



....

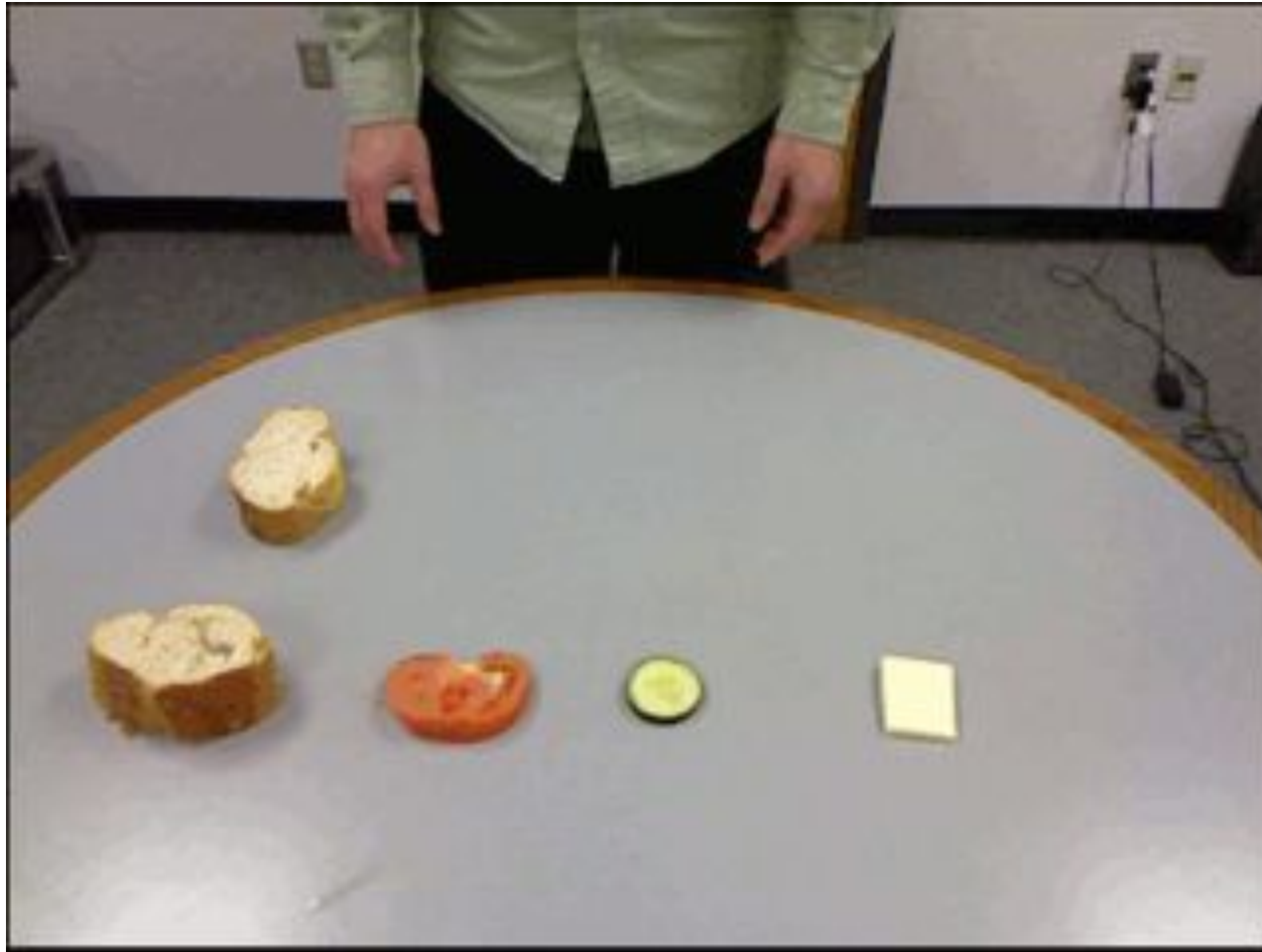
# Degrees of freedom problem



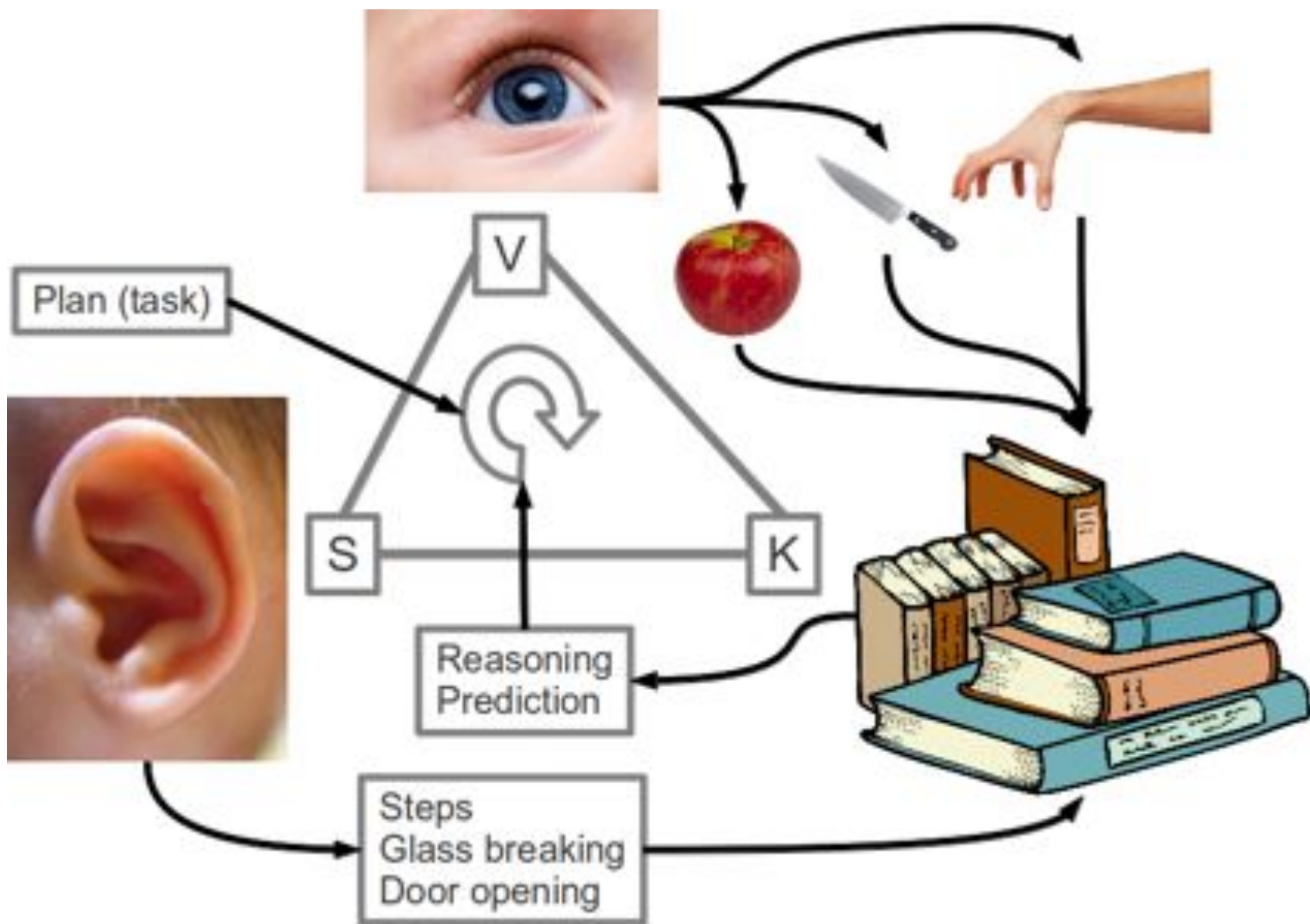
# Perception-Action integration



# Action Sequencing



# Learning



# Actions contain

Objects, tools, hands, movements

# Plan of the presentation

1. Attention, finding/recognizing objects using contours.
2. Recognizing goals (action consequences).
3. Parsing action sequences using the manipulation grammar.



# The Bottom up Attention Mechanism



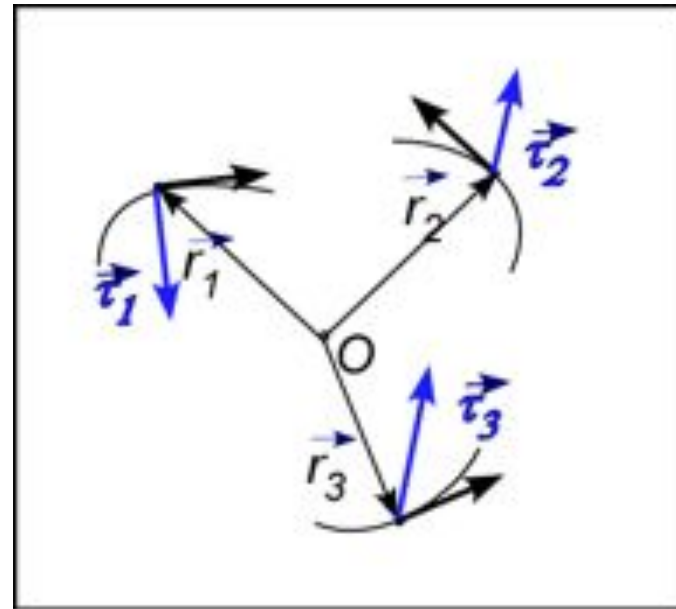
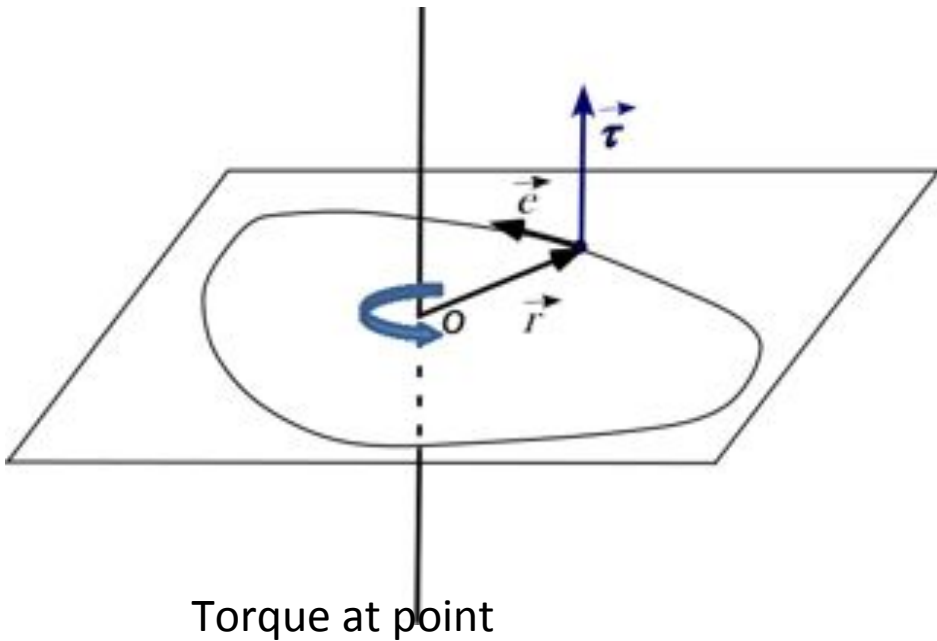
# A Motivating Example



*Where are my scissors?*

# The torque

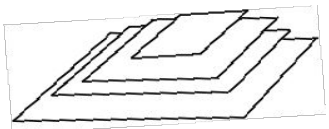
- encodes enclosed, connected regions



Value of torque :  $\vec{r} \times \vec{c}$



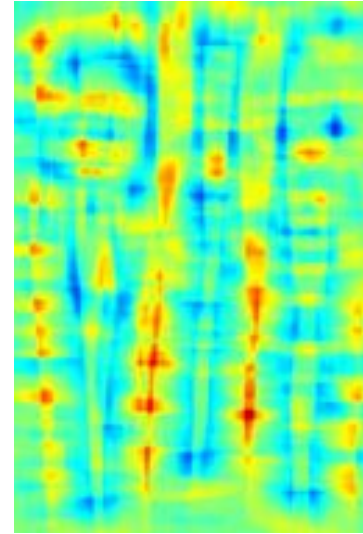
# Torque map



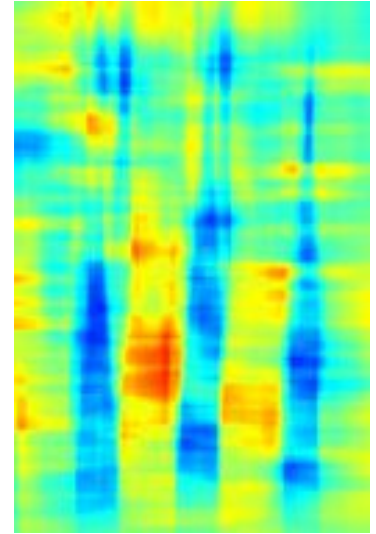
Torque at multiple Window sizes



size = 5 x 5

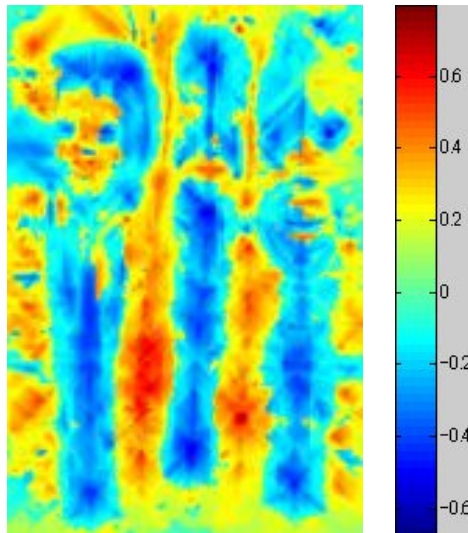


21 x 21



45 x 45

Combined Torque map



Extrema over scale



Canny edges



Use torque map to strengthen edges

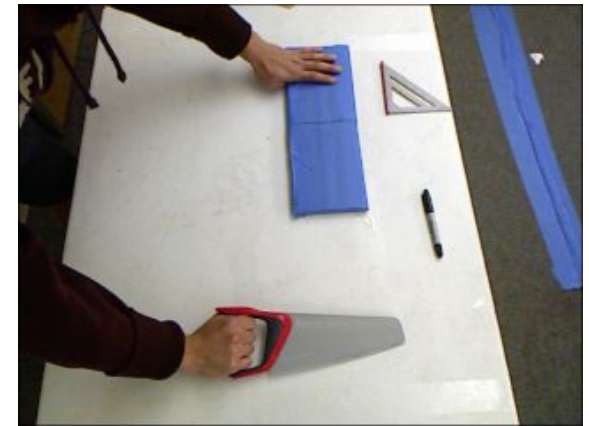
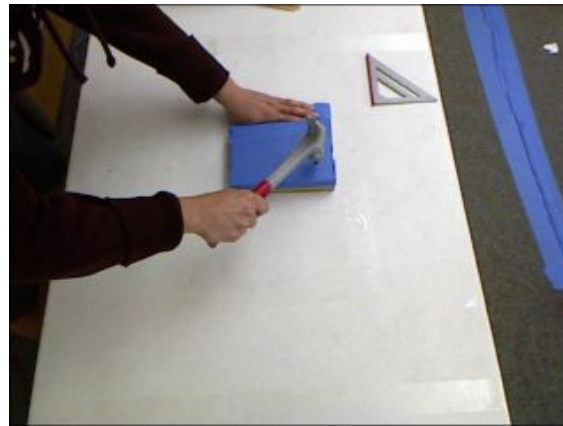
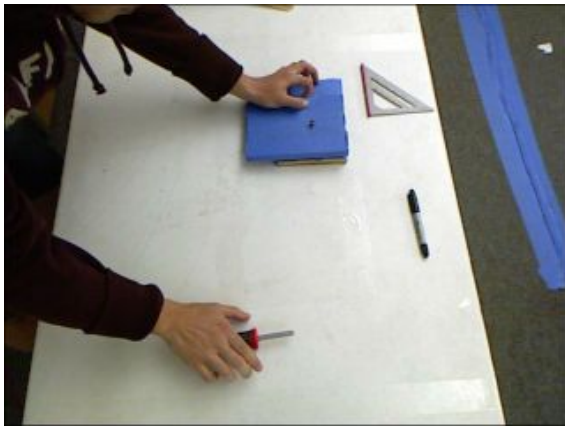
# Object Recognition in clutter/occlusion

Typical scenes from manipulation/search tasks:

1) Clutter with partial occlusions and viewpoint changes



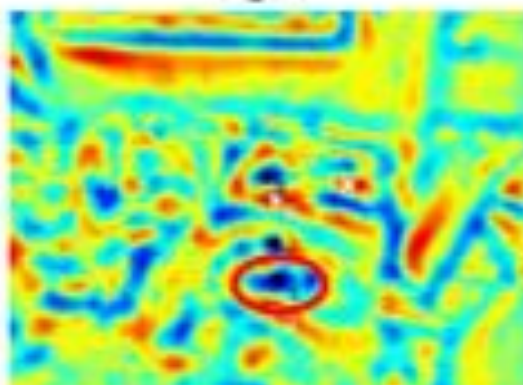
2) Occlusions due to hands, objects in various orientations



Input RGB + (Depth)



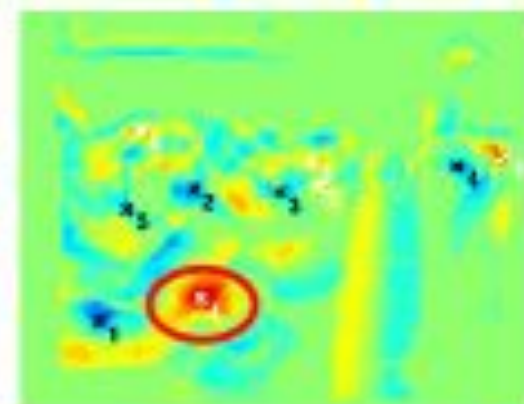
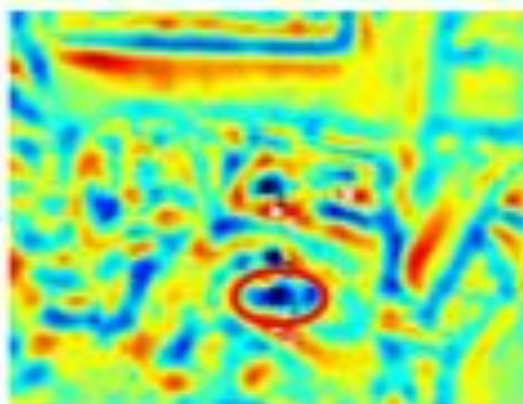
Flashlight



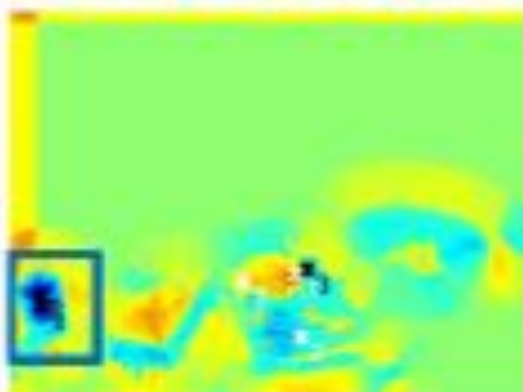
$\tau_p''$



Cap

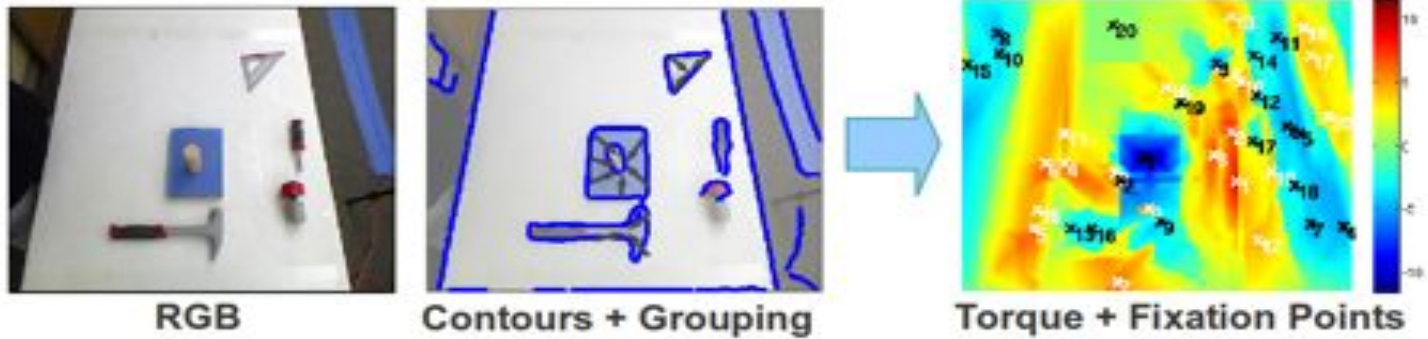


Tissue Box

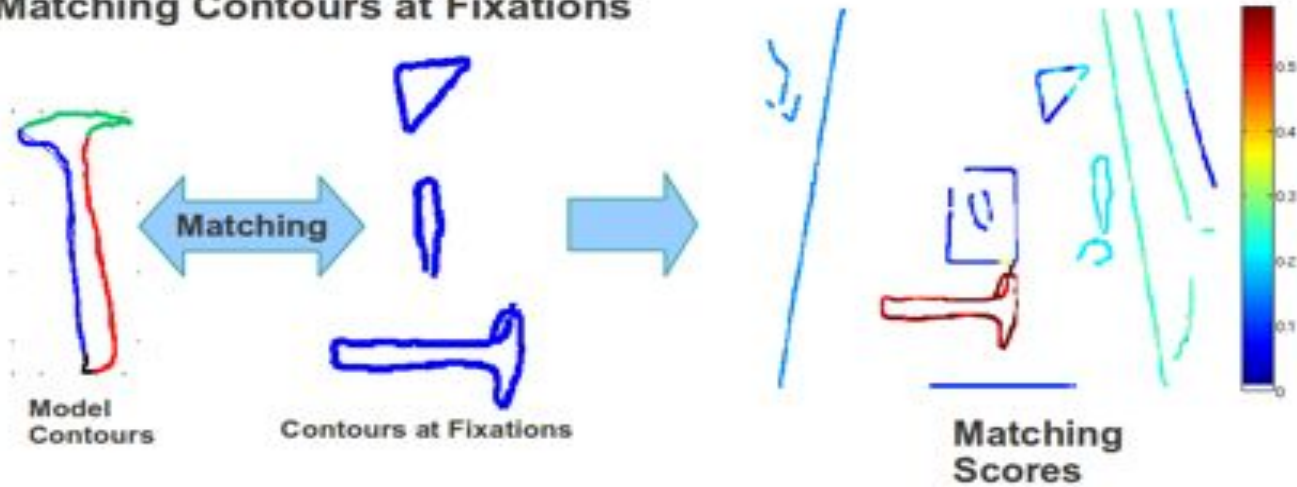


# Solution (3 steps)

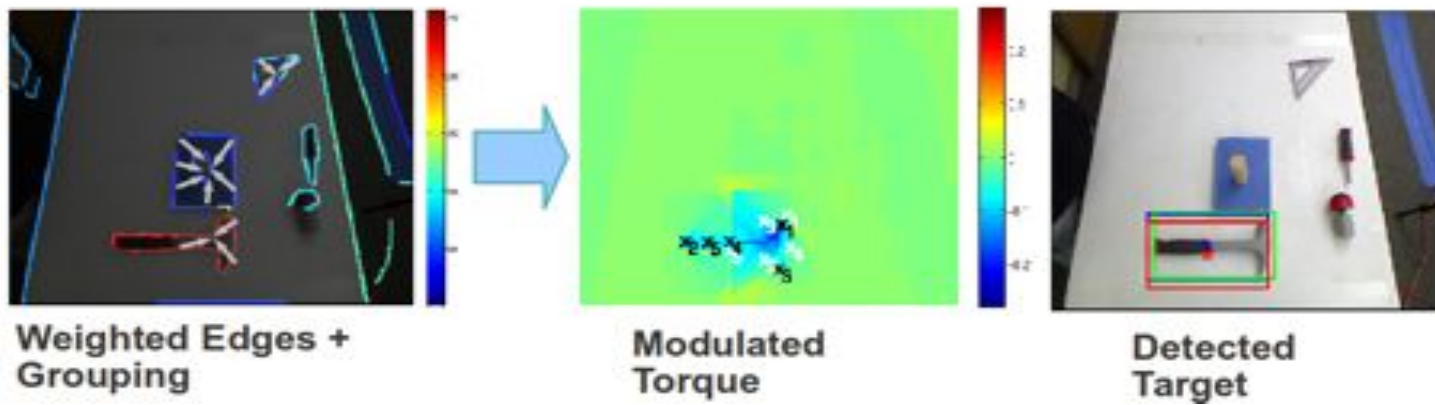
## (I) Determine Fixations



## (II) Matching Contours at Fixations

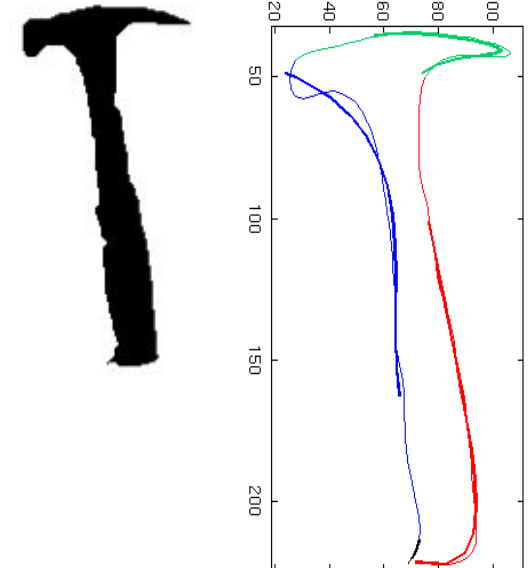
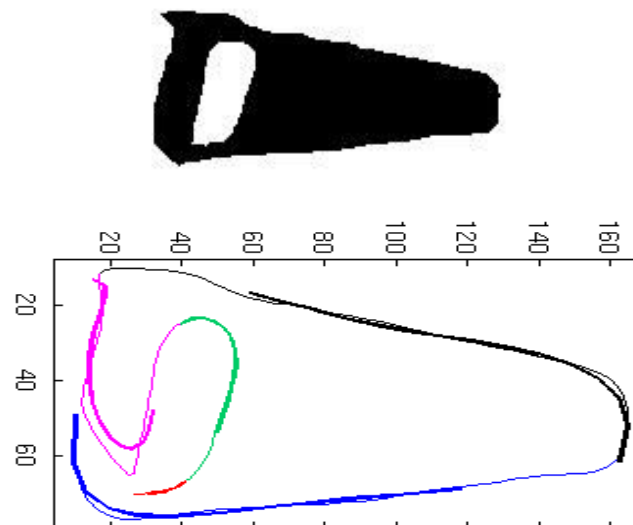
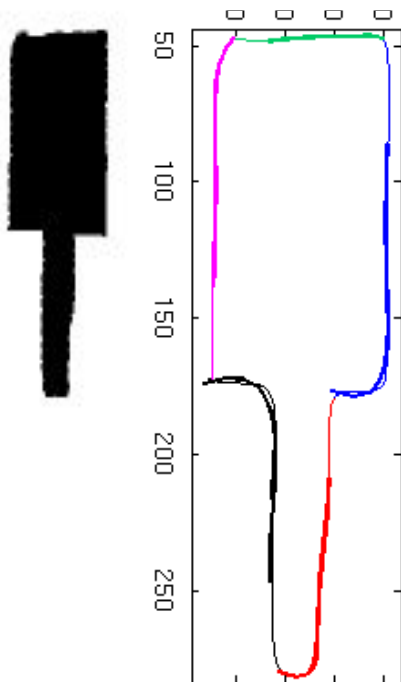


## (III) Target Object Detection





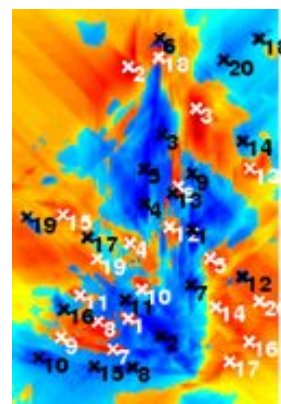
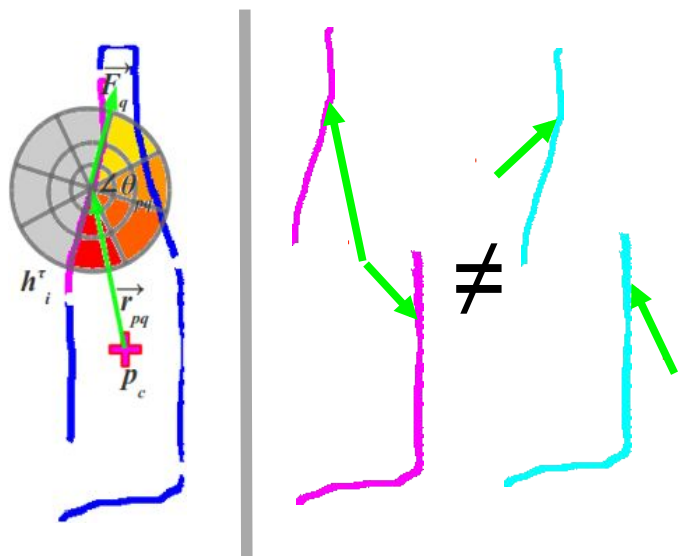
- **Learn** salient contours from annotated examples
- **Match** partial contour fragments to models
- **Reweigh** torque based on matching scores
- Learned contours examples:





# Modulate Torque for Object Similarity

Shape Context Matching with Torque Centers



Matching



*Shape-Context Only*



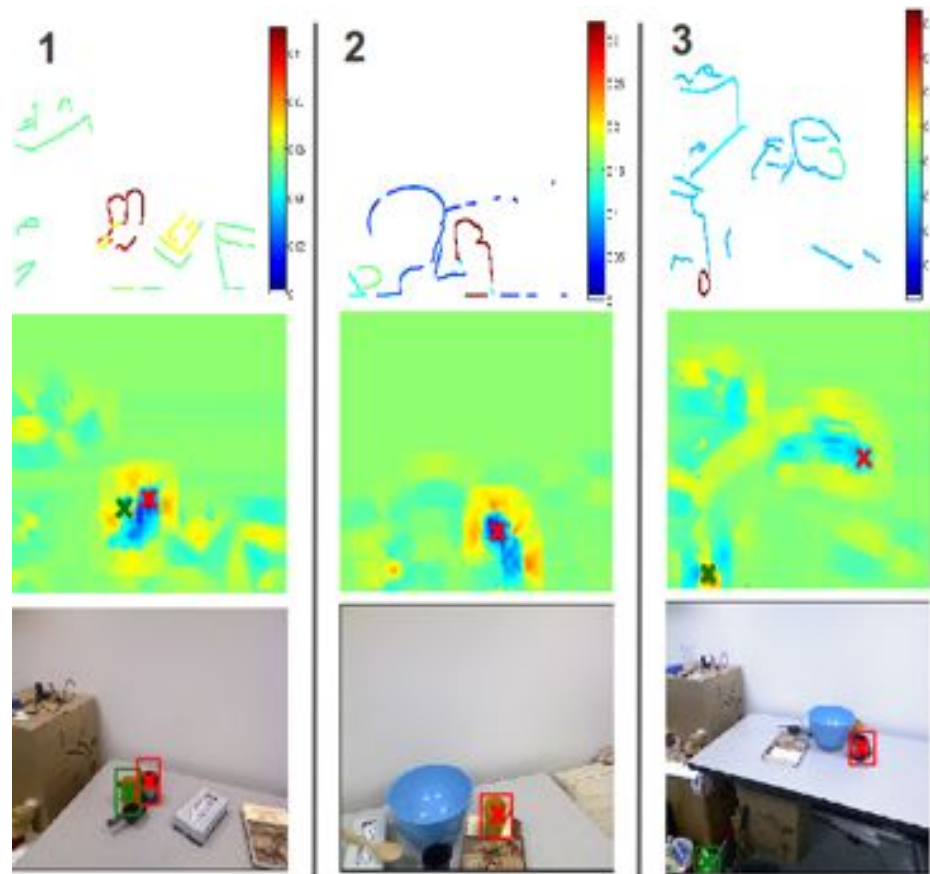
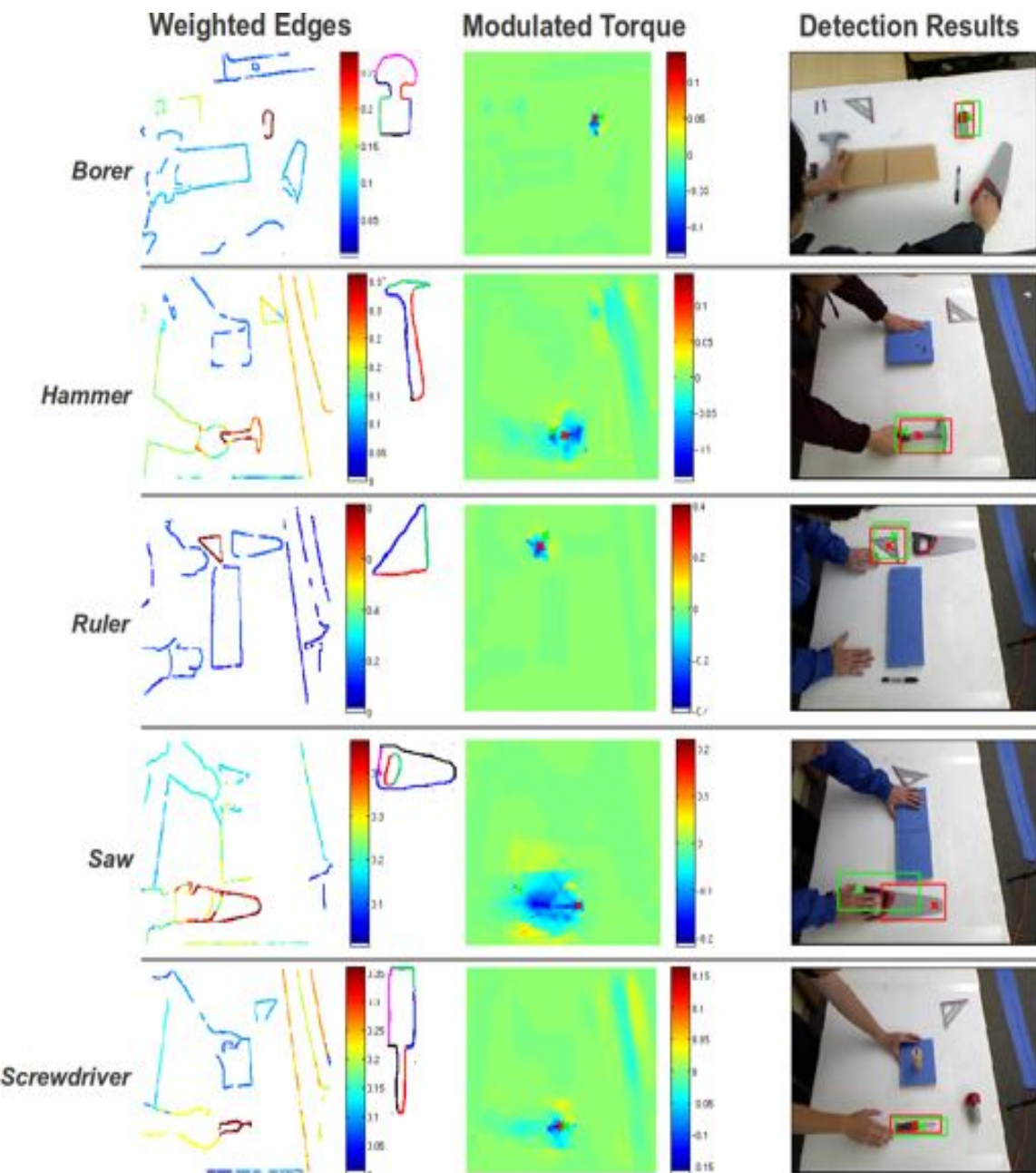
*Torque Shape-Context*

Edge Weights



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

# Results



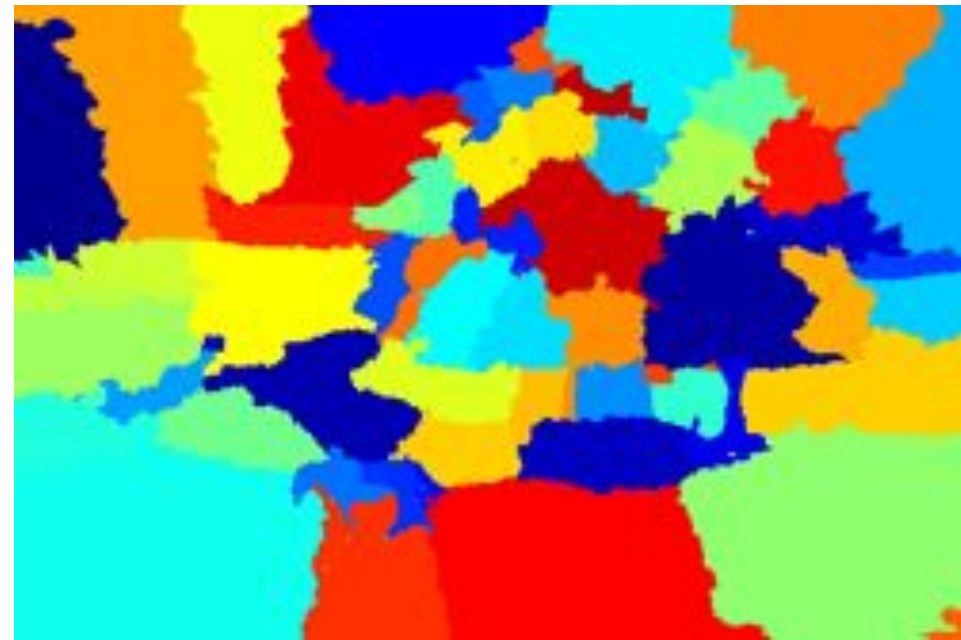
# Segment the scene



# Segmentation with Normalized Cut



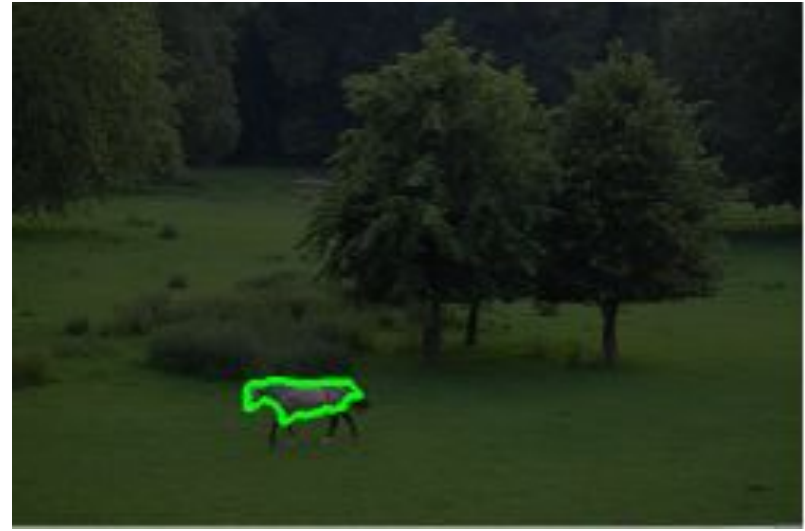
Number of regions = 10



Number of regions = 60

- Which is the appropriate segmentation?
  - ☐ Left (trees), Right (horse)

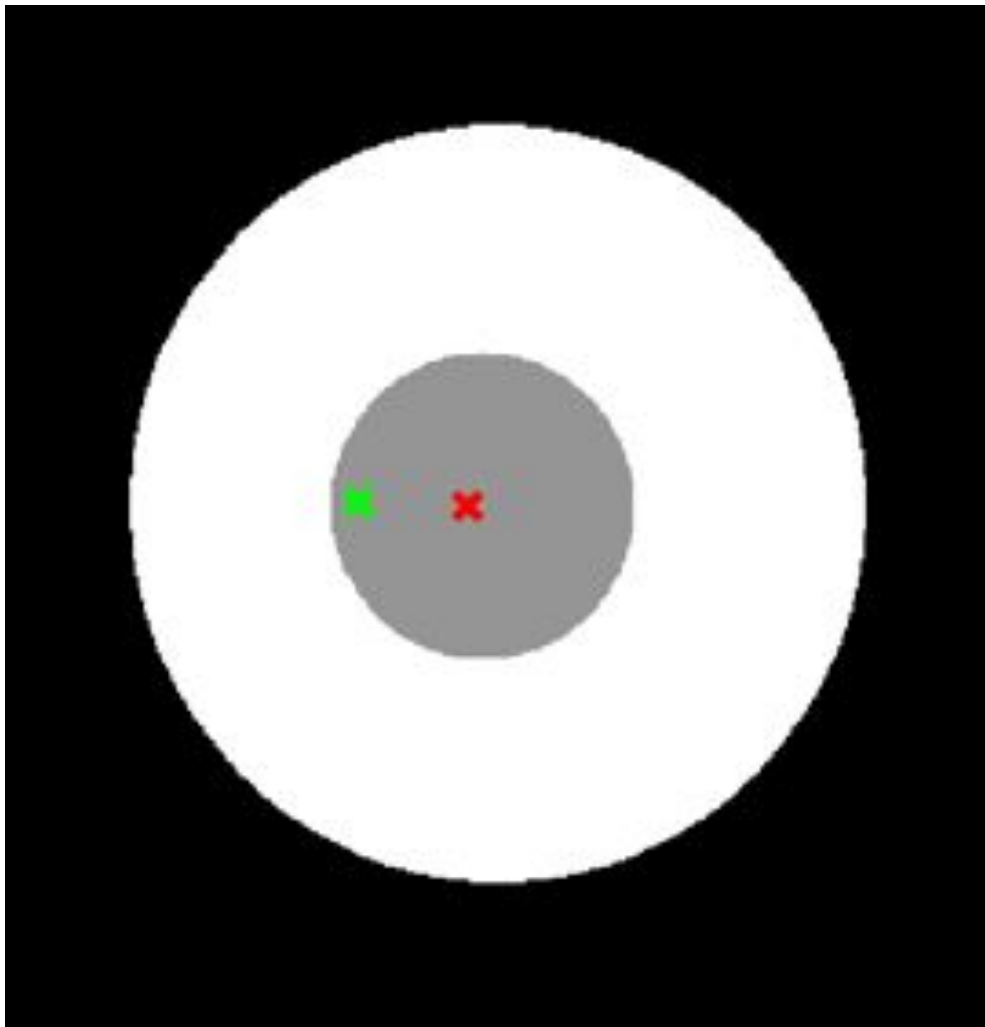
# Segment what you look at



•



•



Internal contour: 0.39; outside contour: 0.78

Internal contour length: 100 pixels

Outside contour length: 400 pixels

Cost of tracing each contour

$88 = (400 (1 - 0.78))$  and  $61 = (100 (1 - 0.39))$

Smaller (dimmer) contour costs less





Internal contour brightness: 0.39; outside contour: 0.78

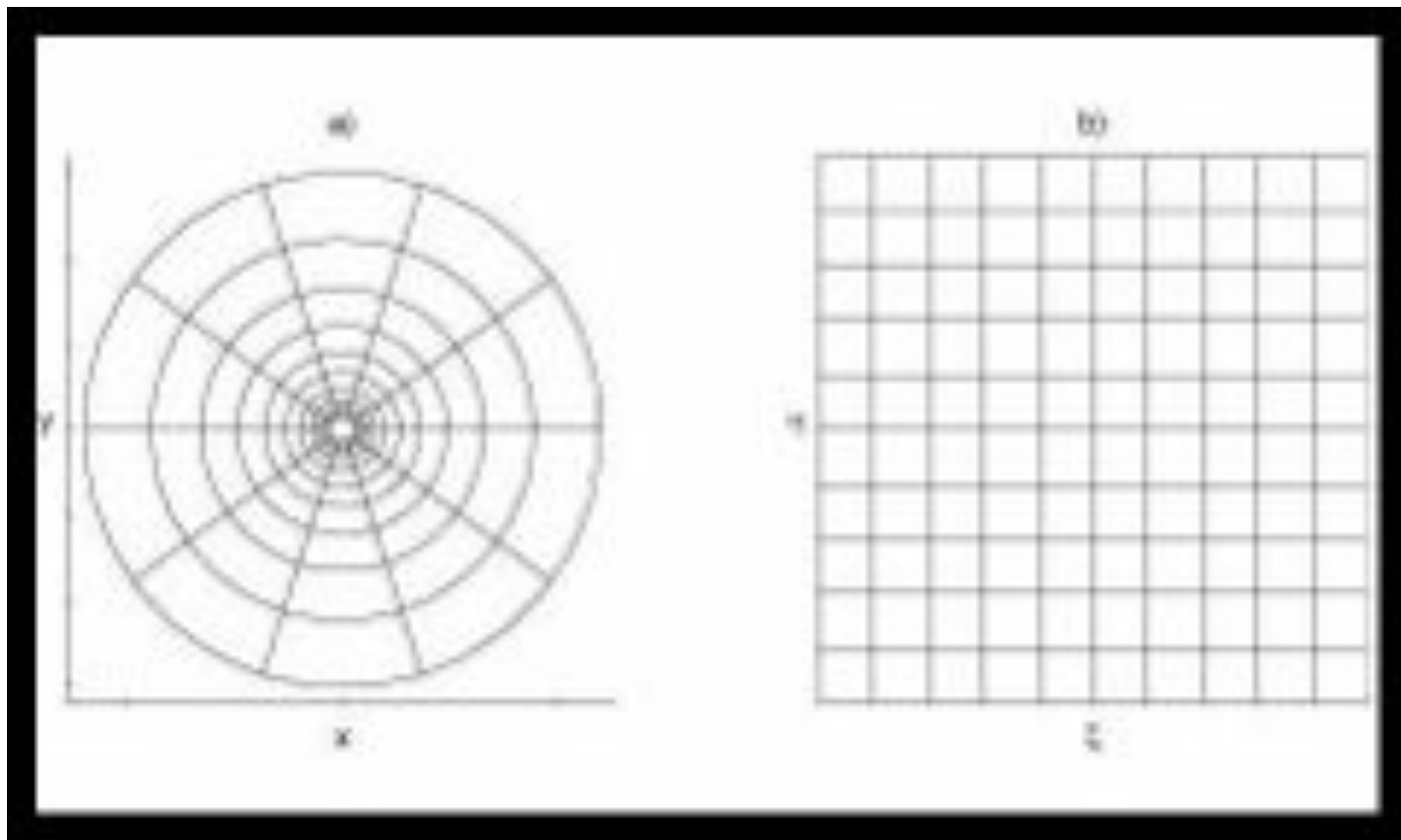
Internal contour length: 361 pixels

Outside contour length: 365 pixels

Cost of tracing each contour

$80.3 = 365 (1 - 0.78)$  and  $220.21 = 361 (1 - 0.39)$

Brighter contour costs less







Color and texture cues find all edge pixels



Motion/stereo and surface cues identifies object boundaries

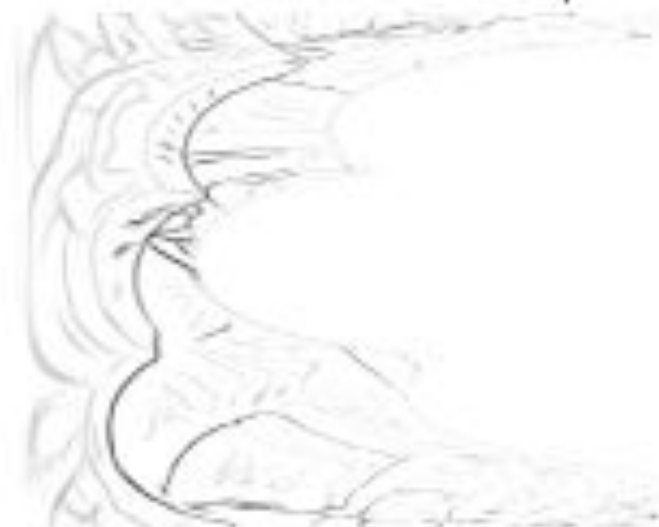


# Probabilistic boundary depth map

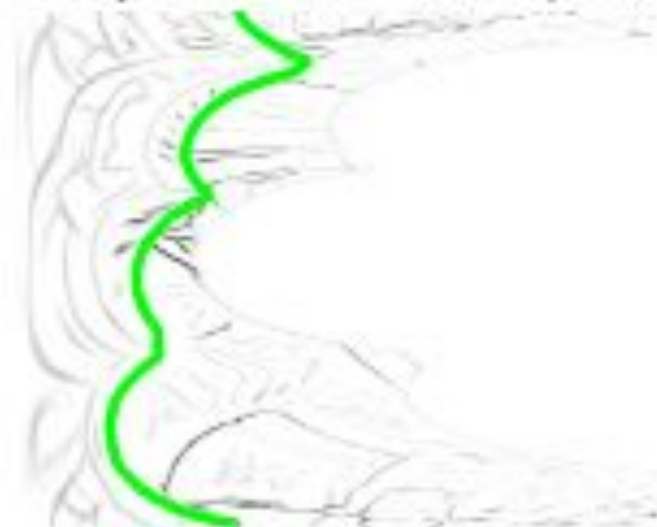


# Segmenting a fixated object

**Step 1:** Cartesian to polar with fixation as the pole

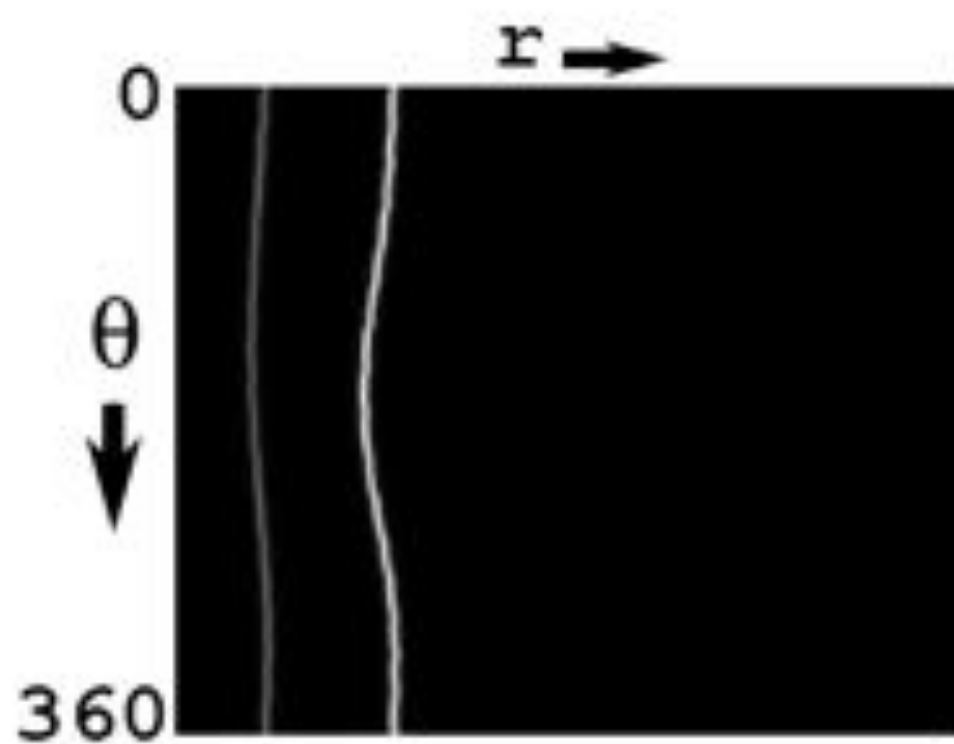
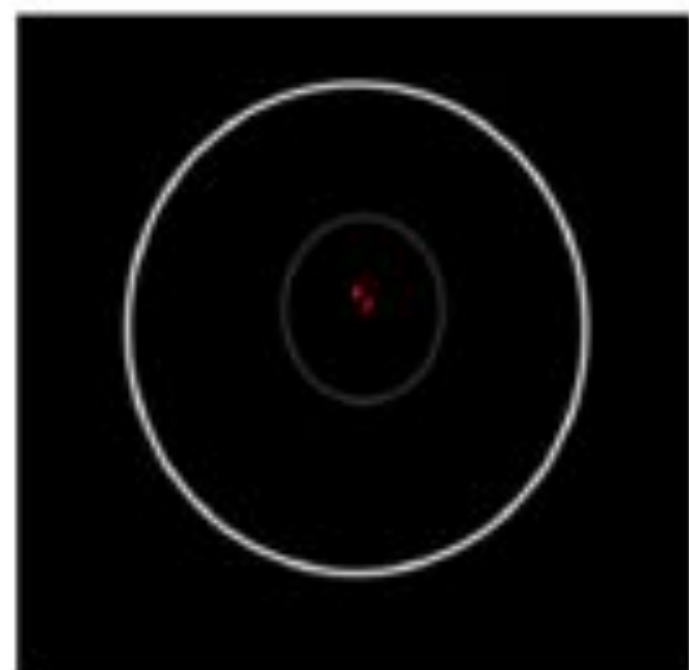


**Step 2:** Find the optimal cut through the polar map



## Step 1: Cartesian to polar conversion

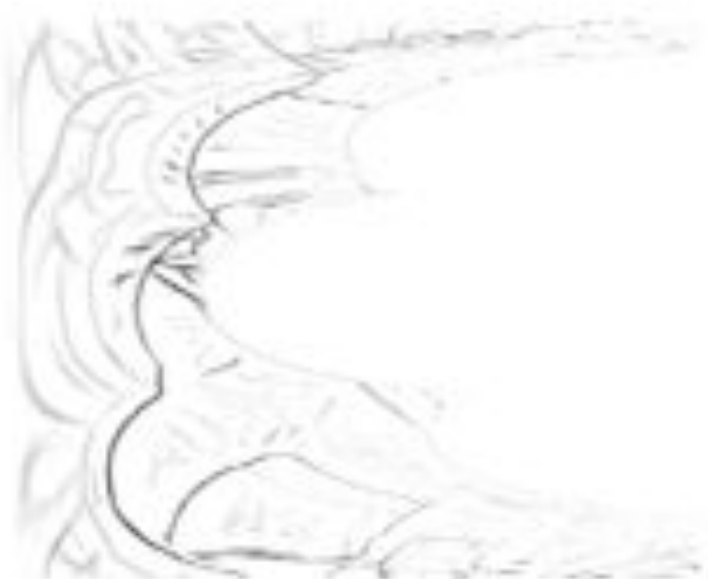
Scale normalization





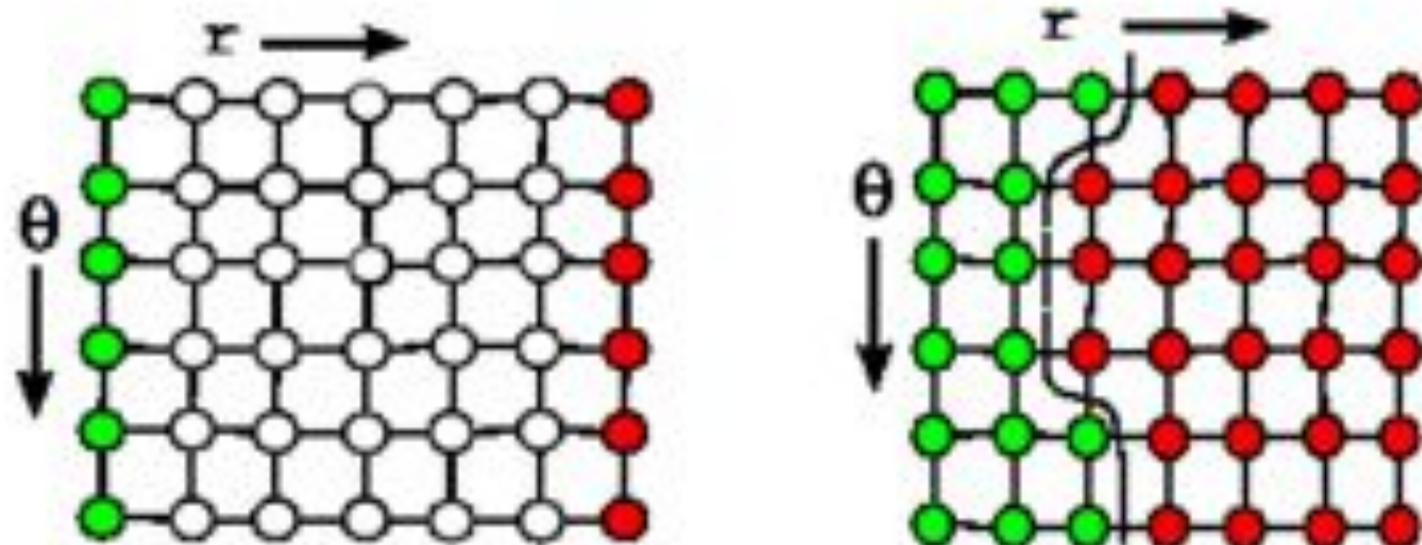
## Step 1: Cartesian to polar conversion

Scale normalization



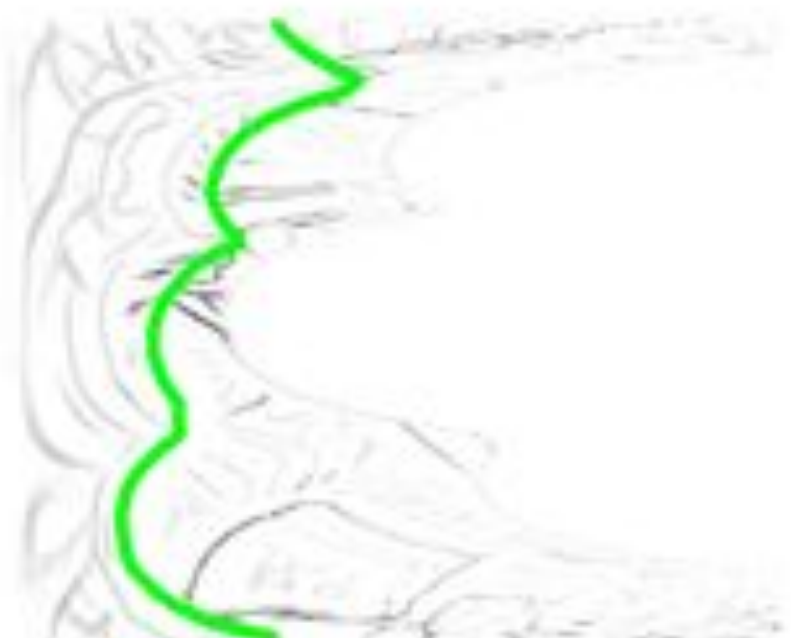
## Step 2: Find the optimal cut

binary labeling problem: Inside Vs outside

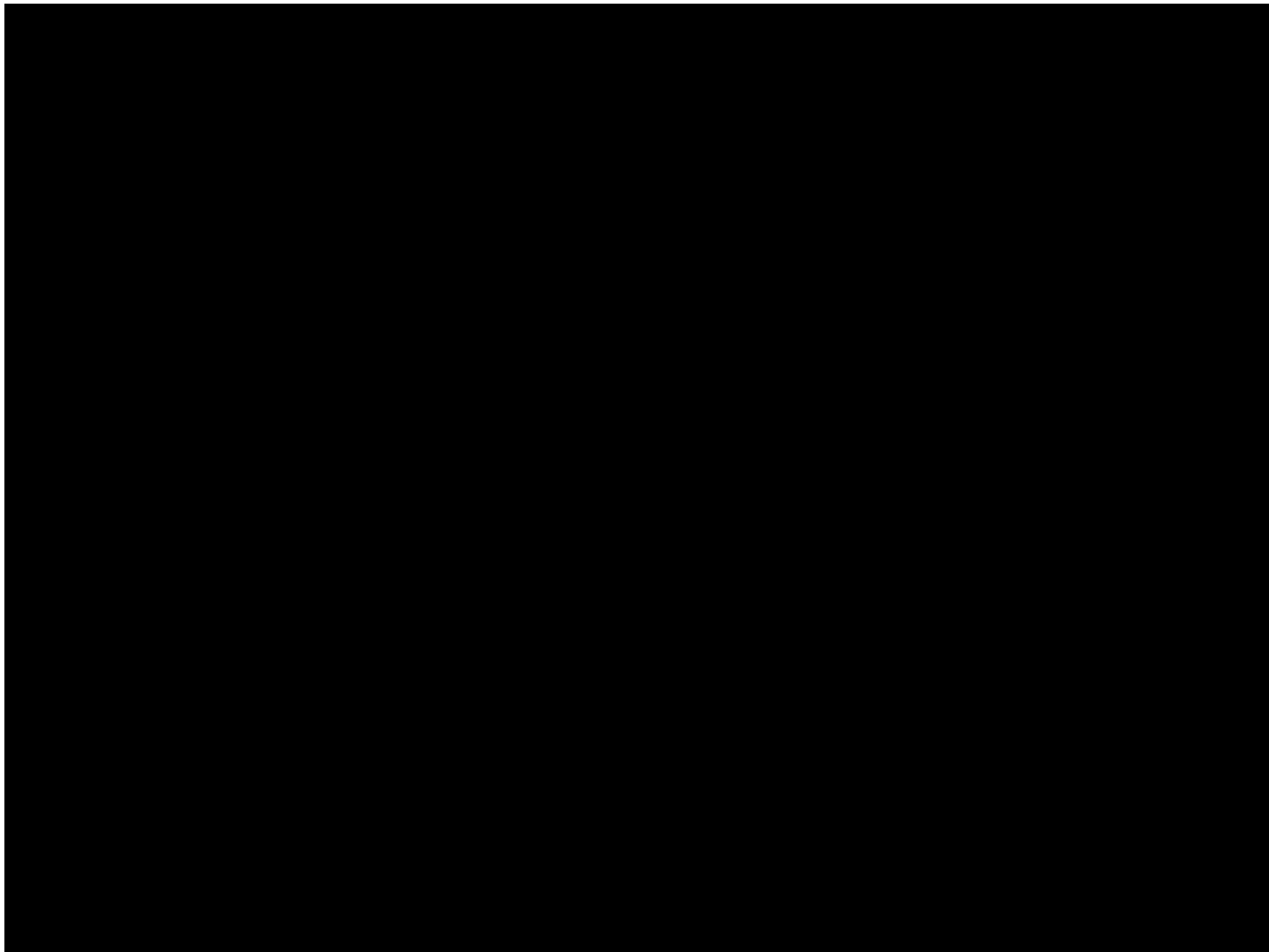


## Step 2: Find the optimal cut

binary labeling problem: Inside Vs outside

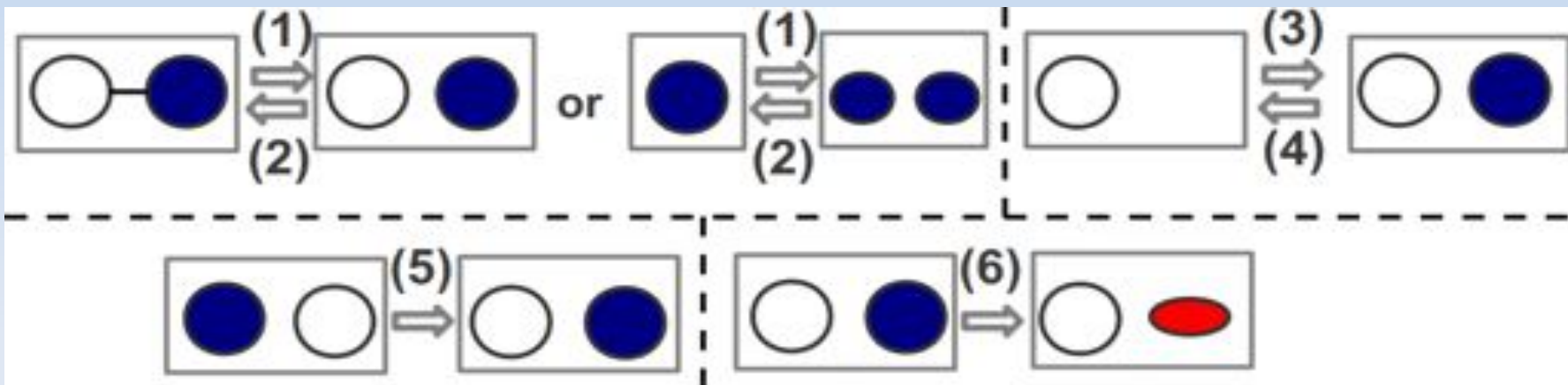




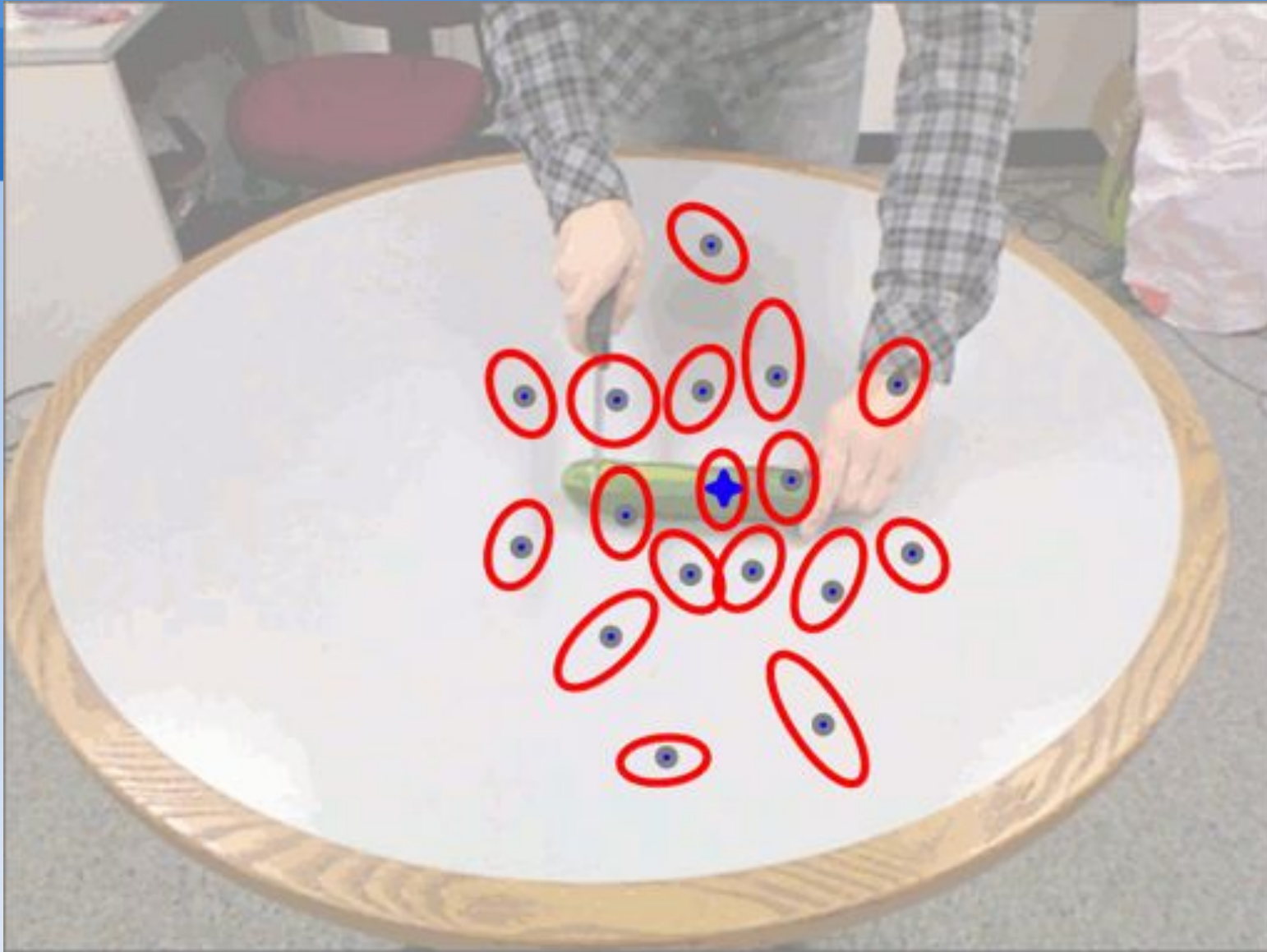


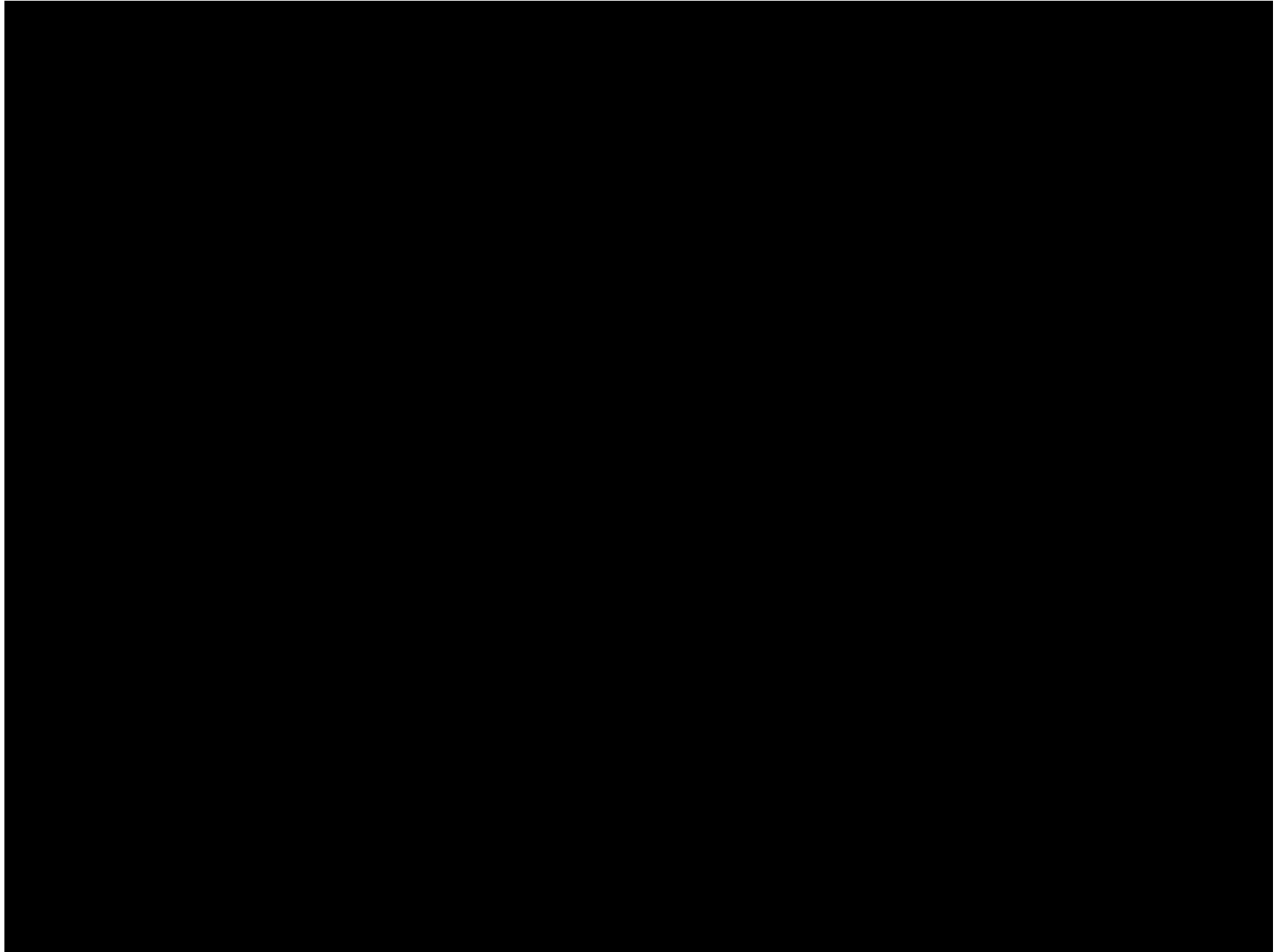
# Definition of SIX Primitive Action Consequences

- (1) ASSEMBLE: Two objects merge into one object, or two objects build an attachment between them;
- (2) DIVIDE: One object breaks into two objects, or two attached objects break the attachment;
- (3) CONSUME: An object is disappeared from the visual space;
- (4) CREATE: An object is brought to, or emerge into the visual space;
- (5) TRANSFER: An object is moved from one location to another location;
- (6) DEFORM: An object has an appearance change.

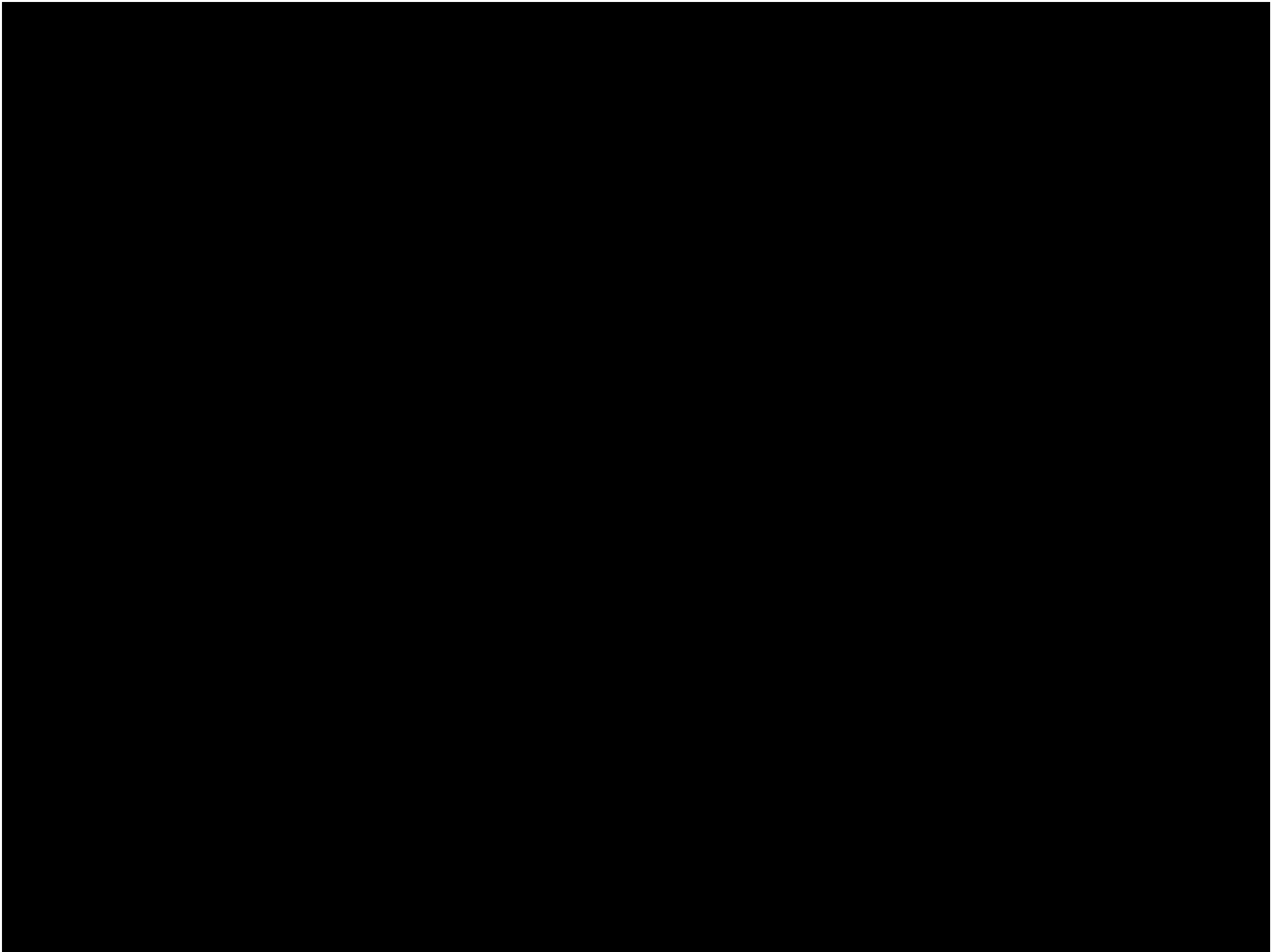


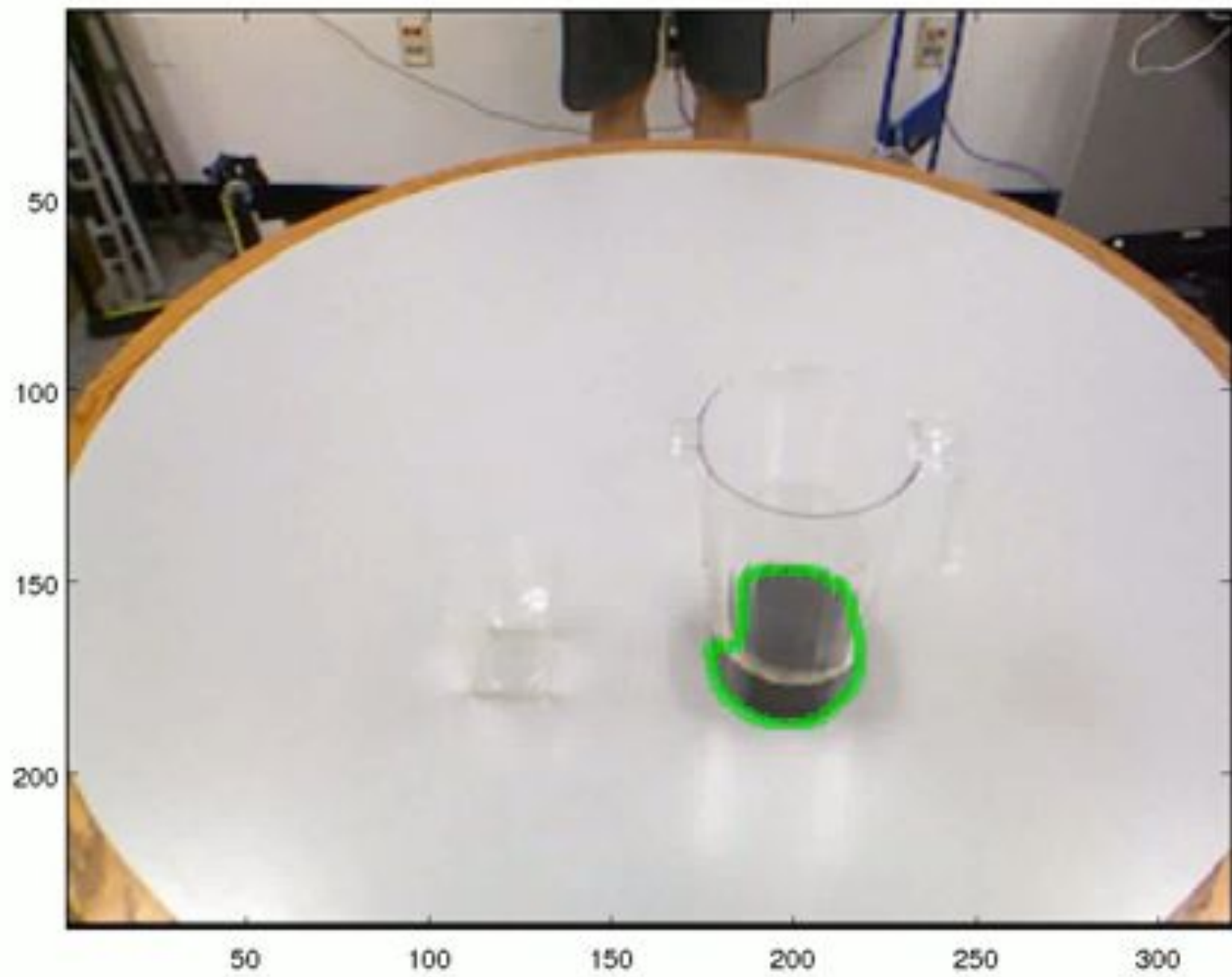
# Action Consequences Detection Through Object Monitoring

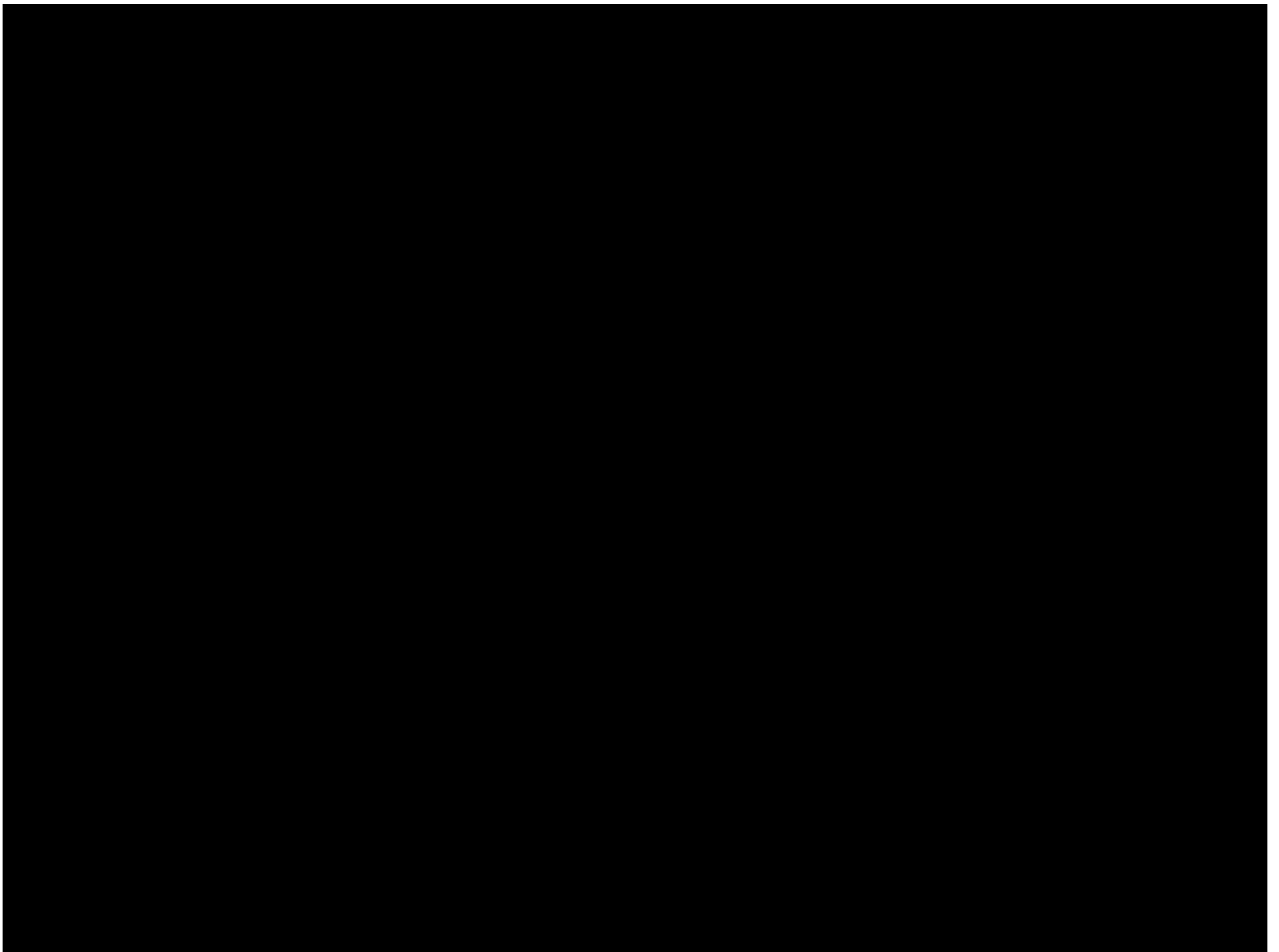




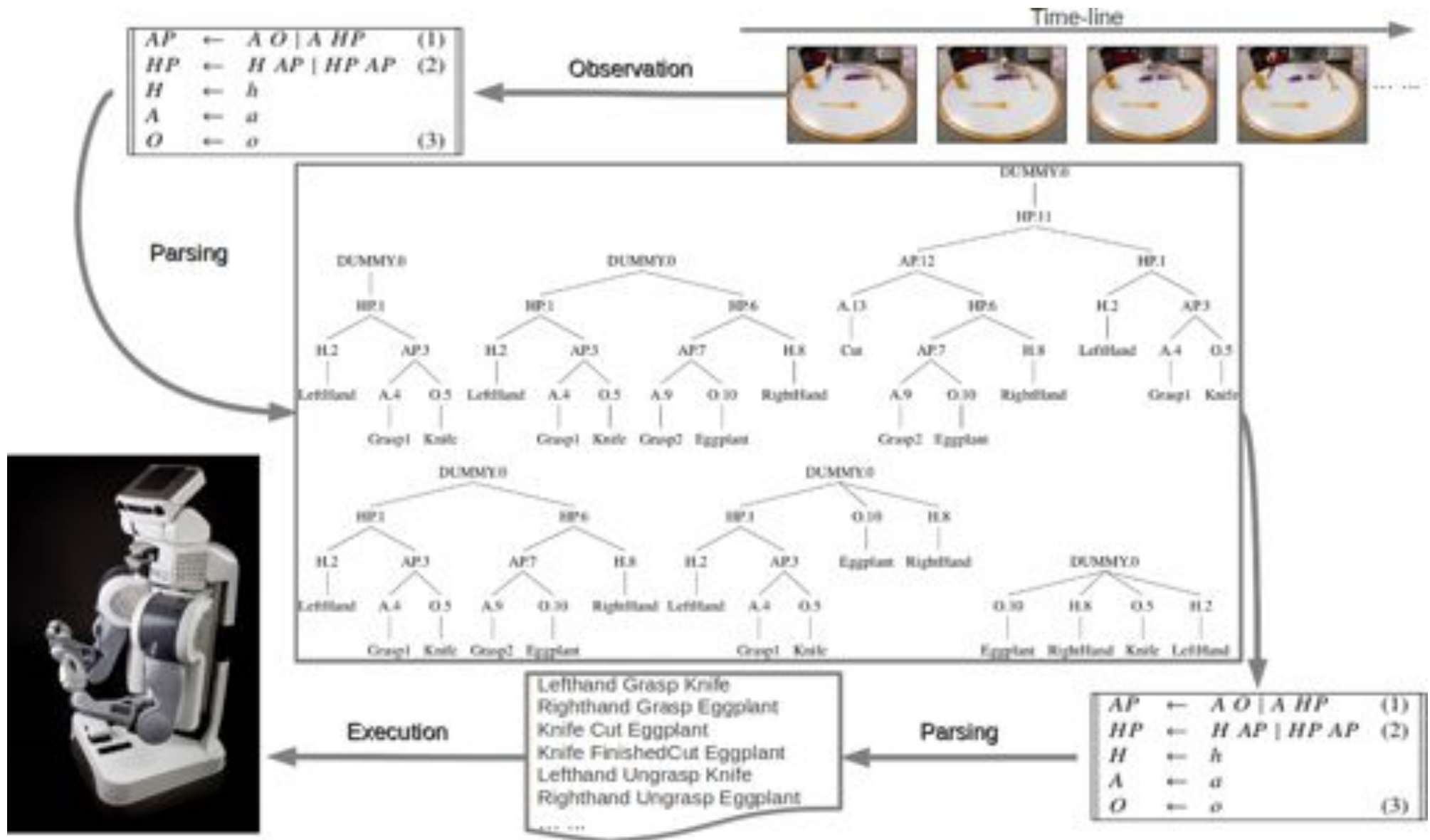


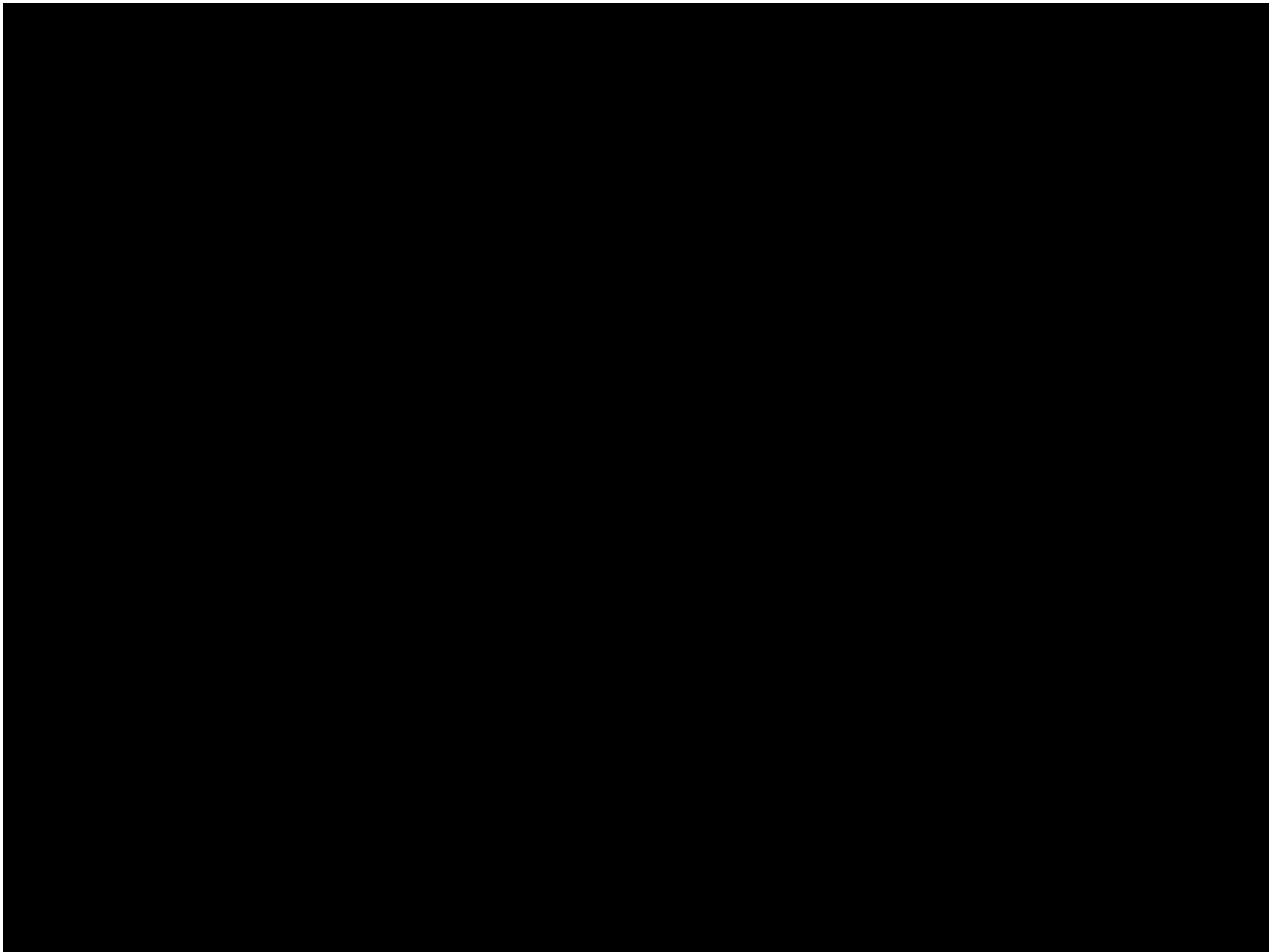




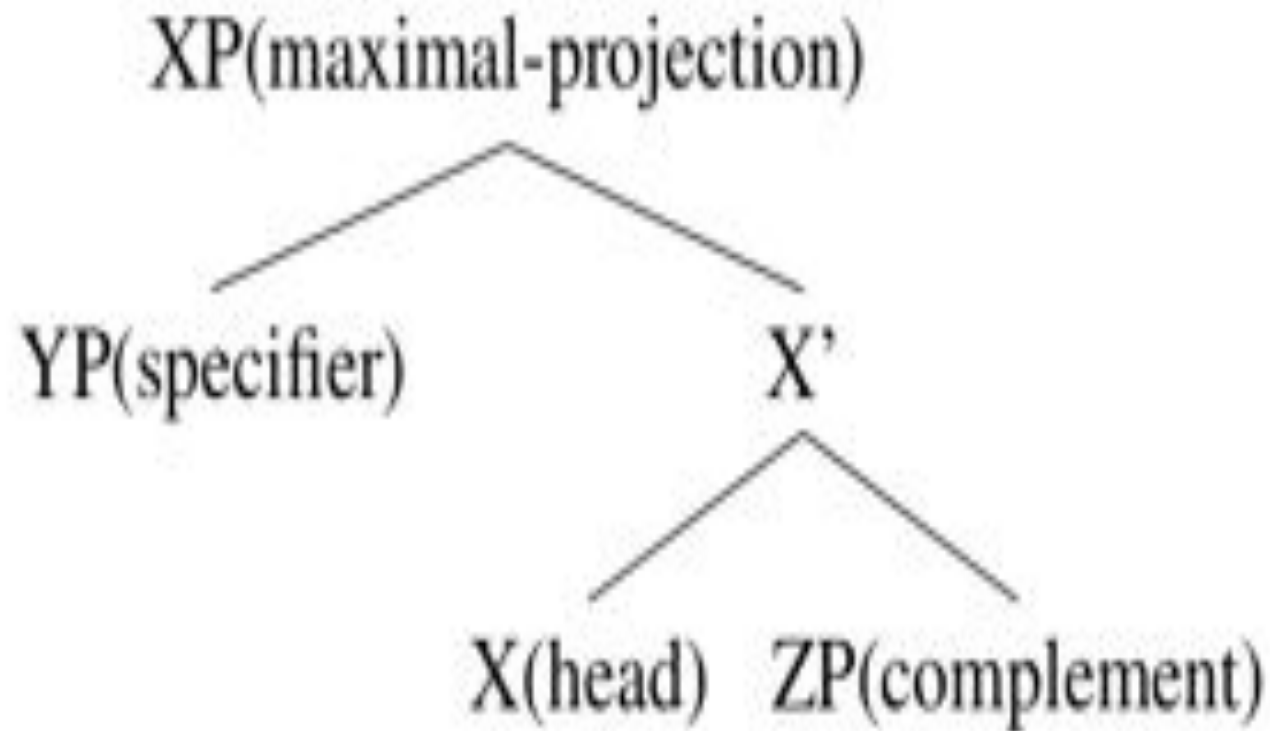


# An Example for Cognitive Robots

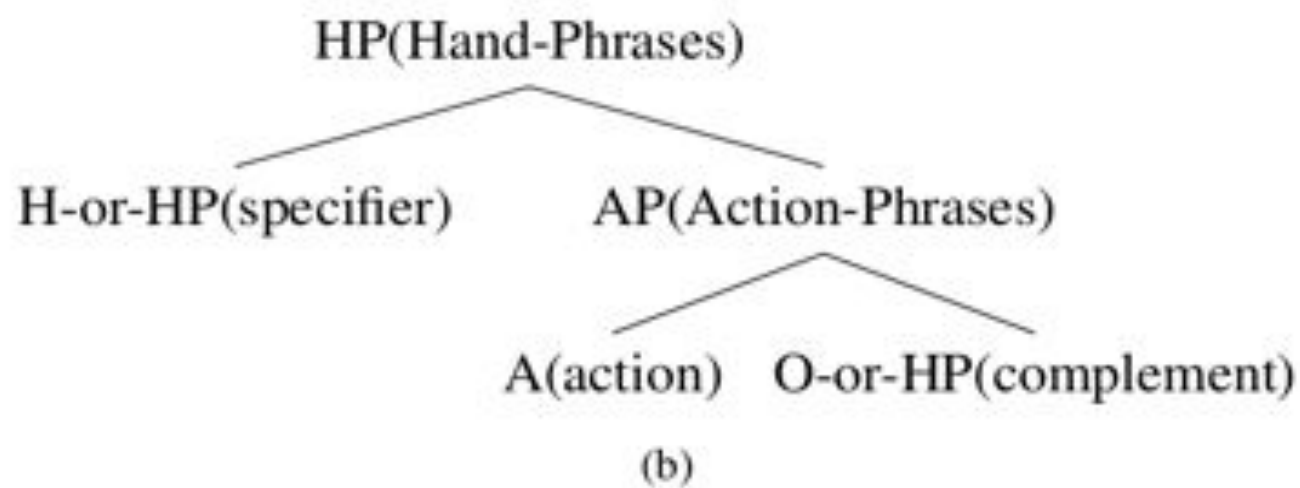
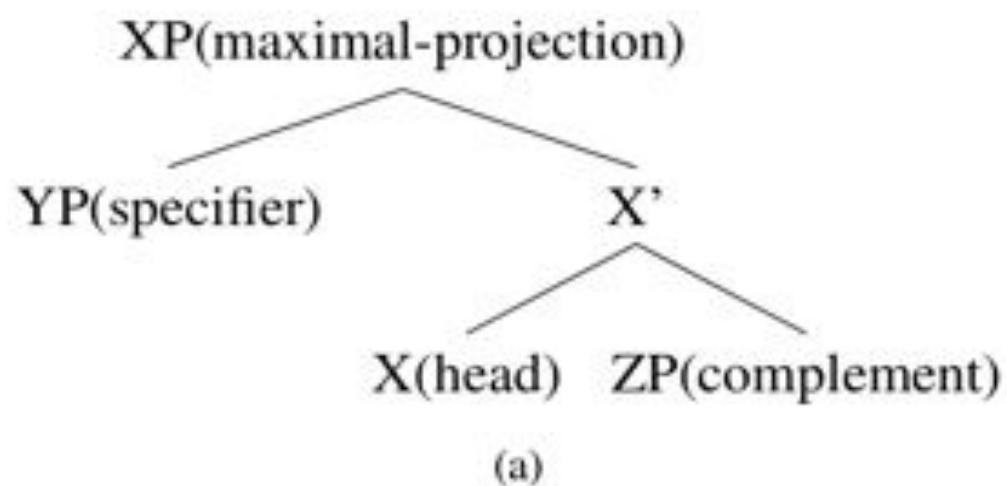




# X-bar Schema









# The Grammar Rules: terminal symbols

*H* ← *h*

*A* ← *a*

*O* ← *o*

H for Hands, A for Actions, O for Objects. AP is Action Phrases, HP is Hand Phrases

# The Grammar Rules: Rules to make Action Phrases (AP)

$$AP \leftarrow A O \mid A HP$$

Insight behind:

- An “Action” (A) can be applied to an “Object” (O) directly,
- or to a “Hand Phrase” (HP), which in turn contains an “Object” (O).

# The Grammar Rules: Rules to make Hand Phrases (HP)

$$HP \leftarrow H AP \mid HP AP$$

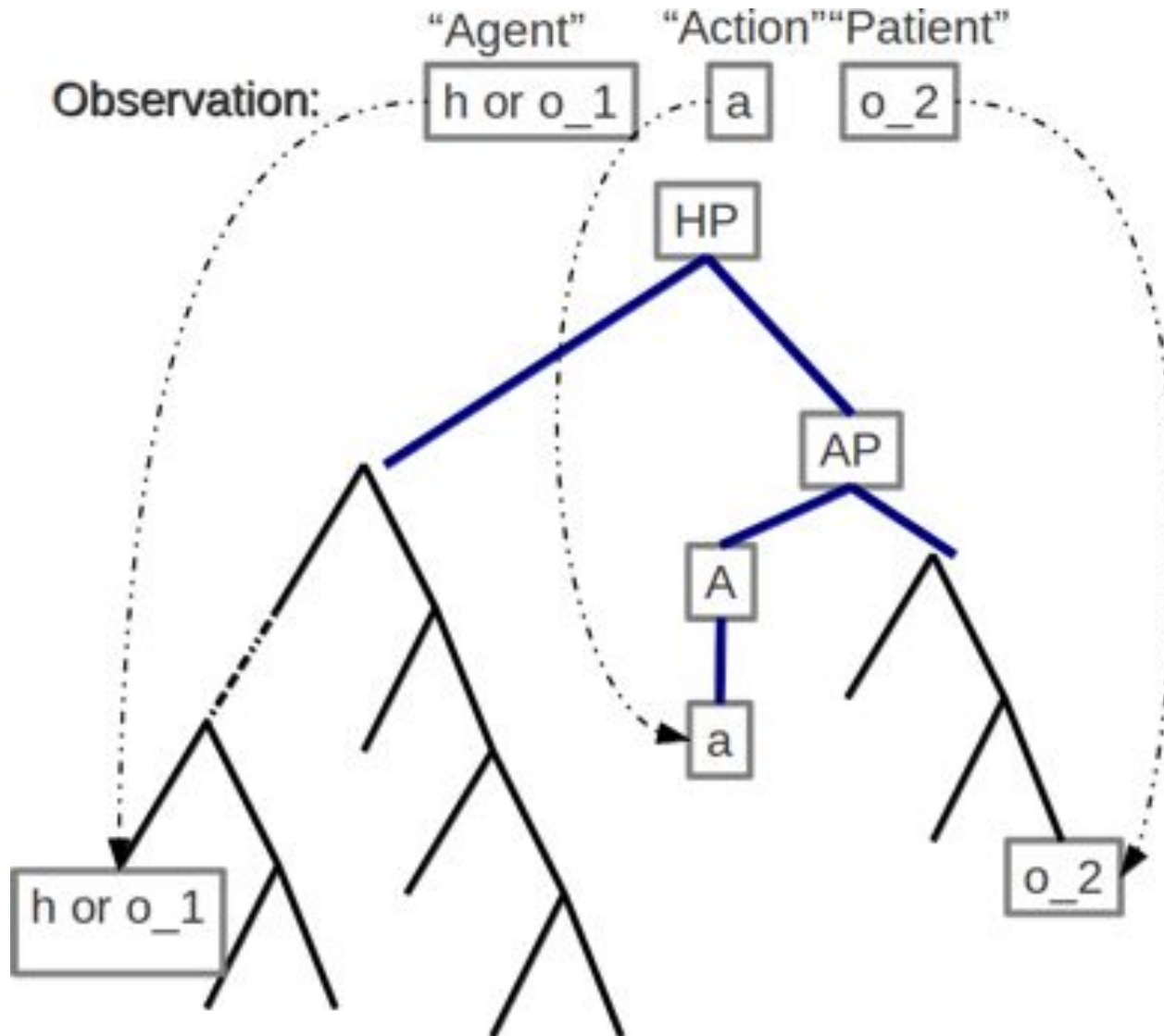
Insight behind:

- An “Action Phrase” (AP) can be combined either with the “Hand”(H), or a “Hand Phrase”, which recursively builds up the “Hand Phrase”.

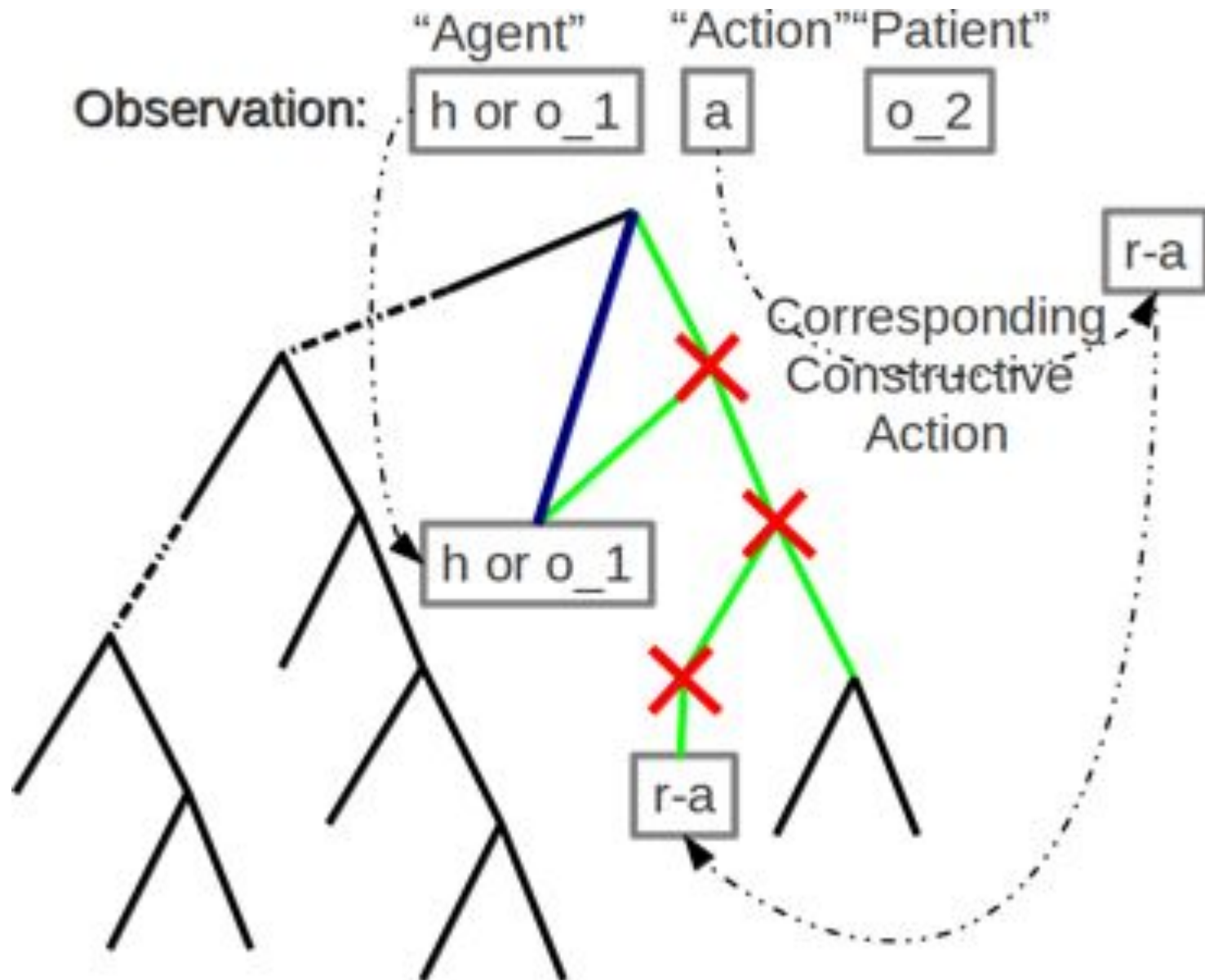
# How to parse?

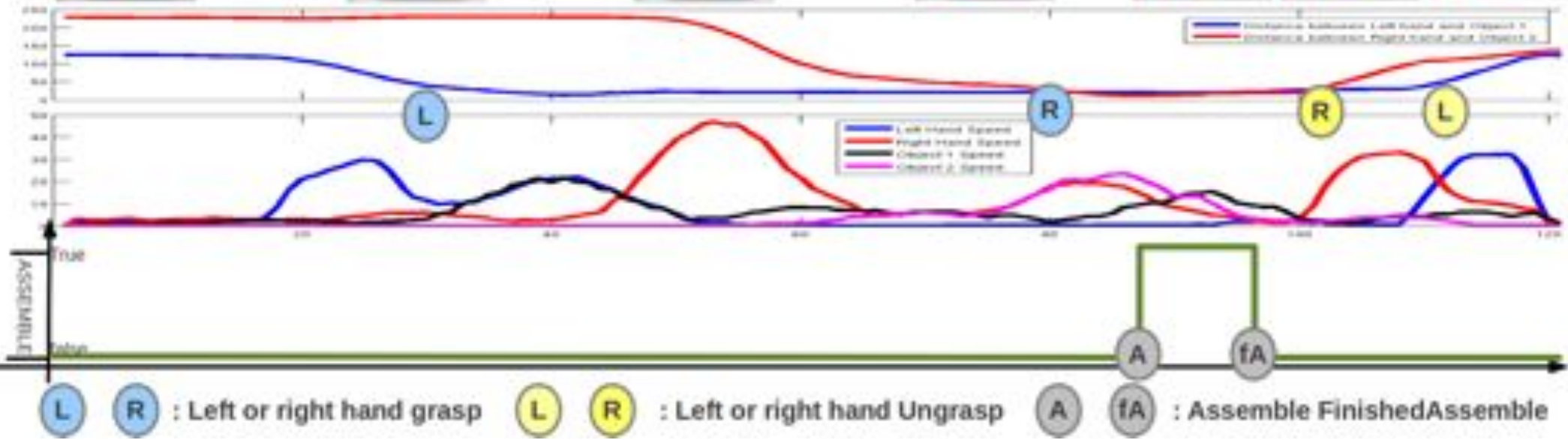
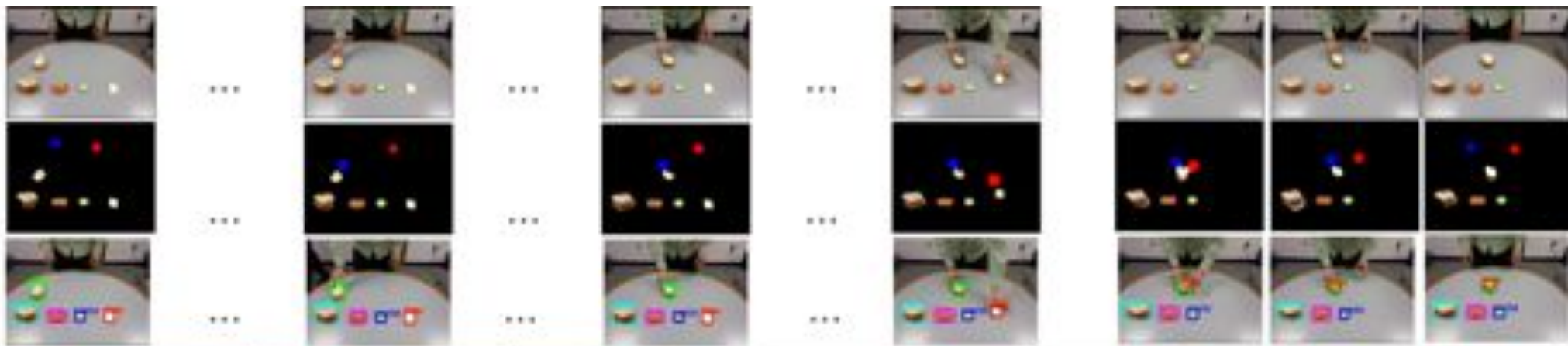
- The visual observations are obtained in a temporal sequence;
- A parsing algorithm is needed for the grammar should to dynamically update the semantic tree structures on the fly;
- Constructive (“Grasp”, “Cut”, etc) and Destructive Actions (“unGrasp”, “finishedCut”, etc).

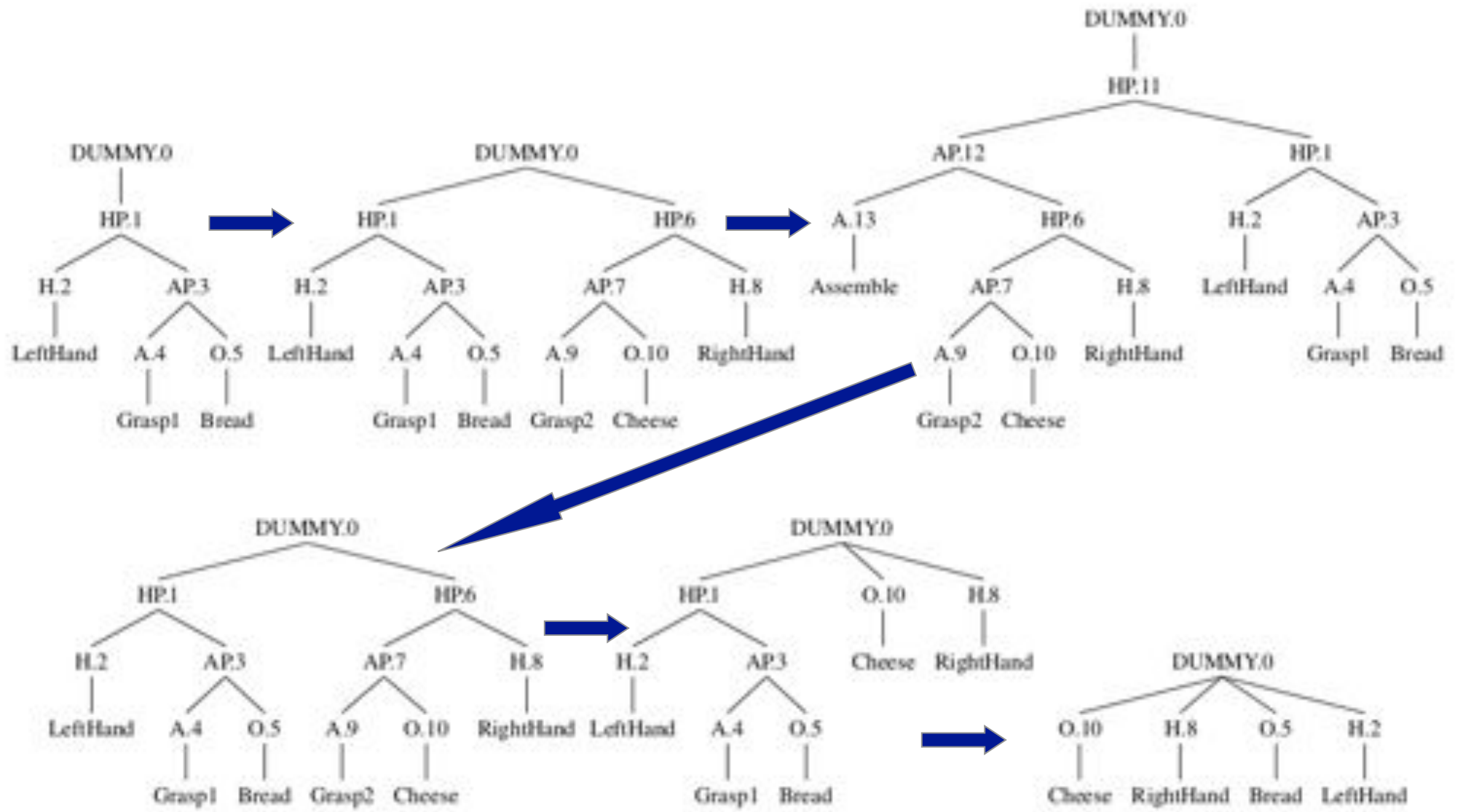
# The Parsing Algorithms: Construction



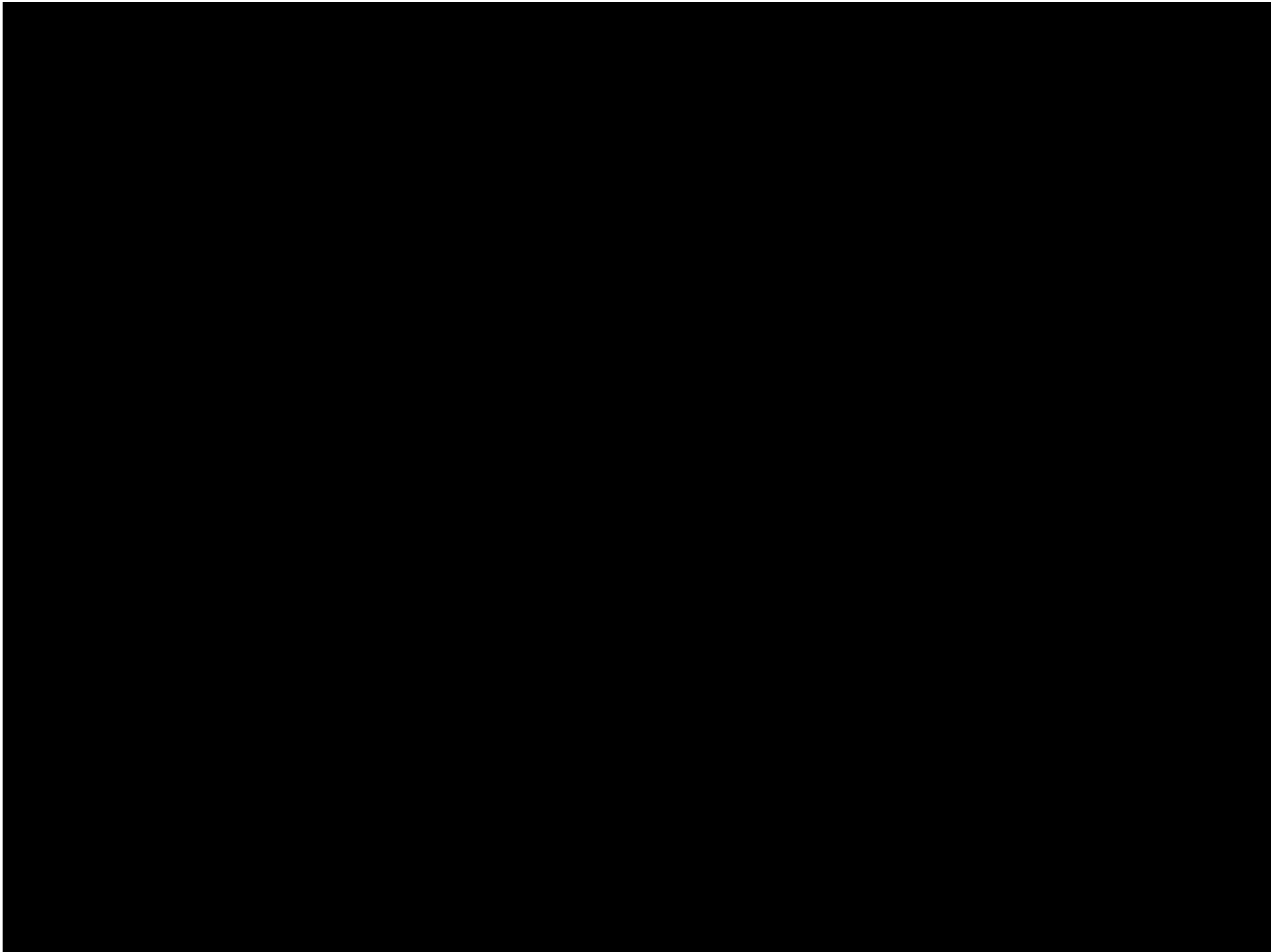
# The Parsing Algorithms: Destruction

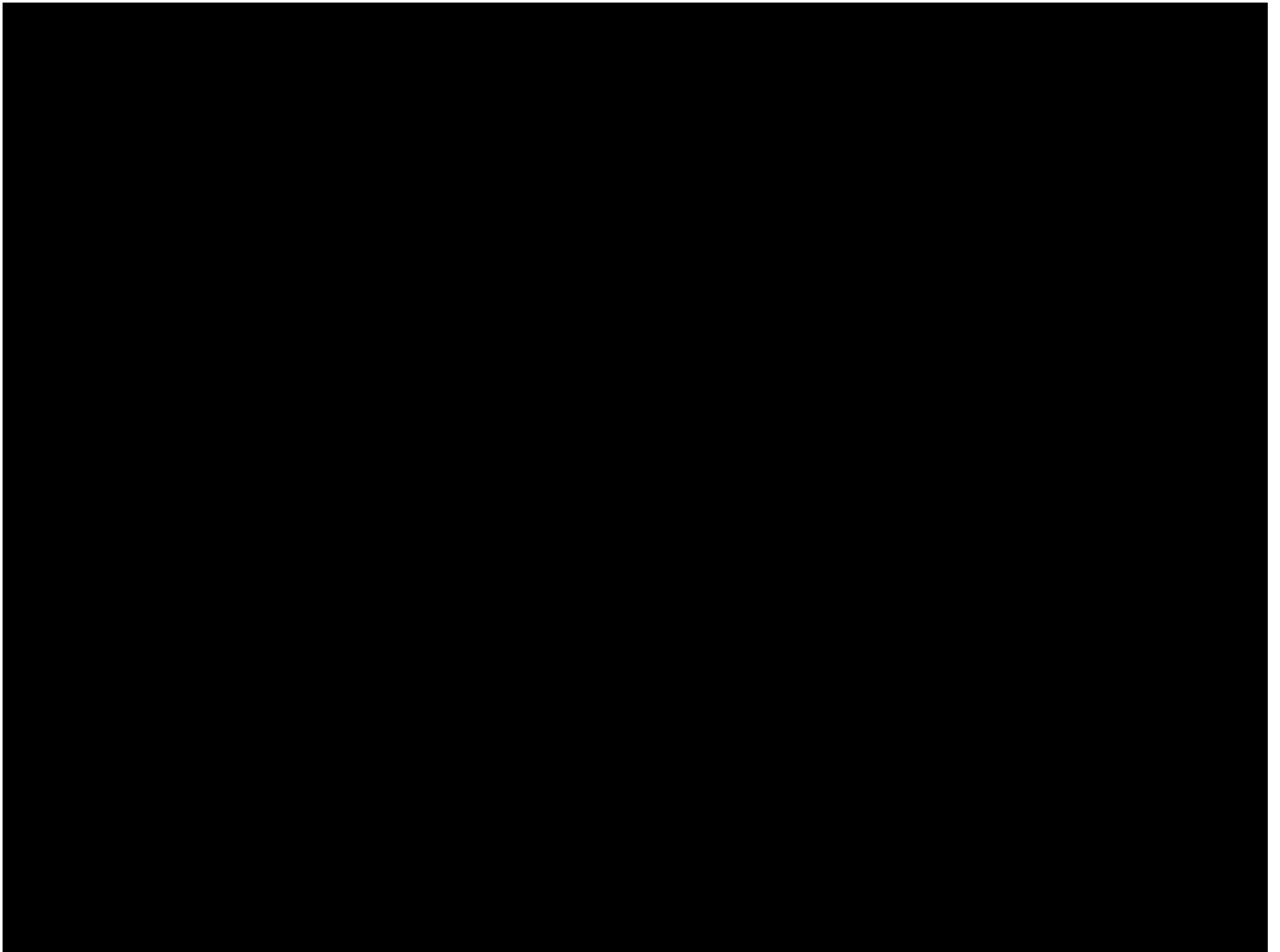












# What's next ?

- How to do the learning?
  - From various observations.
  - From raw language resources like recipes or wikipedia.
  - From human instructions.
- Human Robot Interaction and Collaboration?
  - Action prediction.
  - Warning for mistakes.
- Primitive action recognition, beyond hand trajectories?

**Draw line**

