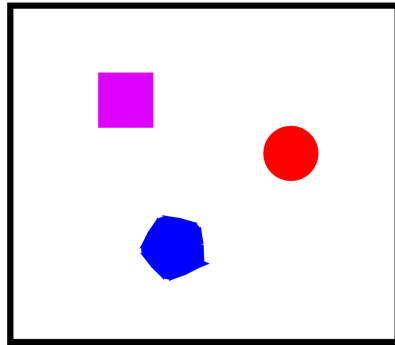6. Test Scene, Until Response

Hollingworth (2005)

- human eye movement is not random (saliency)
- visual search highlights target objects
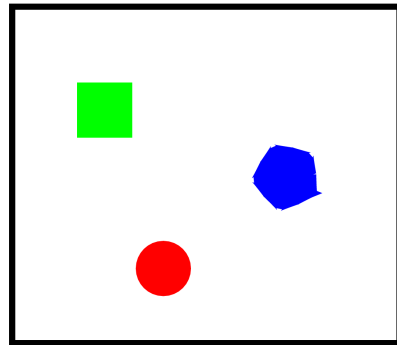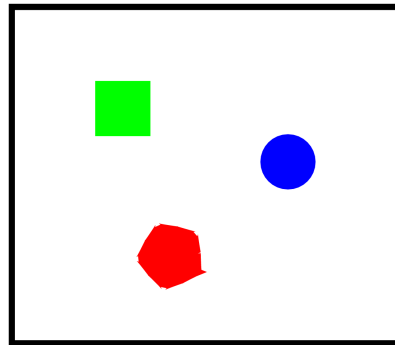
- details of a scene are kept in memory

# Scene Representation in Humans
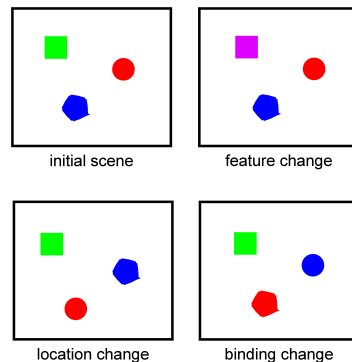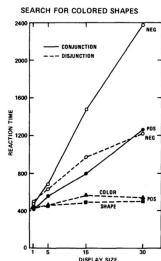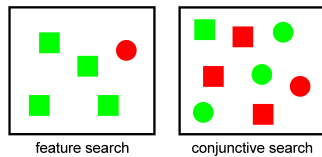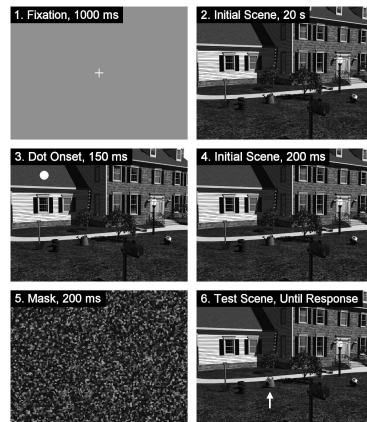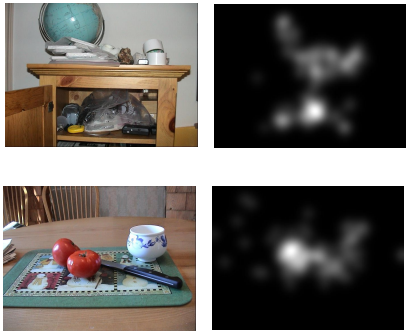


initial scene

feature change

location change

binding change

- human eye movement is not random (saliency)

- visual search highlights target objects

- details of a scene are kept in memory
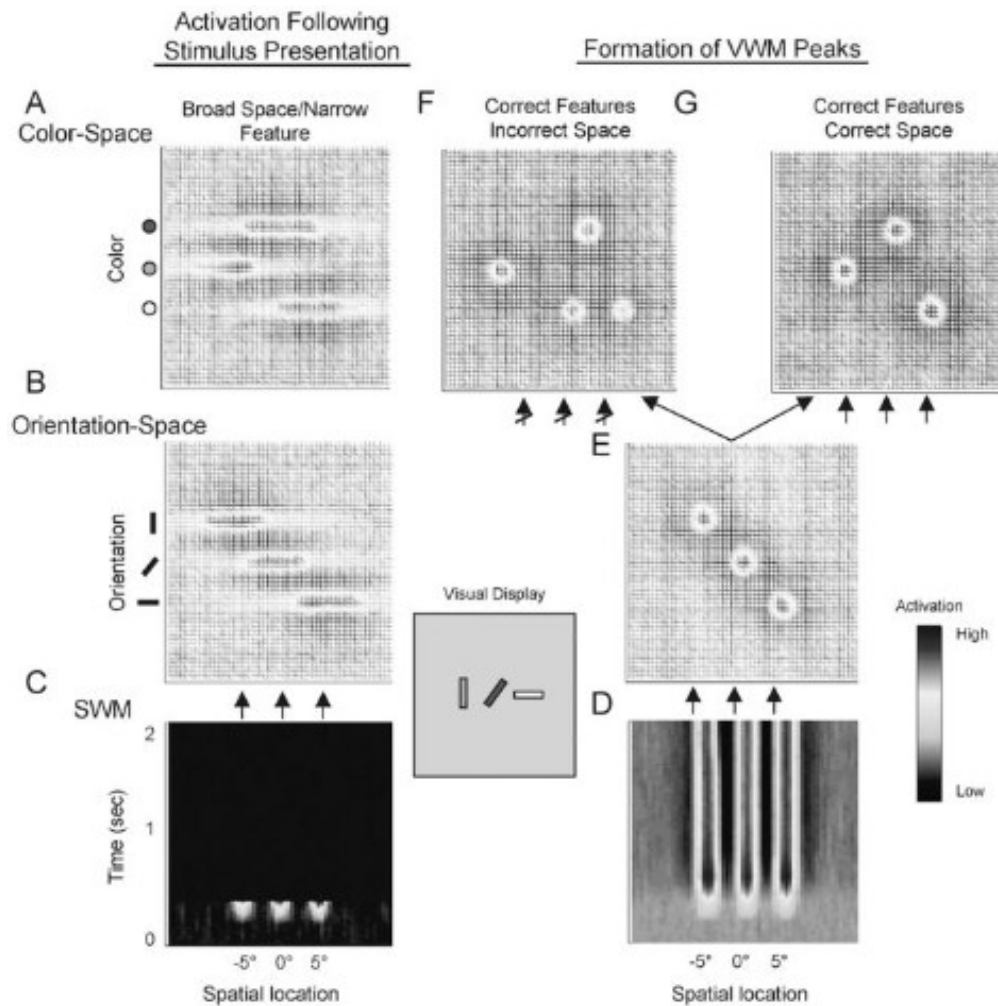
- no full representation (change blindness)

# Scene Representation in Humans



- human eye movement is not random (saliency)

- visual search highlights target objects

- details of a scene are kept in memory
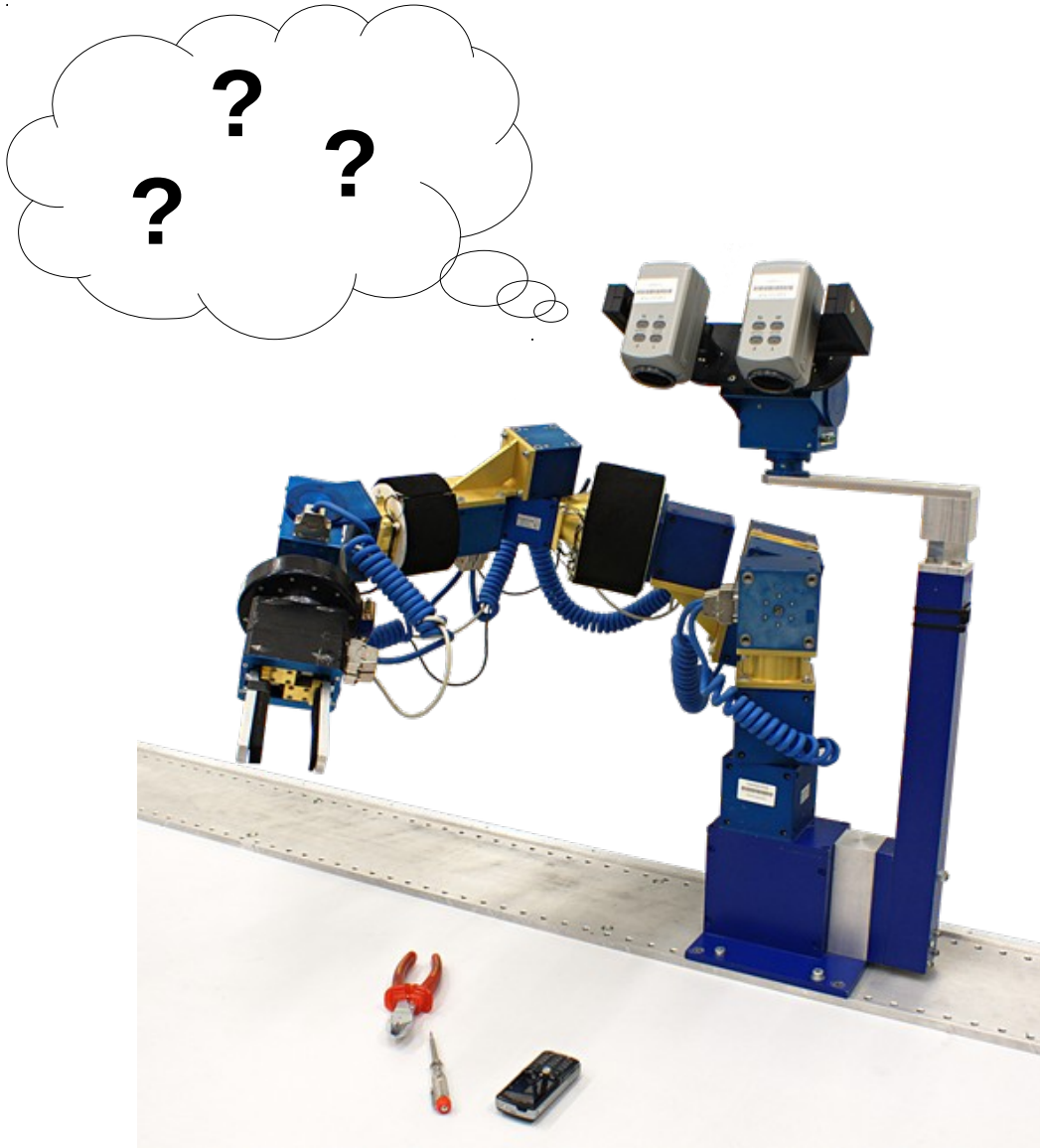
- no full representation (change blindness)

**attention is a key theme in visual processing**

# DFT Model



Activation Following Stimulus Presentation

Formation of VWM Peaks

A Color-Space — Broad Space/Narrow Feature

F Correct Features Incorrect Space

G Correct Features Correct Space

B Orientation-Space

C SWM

E Visual Display

D

Activation: High / Low

Spatial location -5° 0° 5°
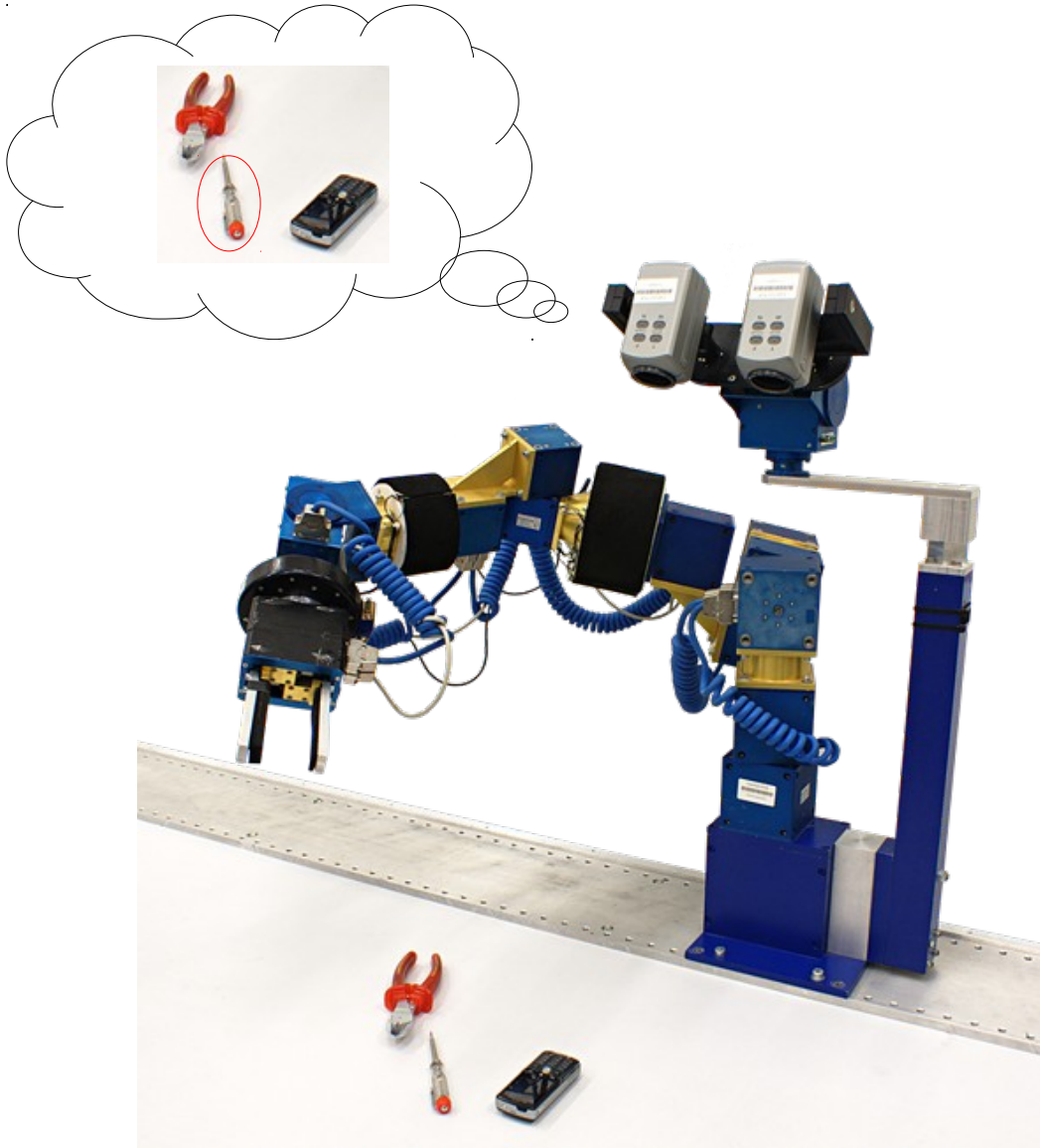
Johnson, Spencer, Schöner (2008)

- model of human visual working memory based on dynamic neural fields

- representation of low-level features

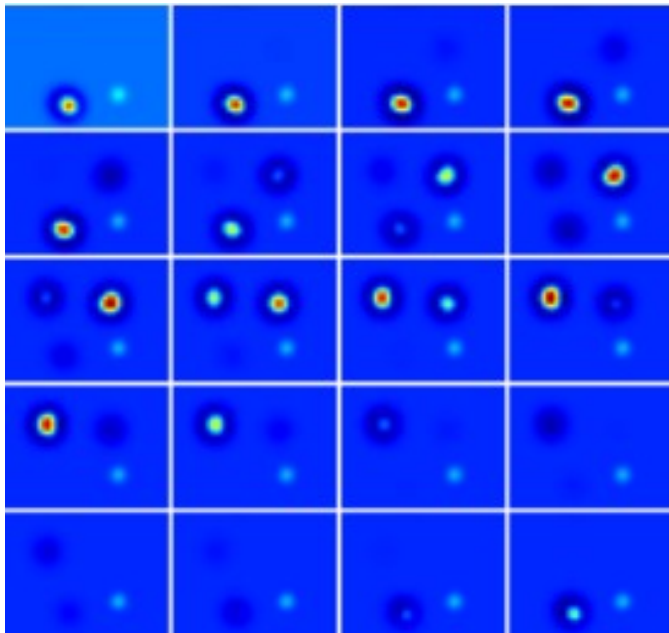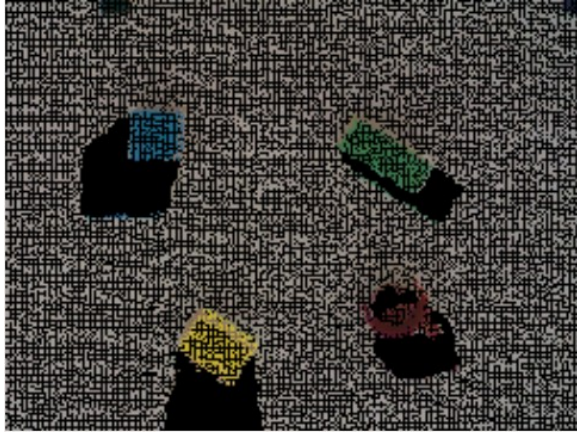- decomposition of features, binding through shared space

# Robotic Scenario



- apply to table-top scenario and human interaction

- use the internal representation for behavior generation (e.g., grasping)

- interact with humans ("hand me the red screwdriver", "what's to the left of the pliers?")

# Robotic Scenario



- explore the environment and store objects and their features internally

- maintain the internal representation

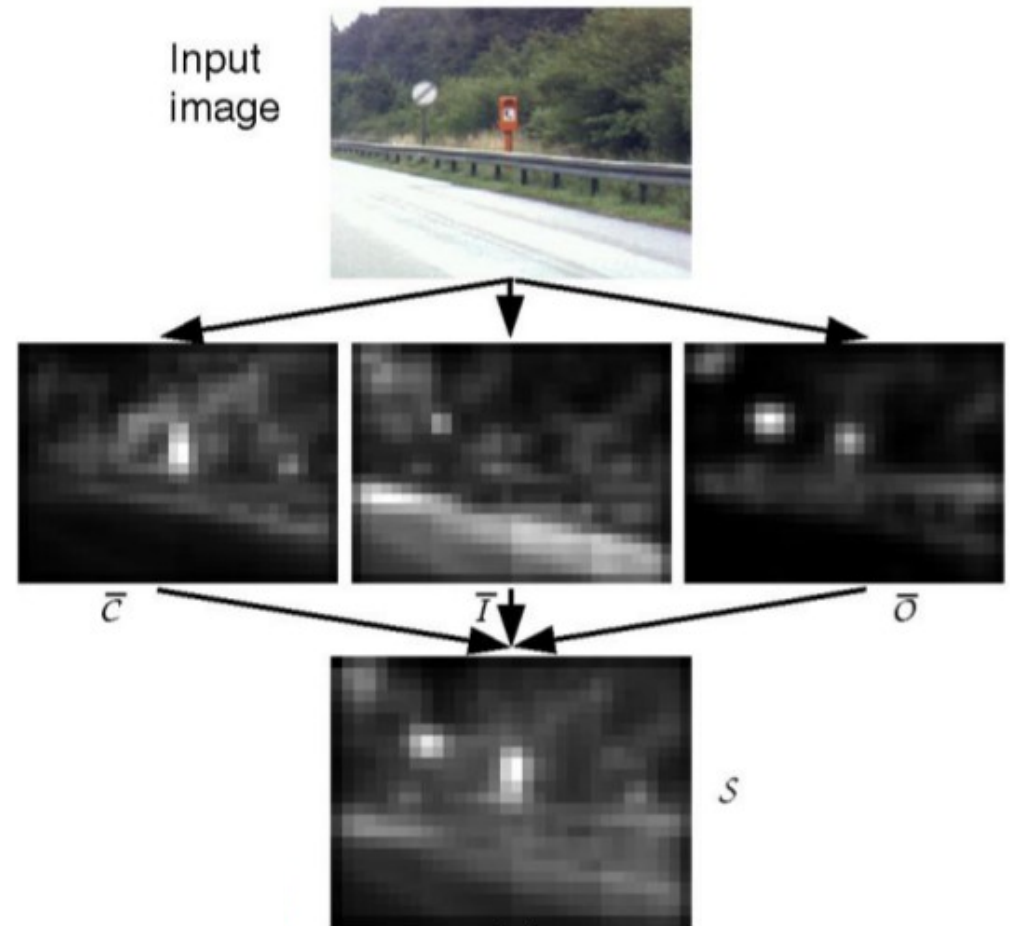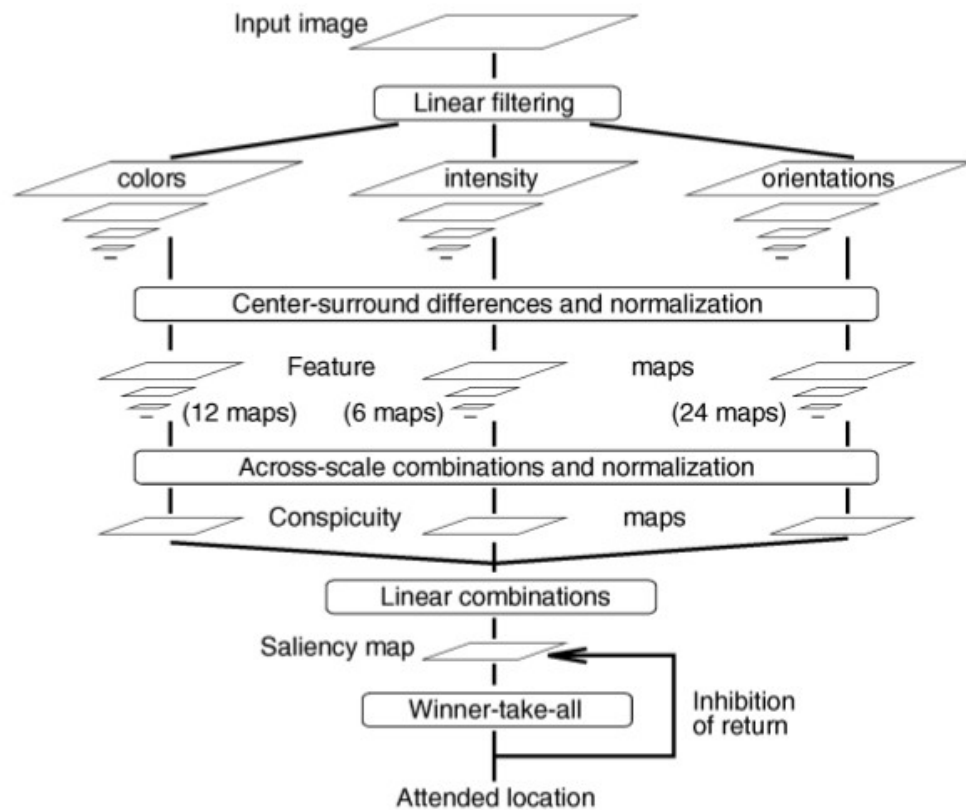- query the representation to create autonomous action based on the representation
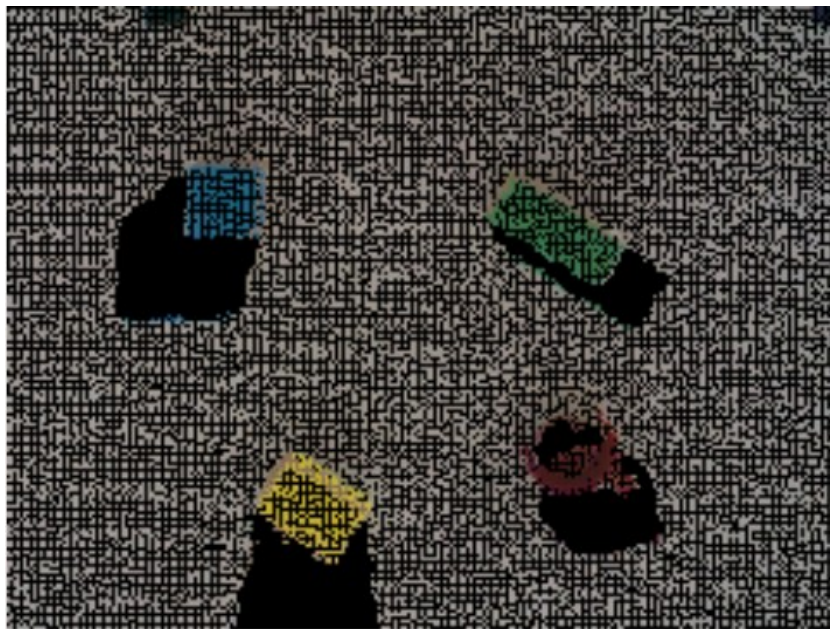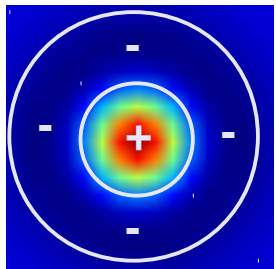
# Robotic Scenario: Challenges



- real sensory input
- moving sensors
- limited field of view
- 3D space
- dynamic scenes
- multiple behaviors
- computational constraints
- and many more ...

# Building Blocks

# Saliency



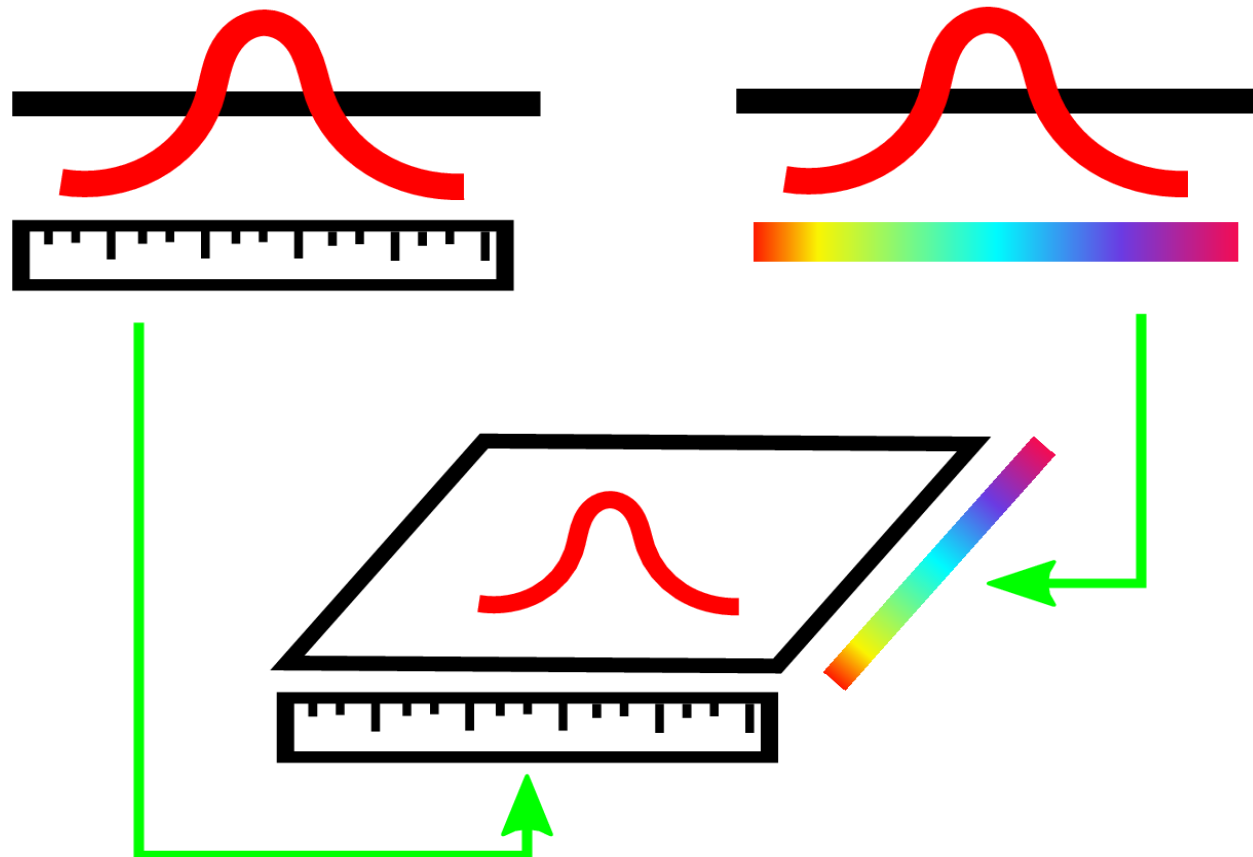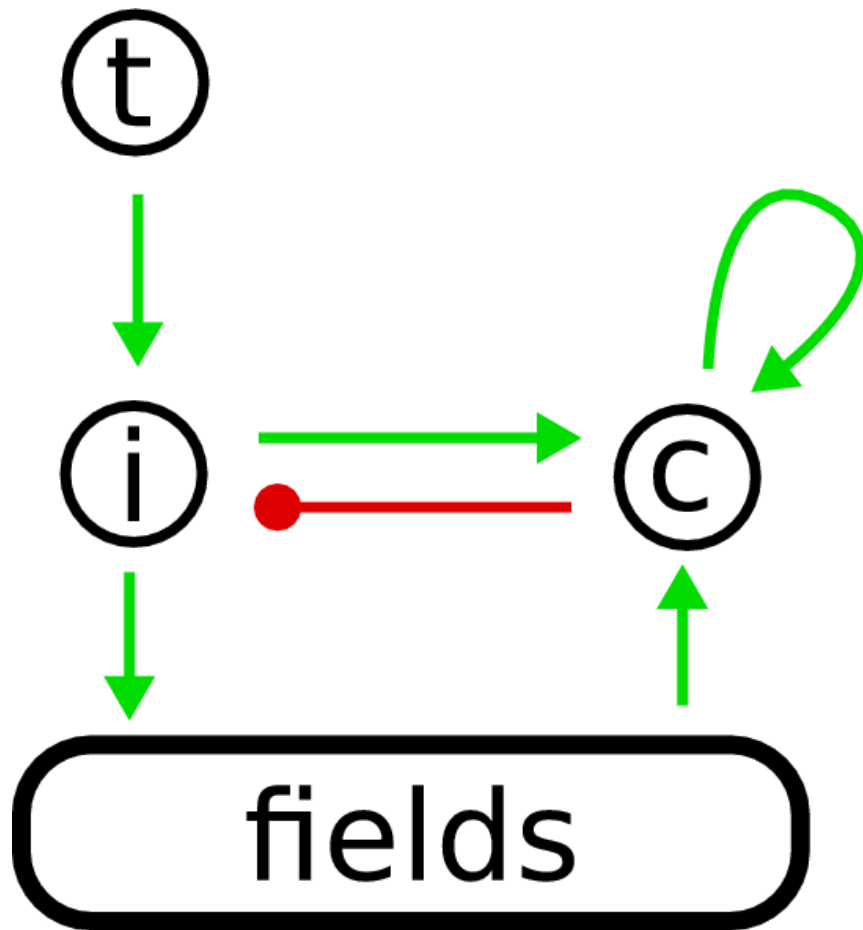Itti, Koch, Niebur (1998)

# Saliency



- on-/off-center responses

- uniform regions result in zero responses

- objects fitting into on-center region produce non-zero responses

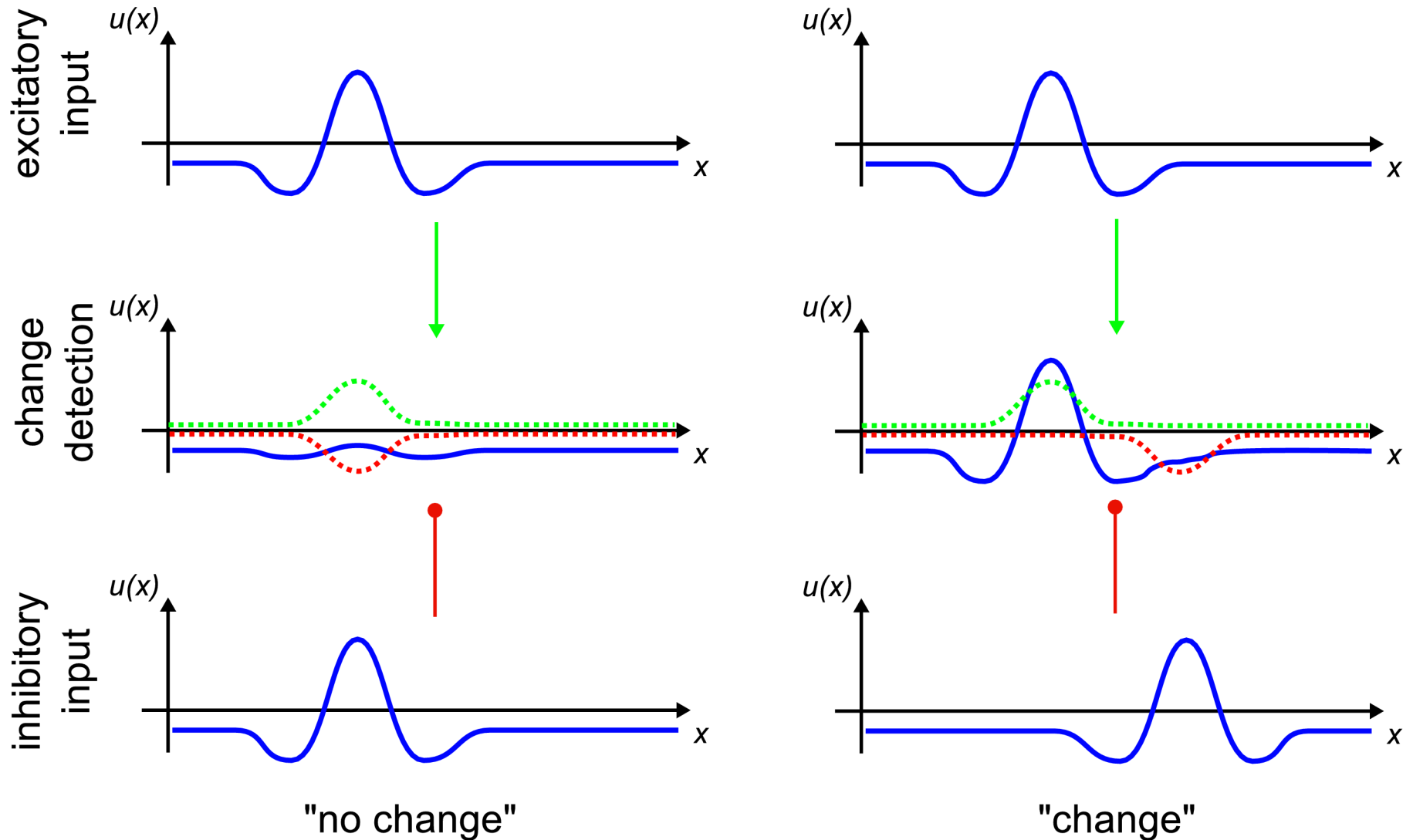- these lead to detection decision in fields

# Space-Feature Links



Feature estimates are linked to spatial positions.
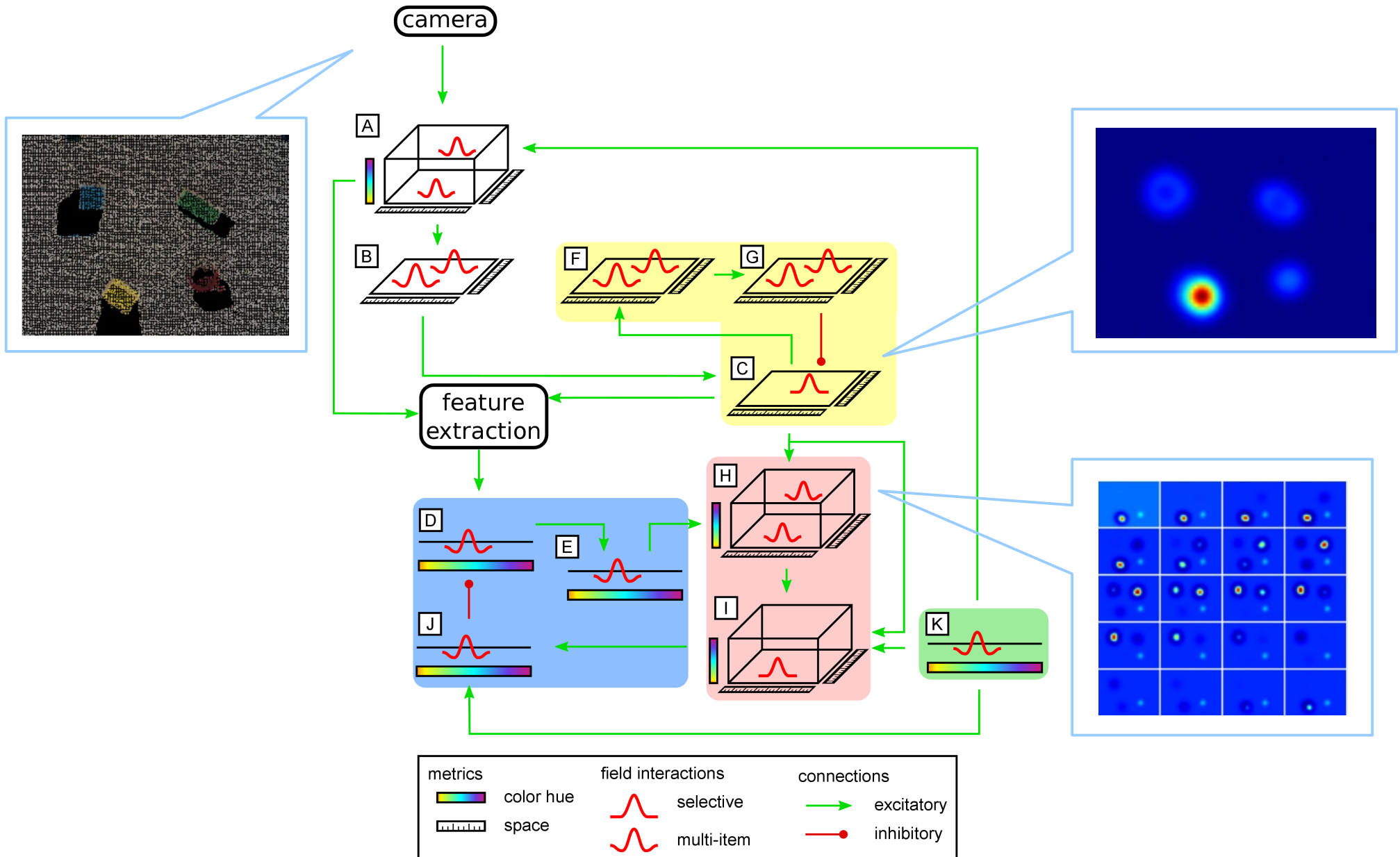
# Autonomy of Behaviors



- elementary cognitive units (ECU)

- intention node boosts fields

- CoS node detects completion of behavior in fields

- sequences or exclusions through precondition and suppression nodes
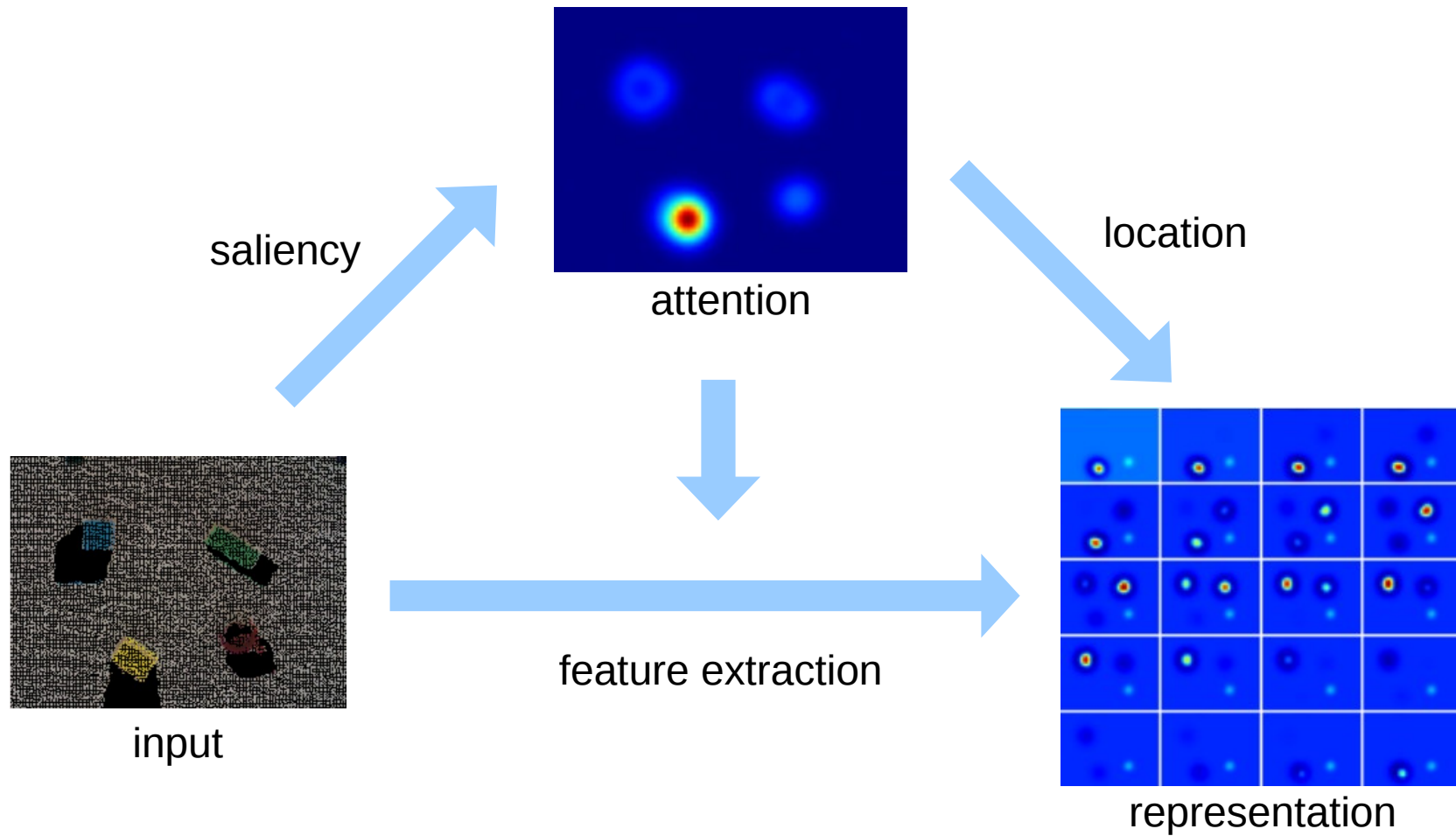
# Change Detection with Fields



excitatory input

change detection

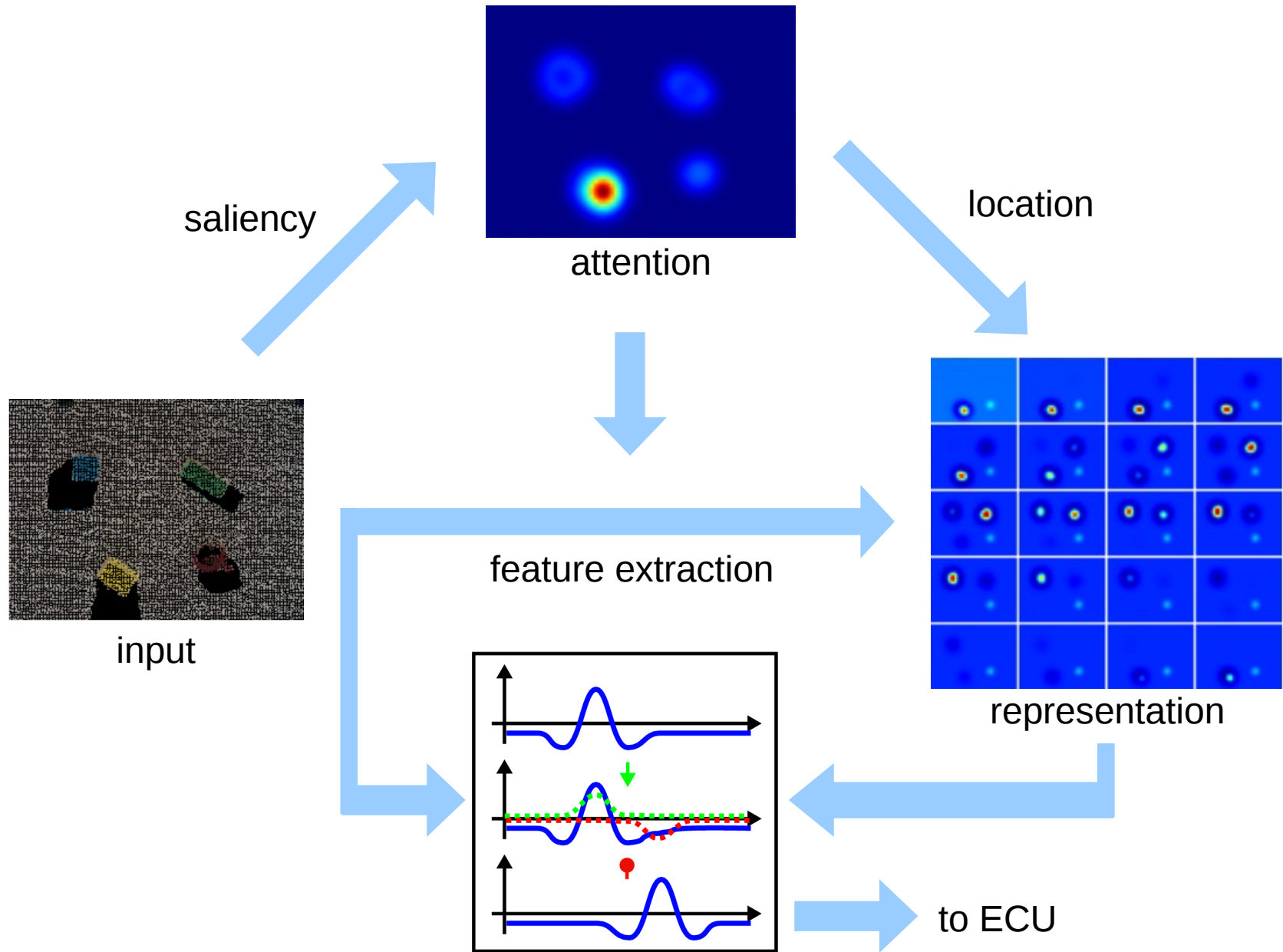inhibitory input

"no change"

"change"

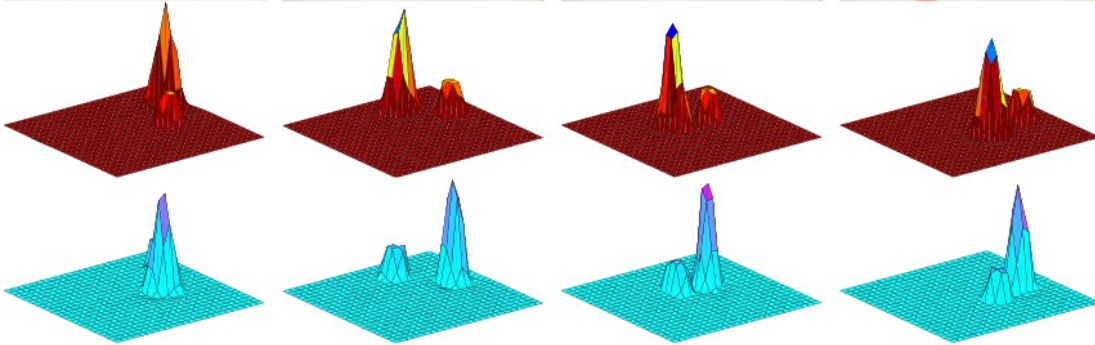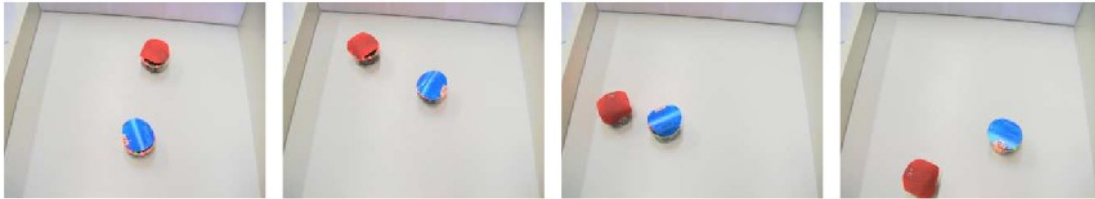# Elementary Behaviors of Scene Representation

# Architecture Overview



metrics

| | |
|---|---|
| color hue | |
| space | |

field interactions

selective

multi-item

connections

excitatory

inhibitory

# Exploration Behavior



saliency

attention

location

input

feature extraction

representation

# Autonomy of Exploration



saliency

attention

location

input

feature extraction

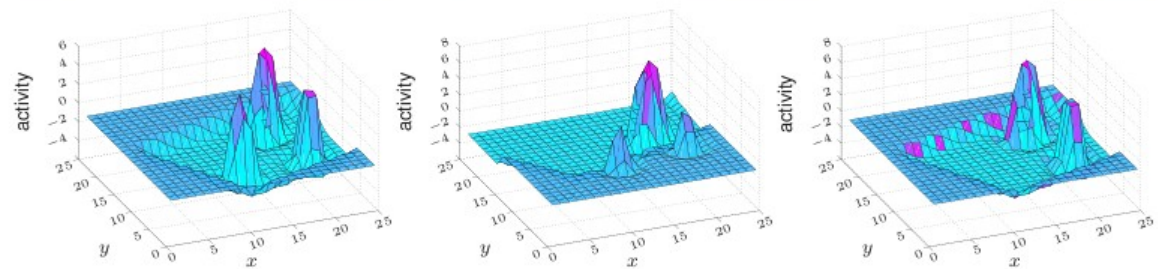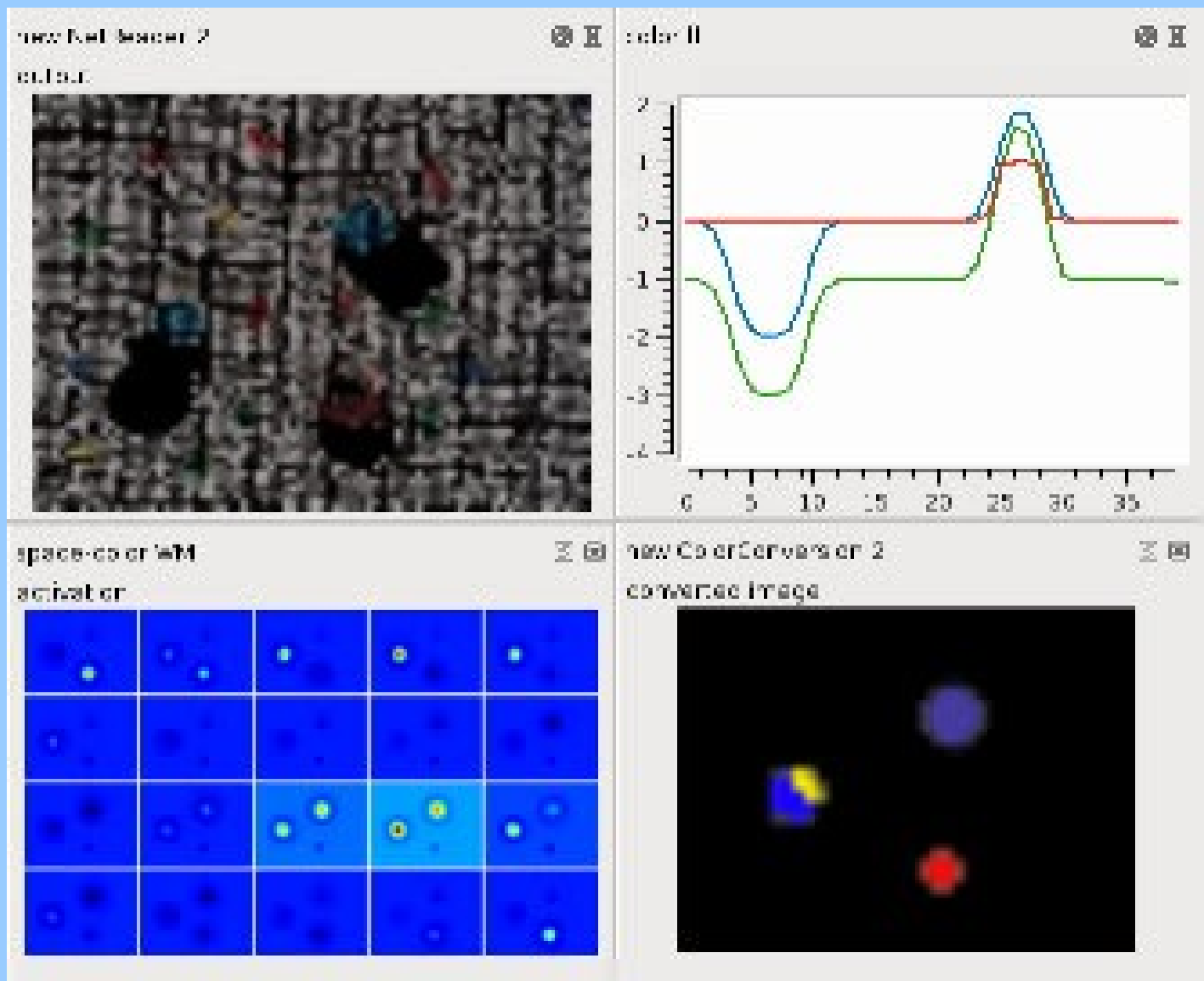representation

to ECU

**Video**

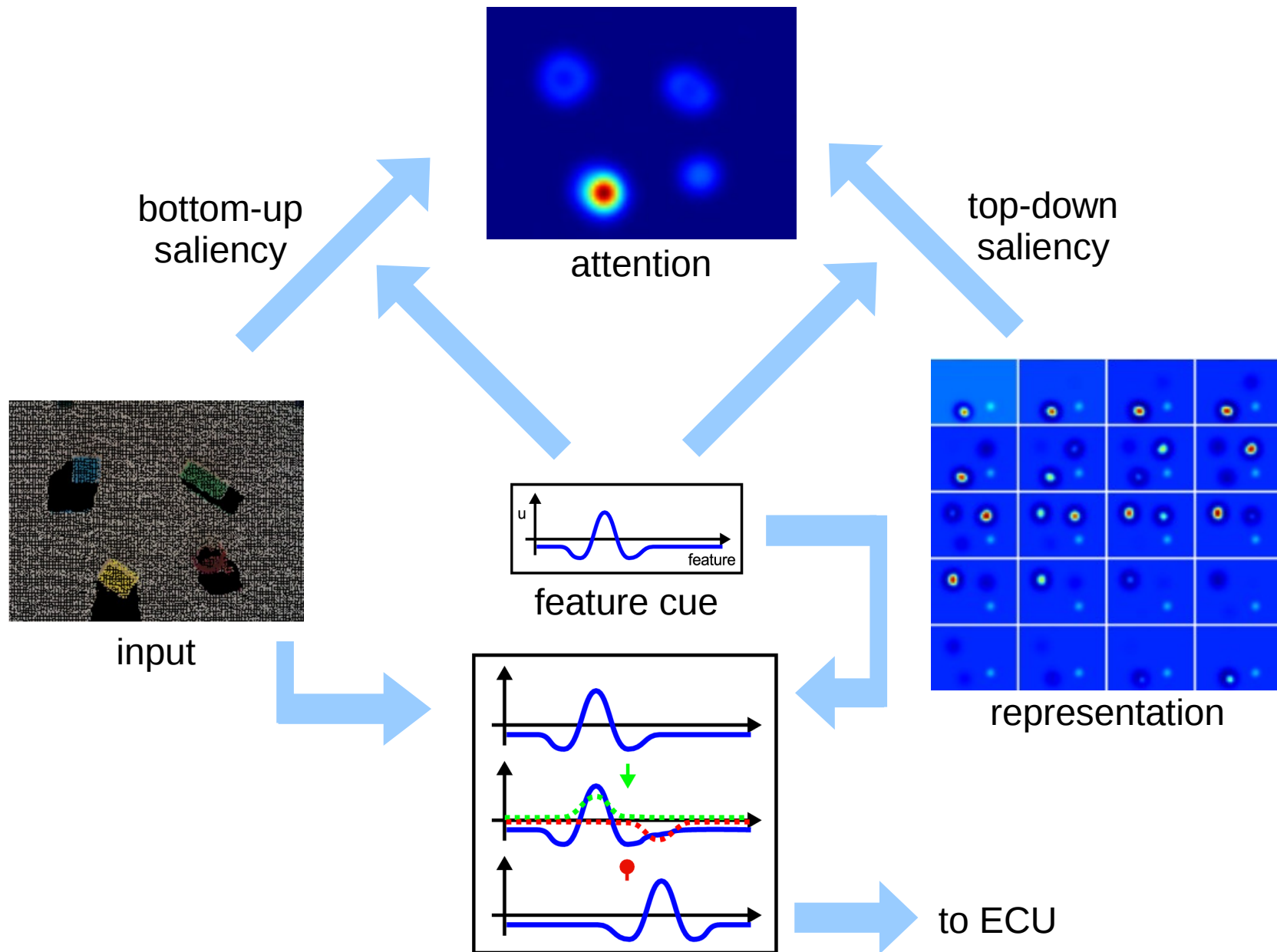# Maintenance Behavior



tracking

working memory
vs.
updating

**Video**

# Query Behavior



bottom-up
saliency

attention

top-down
saliency

input

feature cue

representation

# Autonomy of Query



bottom-up
saliency

attention

top-down
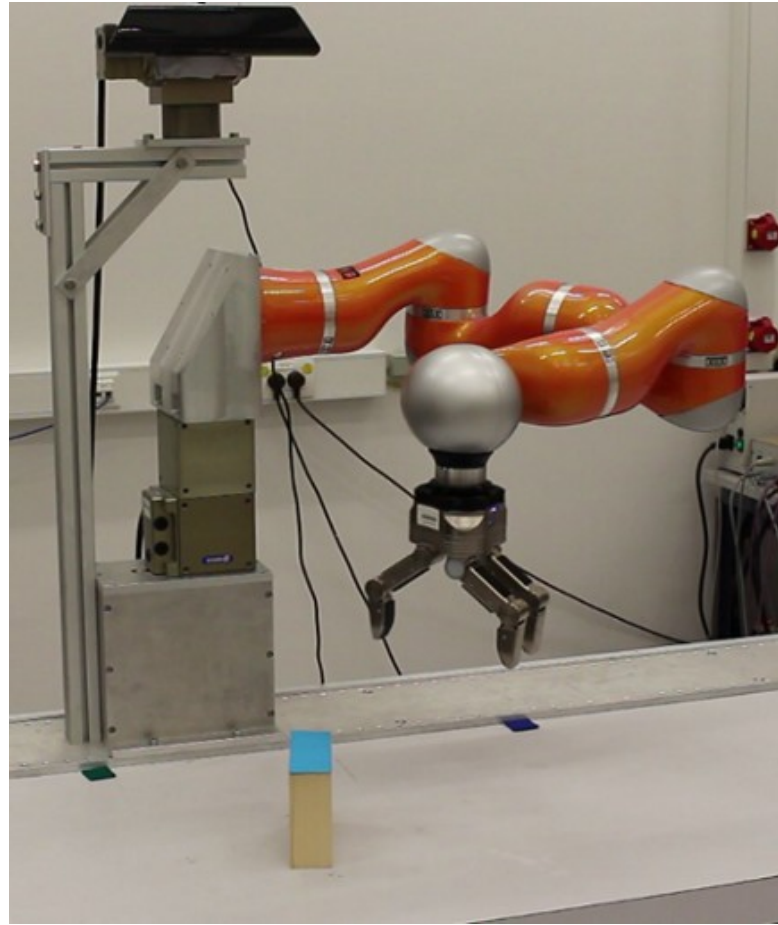saliency

input

feature cue

representation

to ECU

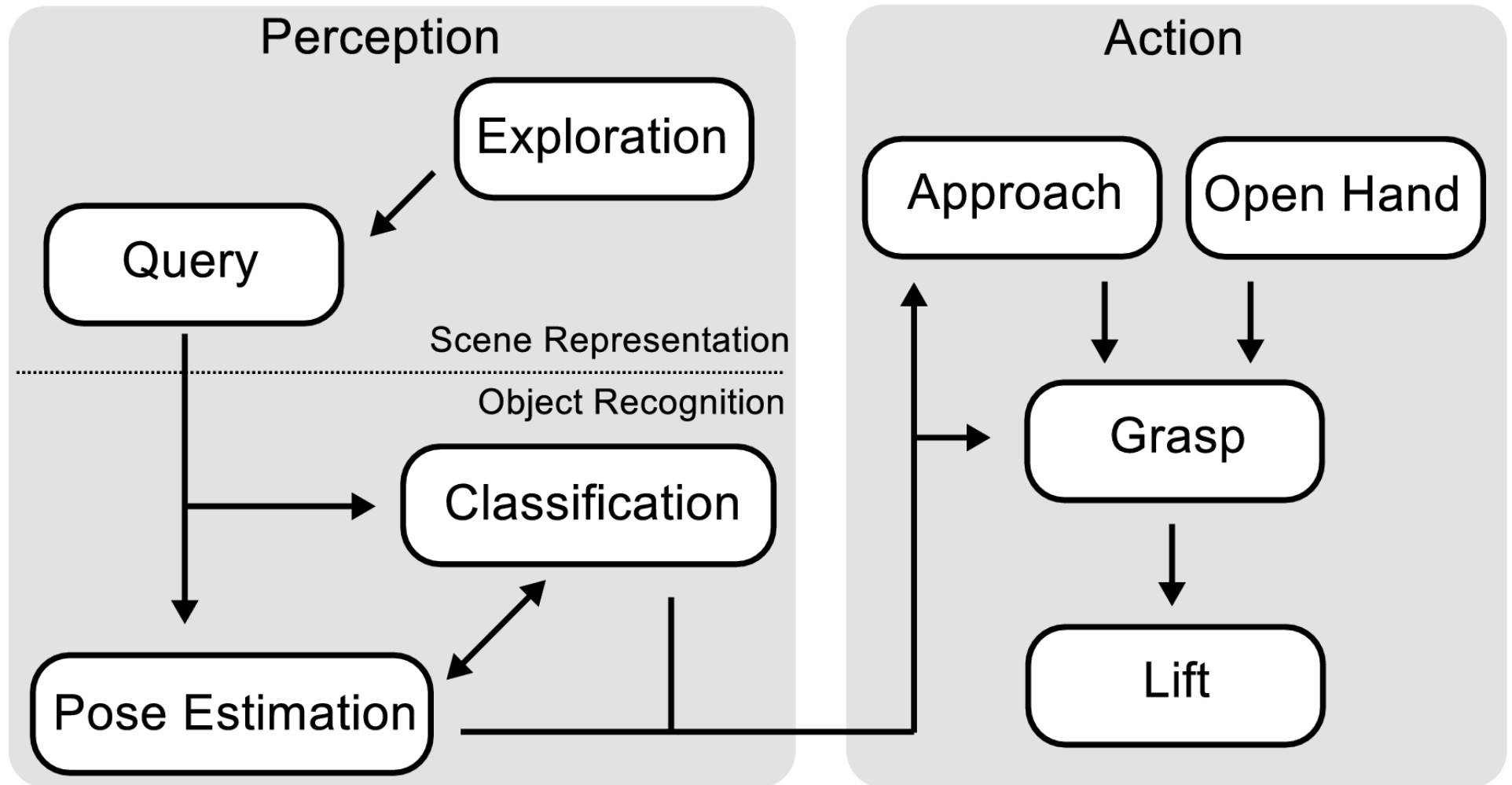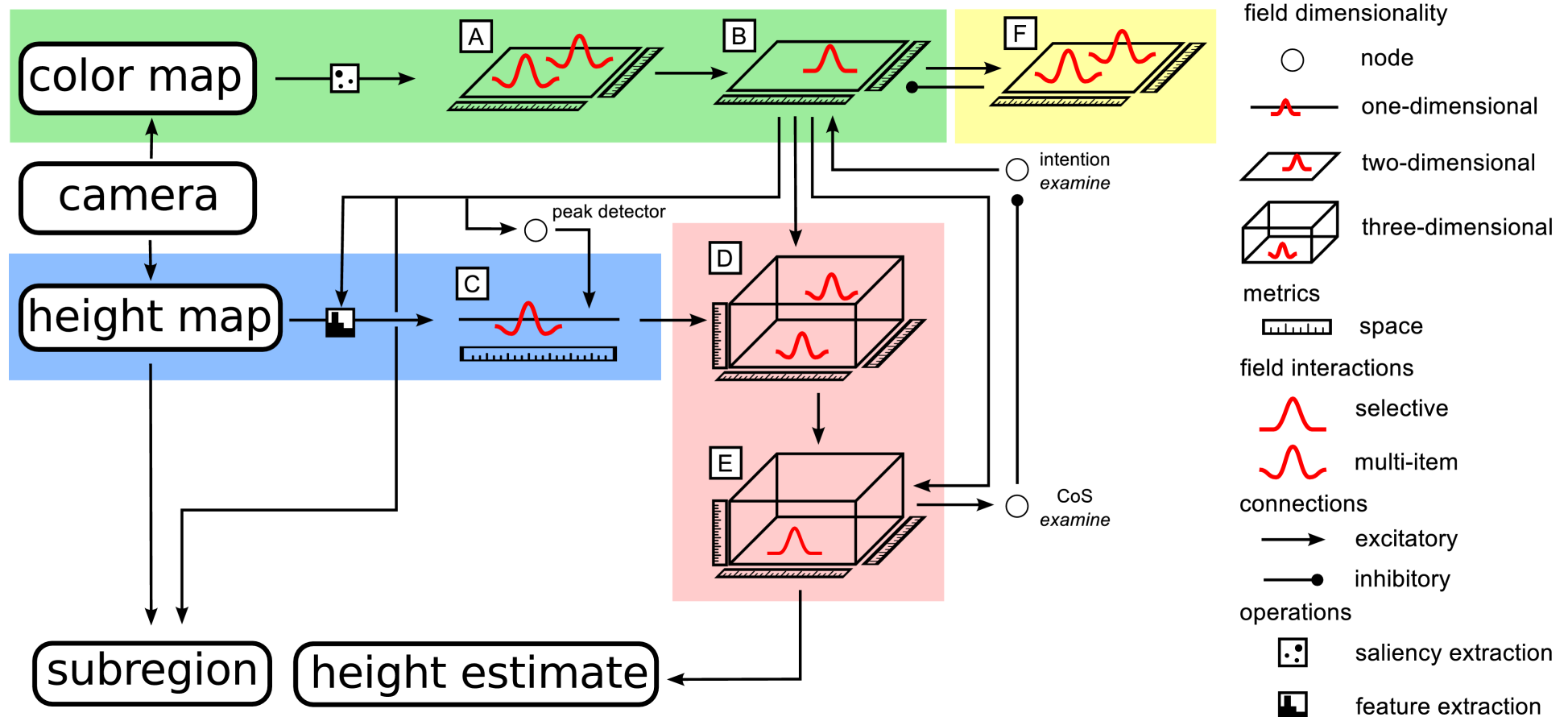# Reaching and Grasping

# Challenges



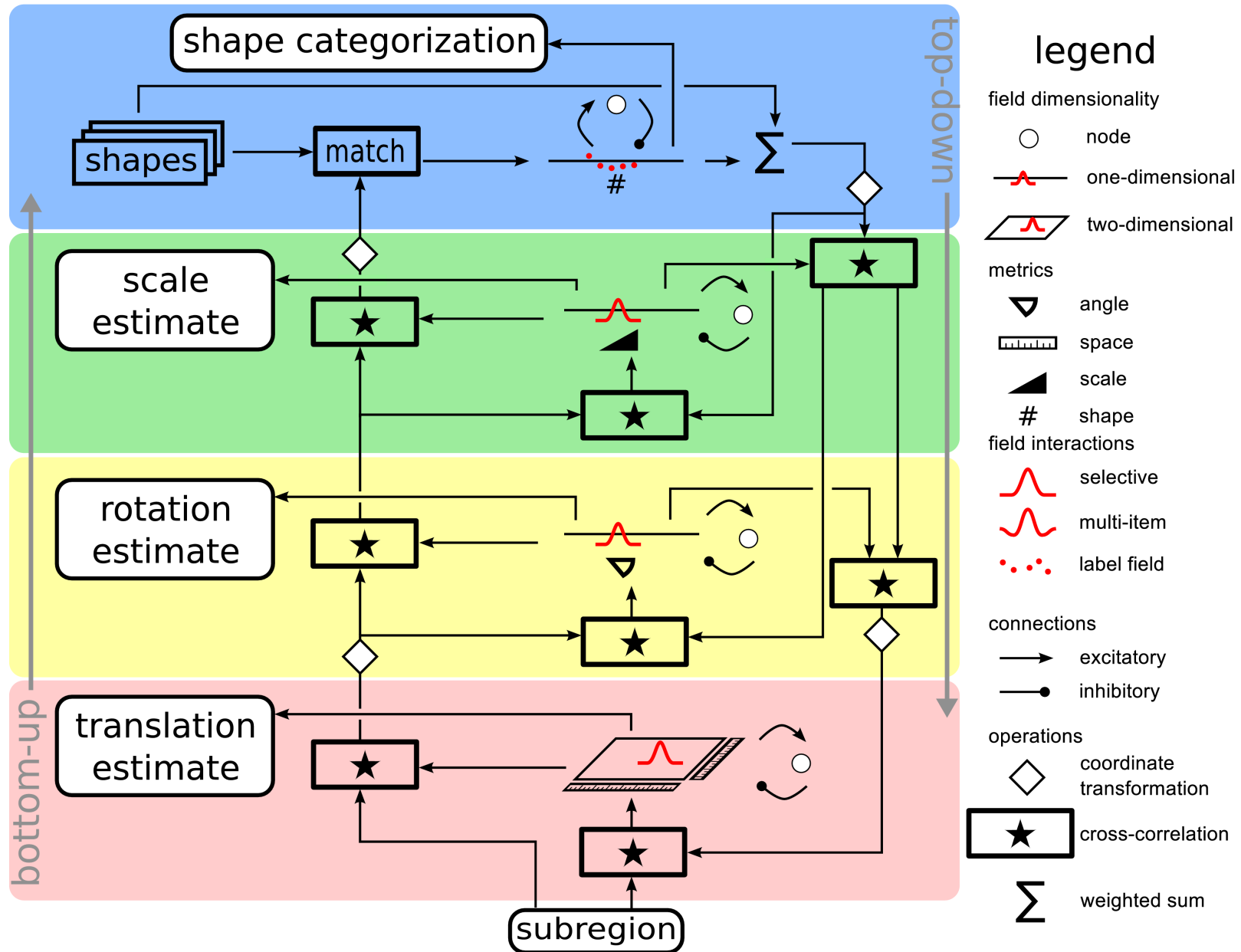focus attention

estimate shape

multiple behaviors

arm movement

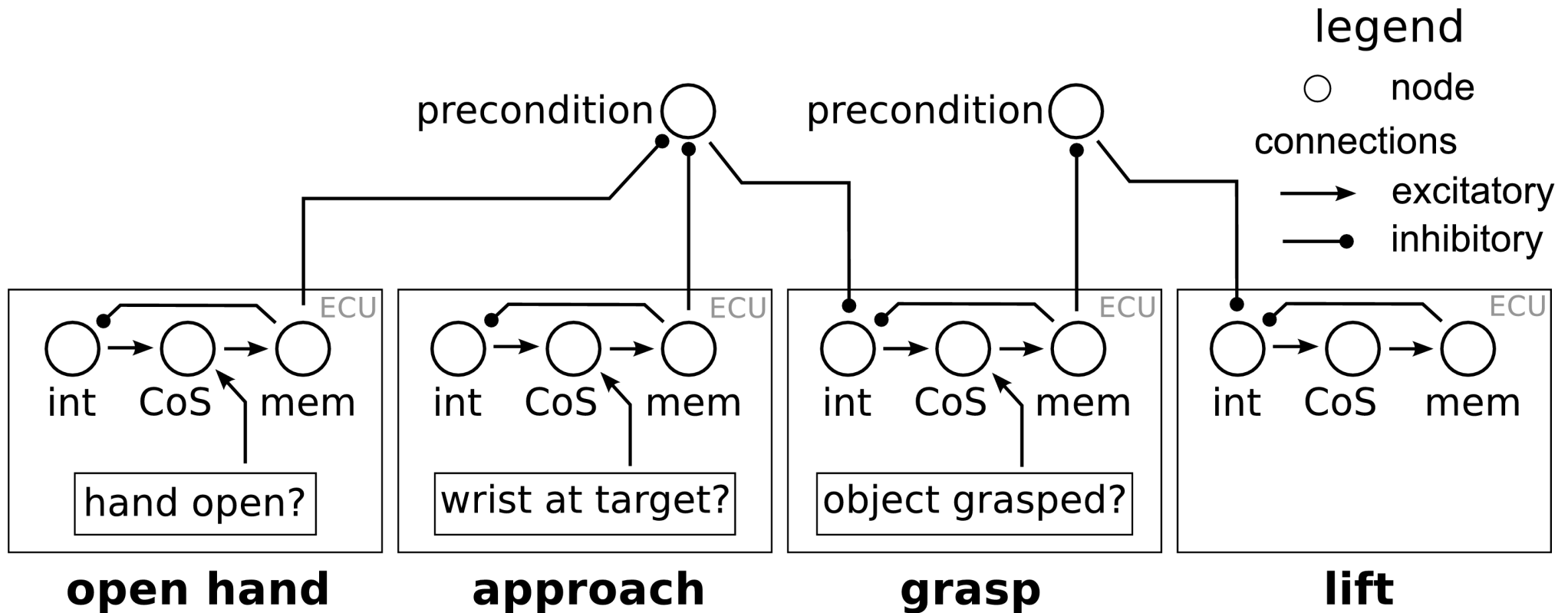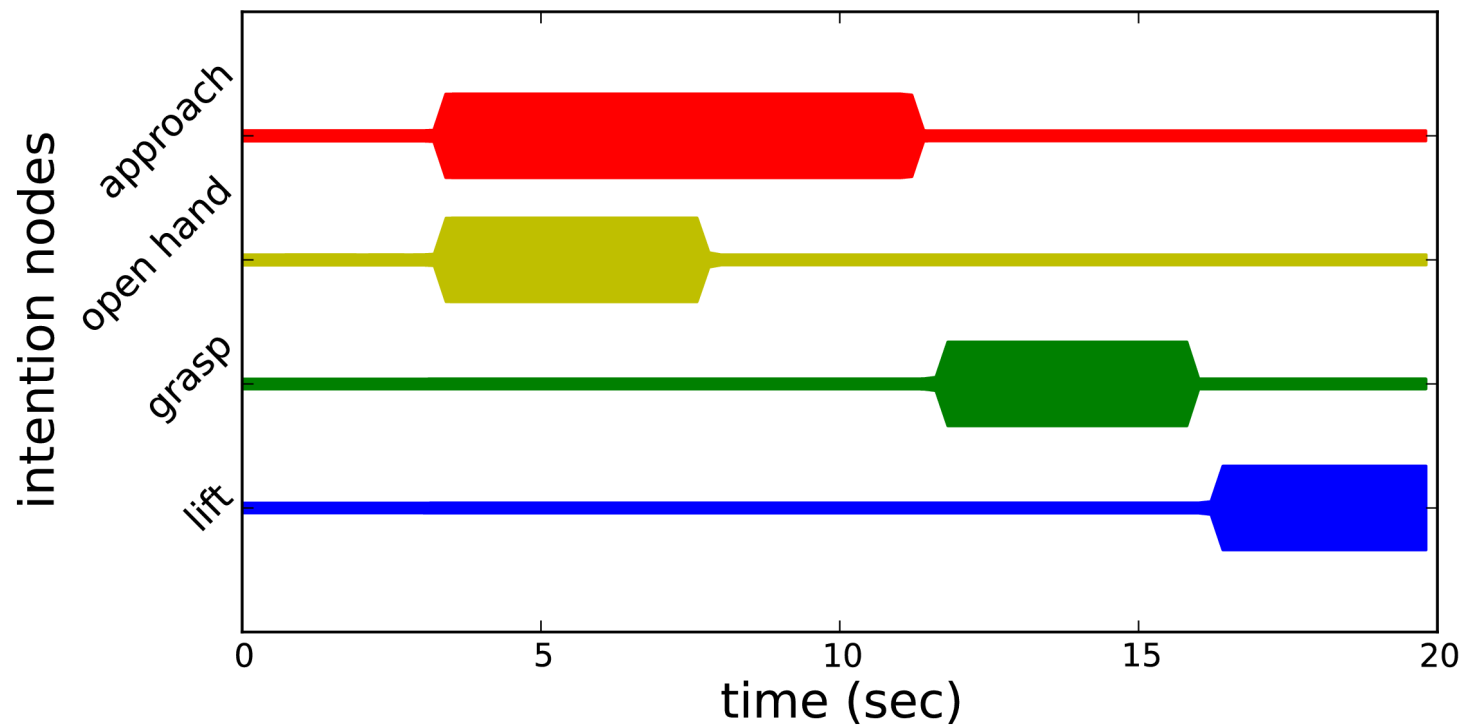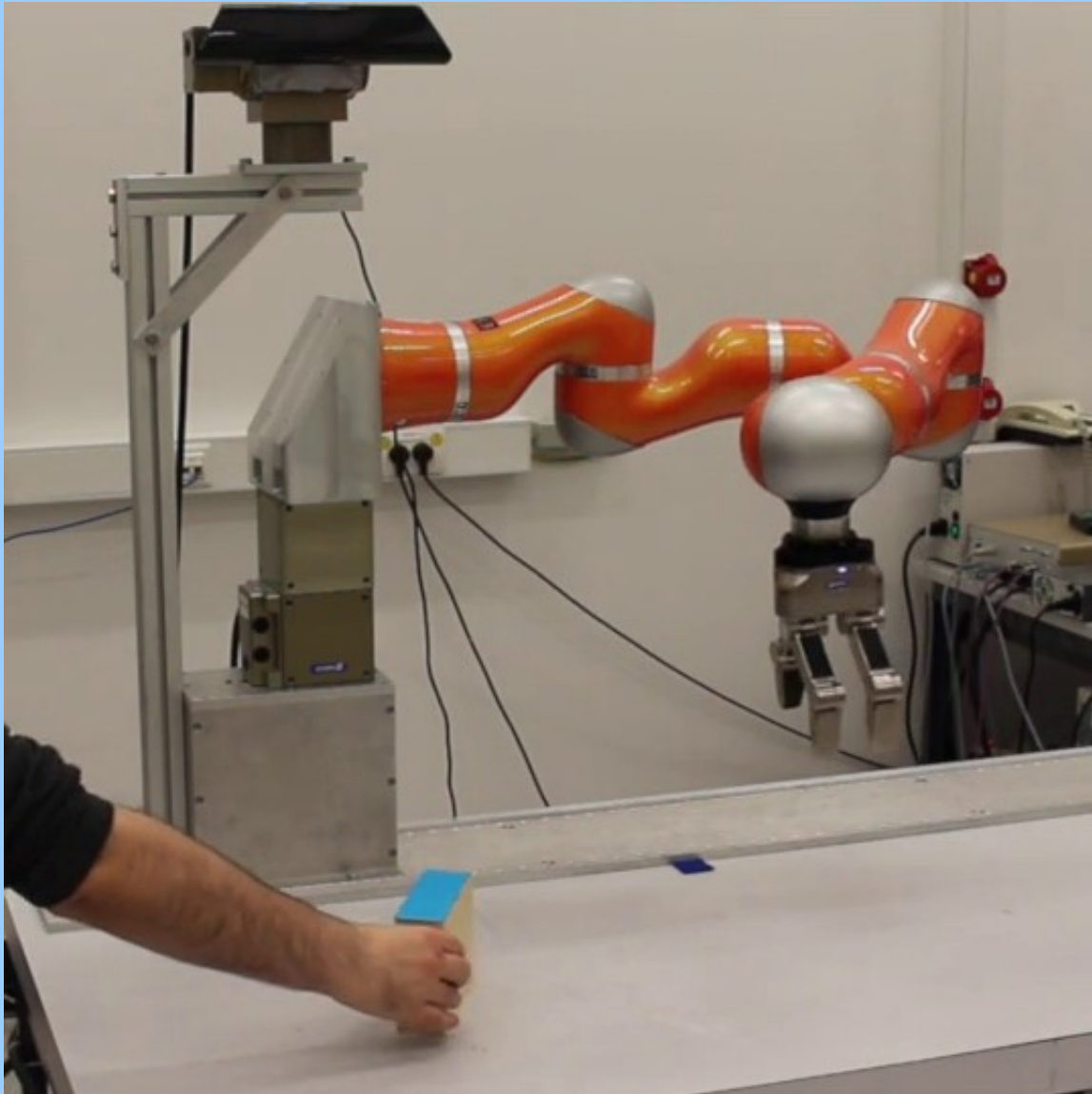# Behaviors Involved in Grasping

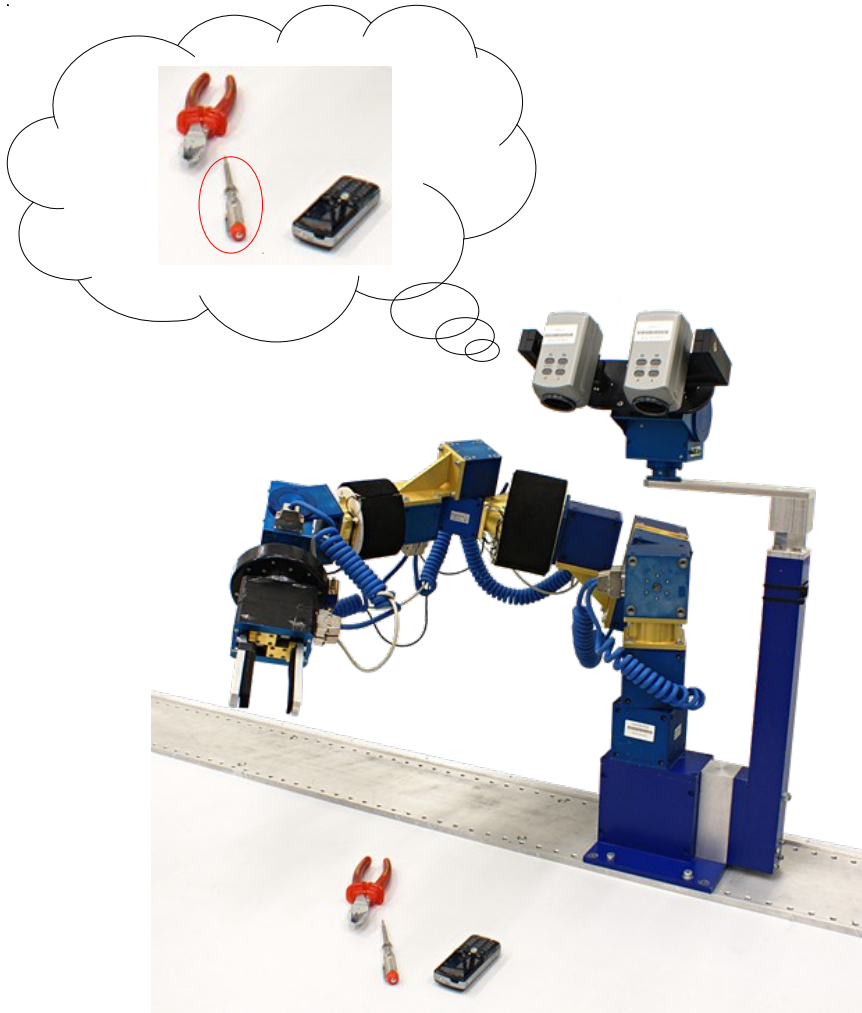# From Height Map to Grasp Parameters

# Grasp Execution

**Video**

# Take-home Message



- exploration, maintenance and query are the core behaviors of scene representation

- change detection is a driving force for autonomy

- integration with other DFT architectures yields complex behaviors such as grasping

- integration is facilitated by DFT framework

# Thanks for your attention!