

# Interactive Learning *and its role in* Pervasive Robotics

Cynthia Matuszek, Nicholas FitzGerald,  
Evan Herbst, Dieter Fox, Luke Zettlemoyer  
*{cynthia | nfitz | eherbst | lsz | fox} @ cs.washington.edu*



University of  
Washington

# Motivation

2

- ◆ **Teachable Robotics**
  - ◆ Systems that **learn** tasks **from end users**
- ◆ Robots becoming more ubiquitous
  - ◆ Cooking, toys, homes...
- ◆ Move away from “mainframe model”
  - ◆ Defined tasks
  - ◆ Fixed environment
  - ◆ Expert interaction
- ◆ An HRI model for low-cost, widely deployed robots.



# Goals

3

- ◆ Make robots capable of natural learning
  - ◆ Interpret and execute upon human input
  - ◆ Interactively learn about physically-grounded objects, attributes, and skills
  - ◆ Many components: language, LBD, active learning...
- ◆ Make ubiquitous robots useful
  - ◆ Robots operating alongside people
  - ◆ Don't try to pre-conceive all possible needs
- ◆ “Grounded Language Acquisition”
  - ◆ Learn language from interacting with people and the world

# Case Studies in Teachable Robotics

4

## 1. Semantic Mapping and Direction Following

- ◆ Following human instructions
- ◆ Real-time world model discovery

ISER 2012 – Matuszek,  
Herbst, Zettlemoyer, Fox

## 2. Learning Object Attributes

- ◆ What colors and shapes exist in the world
- ◆ How people refer to them

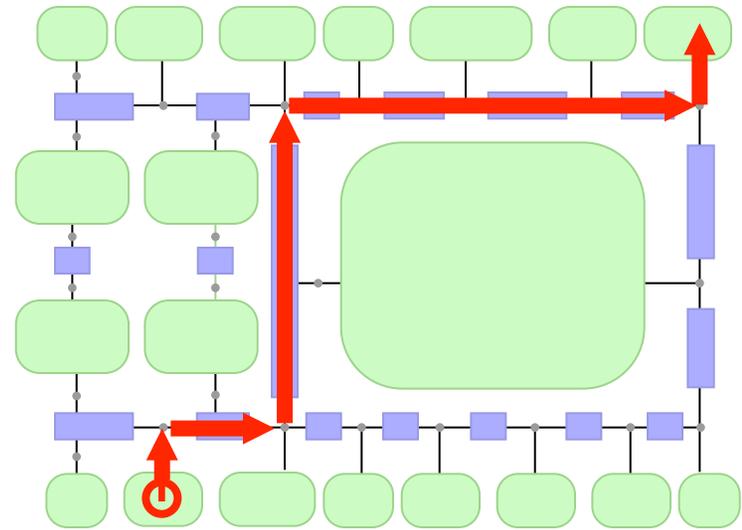
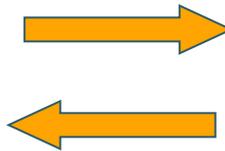
ICML 2012 – Matuszek,  
FitzGerald, Bo,  
Zettlemoyer, Fox

- ◆ Lessons learned
- ◆ Challenges and next steps

# Route Instruction Following

5

"Leave the room and turn right, take the first left, go past the big room and go right, then go to the end of the hall and turn left."



- ◆ Humans are pretty bad at giving/following instructions
  - ◆ Missed turns, left/right errors, ...
  - ◆ *Humans* can only follow human instructions ~70% of time\*
- ◆ Many sources of uncertainty
  - ◆ Map labeling errors, parse errors,
  - ◆ This is a **well-studied** task, but not an **easy** one.

\* Reisbeck et al, 1985;  
MacMahon et al, 2006;  
Matuszek et al, 2010

# Approach: Semantic Parsing

6

- ◆ Traditionally: NL  $\longleftrightarrow$  NL
  - ◆ Use sentence pairs to learn translation model  
"Ich bin müde"  $\longleftrightarrow$  "I'm tired"
- ◆ Our approach: NL  $\longleftrightarrow$  Robot Control Language
  - ◆ Source: natural language directions
  - ◆ Target: Robot Control Language grounded in map
- ◆ Learned parser based on many sentence pairs:

"Go down the hall  
to the intersection,  
then take a left."

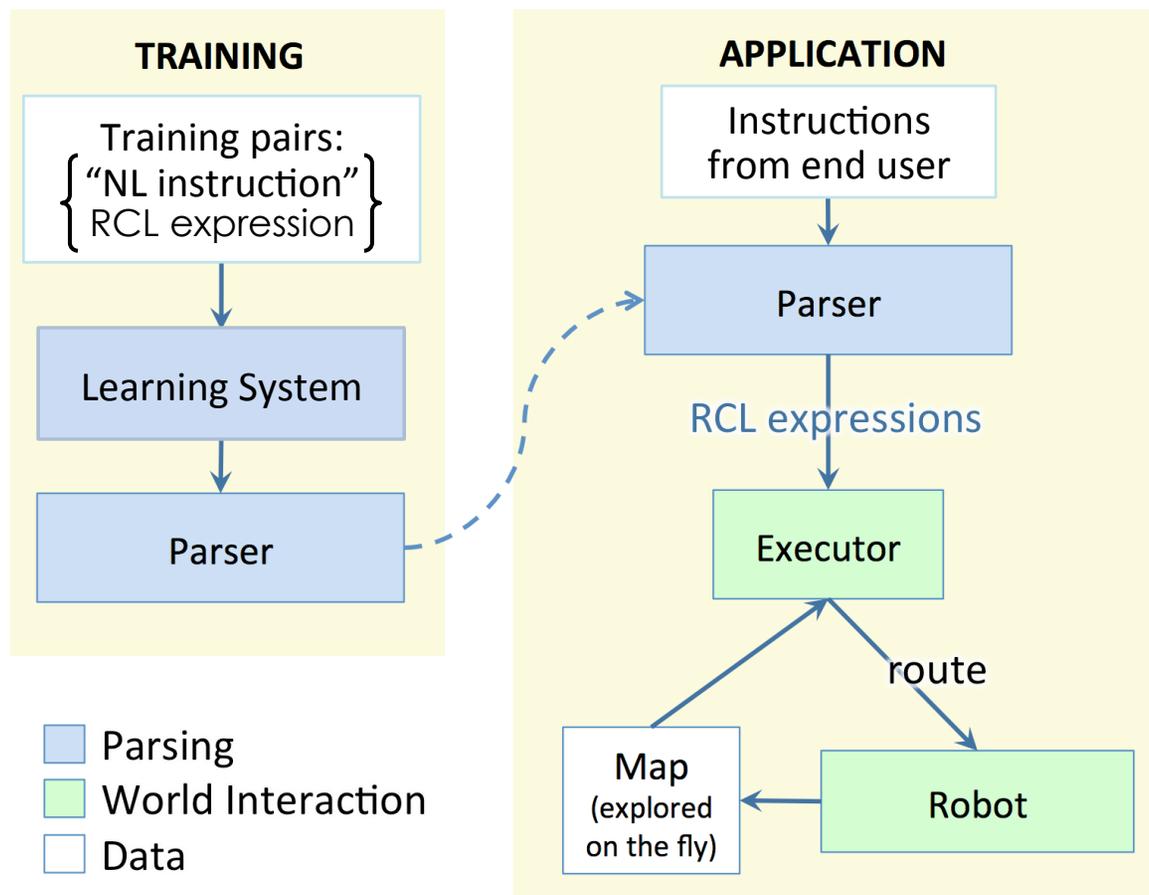


(do-sequentially  
(do-until  
(junction current-loc)  
(move-to forward-loc))  
(turn-left)))

# Robot Control Language

7

- ◆ Formal robot control language (lambda-calculus)



# Example Commands

8

“Go left to the end of the hall.”

```
(do-sequentially
  (turn-left current-loc)
  (do-until
    (or
      (not (exists forward-loc))
      (room forward-loc))
    (move-to forward-loc)))
```

“Go to the third junction and take a right.”

```
(do-sequentially
  (do-n-times 3
    (do-sequentially
      (move-to forward-loc)
      (do-until
        (junction current-loc)
        (move-to forward-loc))))
  (turn-right current-loc))
```

- ◆ Humans generate English; our system generates RCL
- ◆ Assumptions: robot can execute actions, recognize objects, and determine conditionals
- ◆ Primitives can encode complex activities

# Application to Route Instructions

9

- ◆ Training corpus
  - ◆ ~1000 paired route instruction sets
- ◆ Testing: instructions on novel maps
- ◆ How often does the robot reach the goal **by the intended path** (of human gold standard)?

Successes:

Short	924/1000	<b>66%</b>
Long	125/200	<b>49%</b>

- ◆ No error correction, so longer paths are strictly harder

# Teaching with Language: Successes and Challenges

10

- ◆ Lessons and Successes
  - ✓ Robots can learn about tasks from end users
  - ✓ Language is intuitive and natural for giving directions
  - ✓ Language grounding ties in NLP community
- ◆ Challenges
  - ✗ Formal annotations still expensive
  - ✗ Need lots of NL training data
  - ✗ Still bound to a pre-defined grammar
    - ✗ Can't teach tasks that can't be expressed in existing RCL
- ◆ Local error recovery is crucial

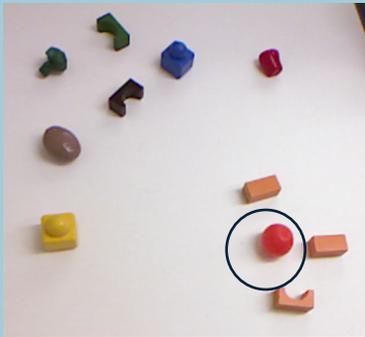
# Case Study: Learning Attributes

11

- ◆ Learn to handle sentences about **novel concepts**
  - ◆ “These are limes”
  - ◆ No longer assuming underlying concepts already exist in RCL
    - ◆ What if ‘green’ is totally new?
- ◆ Requires:
  - ◆ Perception models
    - ◆ **green; round**
  - ◆ Language model
    - ◆ How words relate to these detectors
- ◆ Need a joint model for learning these together!


$$\lambda x . [\text{green}(x) \wedge \text{round}(x)]$$

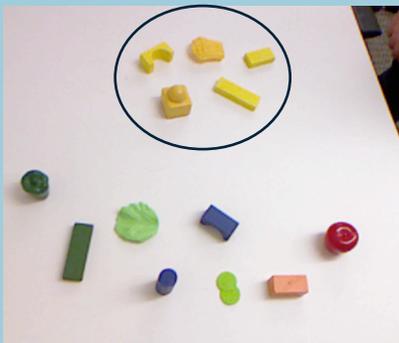
### 1: Initialization



“This is an orange ball.”

$\text{obj-color}(x, \text{color-orange}) \wedge \text{obj-shape}(x, \text{shape-round})$

### 2: Training



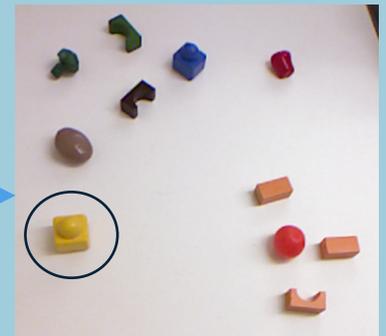
“All of these toys are yellow.”

### 3: Testing



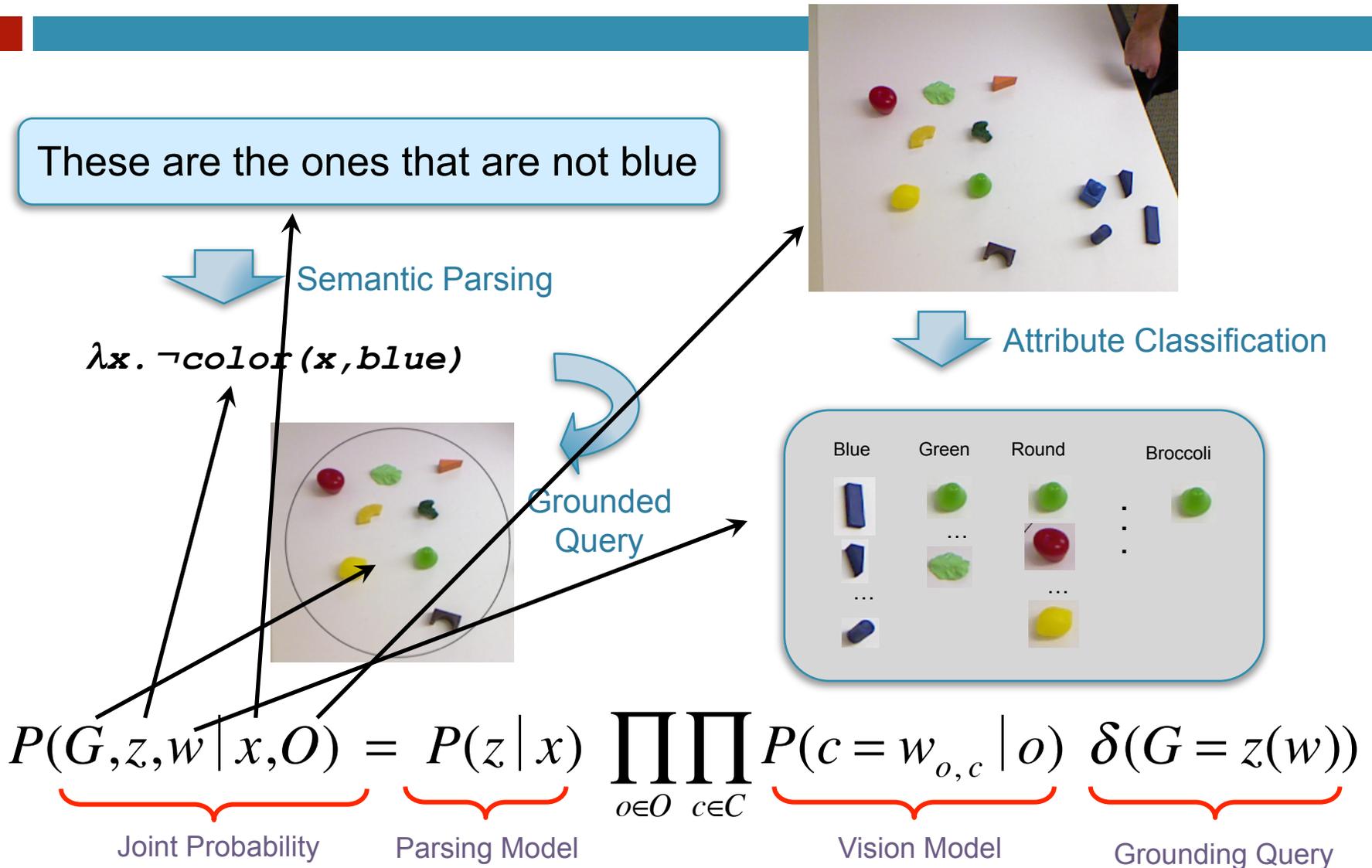
“It’s the yellow one.”

$\text{obj-color}(x, \text{color-NEW})$



# Joint Language / Perception Model

13



# NL Data from Crowdsourcing

## What is the Parent Saying?

Watch the video, then **describe what the parent is saying to the child**, in complete sentences.



- Pretend you are a parent teaching a child about something.
- The question is:

*How does the parent describe this group of objects?*

Your answer should be the sentence(s) the parent said while pointing to these things.

Submit

“This one’s an orange ball.”

Showing HIT 1 of 3

Next HIT

$\lambda x. \text{obj-color}(x, \text{color-orange}) \wedge \text{obj-shape}(x, \text{spheroid})$

# Percepts Reduce Uncertainty

15

- ◆ New language can be ambiguous. Maintain **hypotheses**:

- ◆ “This is *<new-word>*”:

- ◆ New attribute (color)?

“This is red.”

$\lambda x.\text{obj-color}(x, \text{color-NEWCOL6})$

- ◆ New attribute (shape)?

“This is arched.”

$\lambda x.\text{obj-shape}(x, \text{color-NEWSHP1})$

- ◆ Synonymy?

“This is peach.”

$\lambda x.\text{obj-color}(x, \text{color-orange})$



After a few scenes, one hypothesis will have best predictive power.

# Experimental Evaluation

16

- ◆ Fully Supervised Training: NL sentences from Mechanical Turk
  - ◆ 1,003 sentence/annotation pairs:
    - ◆ “These are yellow” /  $\lambda x.obj-color(x, color-yellow)$
  - ◆ 142 **scenes** (image, with a circle showing positive data)



- ◆ Training
  - ◆ 3 colors, 3 shapes used for bootstrapping
  - ◆ Train on previously unseen 3 colors and 3 shapes
- ◆ Testing: novel English, trained attributes, and novel scenes

# Experimental Evaluation

17

- ◆ 18 trials
- ◆ Of dataset:
  - ◆ Average:
    - ◆ 502 supervised training triplets
    - ◆ 401 weakly supervised training pairs (NL/annotation)
    - ◆ 100 test NL/annotation pairs

$P = 0.82$ ;  $R = 0.71$ ;  $F = 0.76$

# Teaching with Language plus Vision: Successes and Challenges

18

## ◆ Lessons and Successes

- ✓ Language + vision = **better interaction**
  - ✓ Robot must be aware of human workspace
- ✓ Experts providing **just initializations** is easier
- ✓ **Crowdsourcing** for collecting non-expert data
- ✓ Completely new concepts can be learned **from users**

## ◆ Challenges

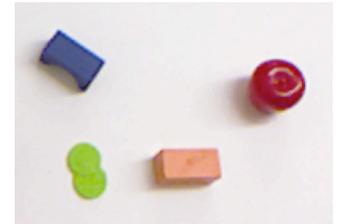
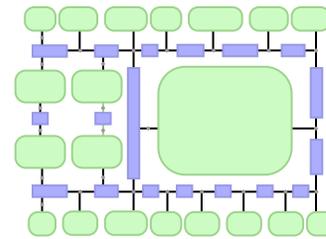
- x Some (expensive) annotation still needed
- x Still providing some **domain-specific** seeding
- x Need to learn more **efficiently**
  - x Still need too much training data to learn in one session

# Case Study Discussion & Ideas

19

- ◆ **Learn** to understand human instructions, descriptions

- ◆ Learn to follow instructions
- ◆ To identify novel world attributes



- ◆ Challenges / suggestions

- ◆ **Crowdsourcing** helps with (inevitable?) data needs
- ◆ Minimize expert involvement
- ◆ Avoid pre-defining uses



$\lambda x. \text{green}(x)$   
 $\wedge \text{round}(x)$

- ◆ **Language grounding** is powerful and applicable

- ◆ Future directions

- ◆ Beyond language: gesture, gaze, body language
- ◆ **(Inter)active** teaching / grounding



# Summary: Teachable Robotics

20

- ◆ Users already know how they want to interact.
- ◆ We won't always know what tasks will be.
- ◆ **Grounded Knowledge Acquisition** lets us learn:
  - ◆ *What* humans need
  - ◆ *How* they want to convey it
- ◆ Instead of designing interaction, design **teaching scenarios**.

