

# Learning to Parse Natural Language Commands to a Robot Control System

Cynthia Matuszek,  
Evan Herbst,  
Luke Zettlemoyer,  
Dieter Fox



*{cynthia/eherbst/lasz/fox}@cs.uw.edu*

# Motivation

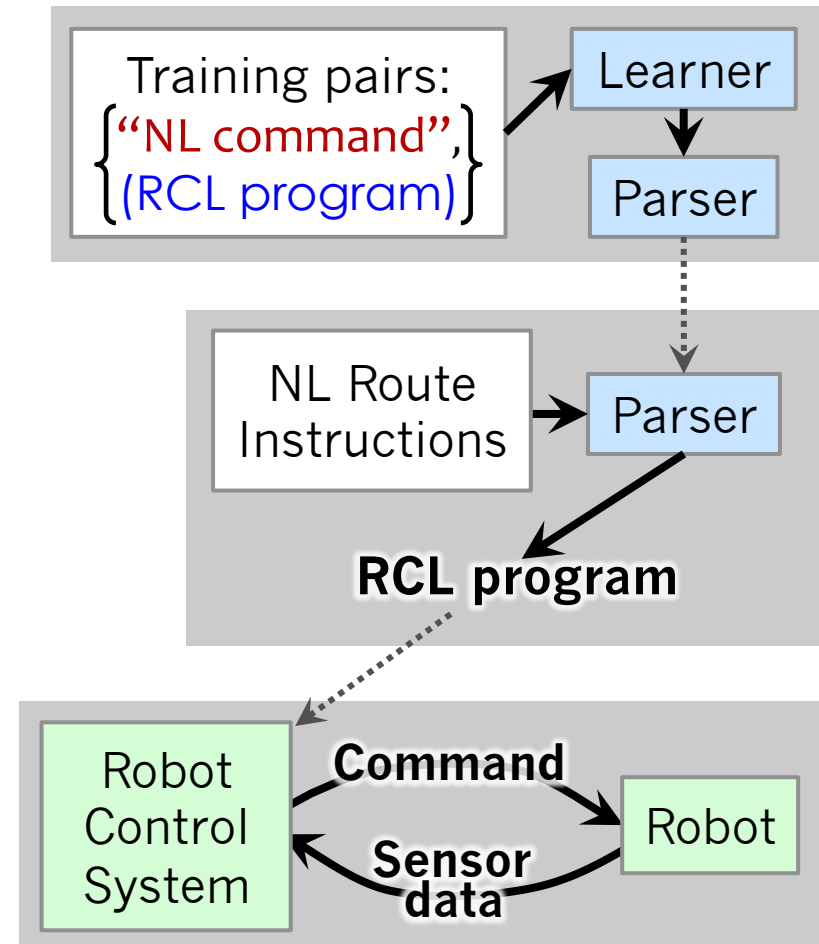
2

- ◆ More natural Human-Robot Interaction (HRI):
  - ◆ In the long term, this entails parsing rich human input: speech, gesture, gaze, ...
- ◆ Taking instructions from users in **Natural Language**
- ◆ Ideally, language understanding **learned** from data
- ◆ Key contributions:
  - ◆ Follow instructions in a previously unseen world
  - ◆ Learn from data to parse natural language
    - ◆ Into robot-executable control system

# Goal

3

- ◆ “Grounded Language Acquisition”
  - ◆ Transform natural language into semantically meaningful representation
  - ◆ Map that information to to perceived world
- ◆ **Learn a parser**
  - ◆ Produces robot-executable commands from NL instructions.



# Some Related Work

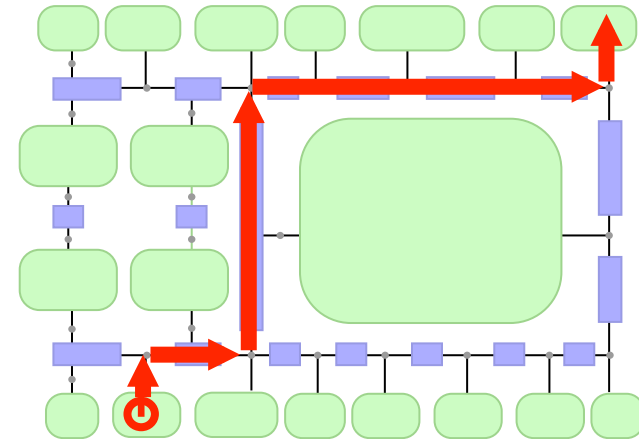
4

- ◆ Logic-based representations for robot control  
[Beetz *et al*, Konolige *et al*, Kress-Gazit *et al*, Dantam-Stilman *et al*, ...]
- ◆ Direction following in rich simulated environments  
[Kuipers-MacMahon-Wong *et al*, Chen-Mooney, ...]
- ◆ Learned semantic parsing [Zettlemoyer *et al*, Liang *et al*, ...]
- ◆ Learn to parse NLP for RoboCup and direction following (with minimal supervision) [Mooney *et al*]
- ◆ Parsing NL in known world and action models: for direction following; for forklift operation  
[Tellex-Kollar-Roy *et al*]

# Testbed: Route Instructions

5

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



- ◆ Previous work grounded instructions directly into the map – no target concepts such as **while**

"Take the second left."

- ◆ Parser must be able to produce many **possible** groundings: →
- ◆ High-level concepts are worse:
  - ◆ "go to the end of the hall,"
  - ◆ "keep turning right until you can't any more", ...

```

1: (go (hall) (4junction 1)
      (hall) (3junction lt 0) (room))
2: (go (room) (4junction 1)
      (room) (3junction lt 0) (room))
3: (go (hall) (4junction 1)
      (hall) (3junction rt 1) (room)
      (3junction lt 0) (room)) ...

```

# Example Commands

6

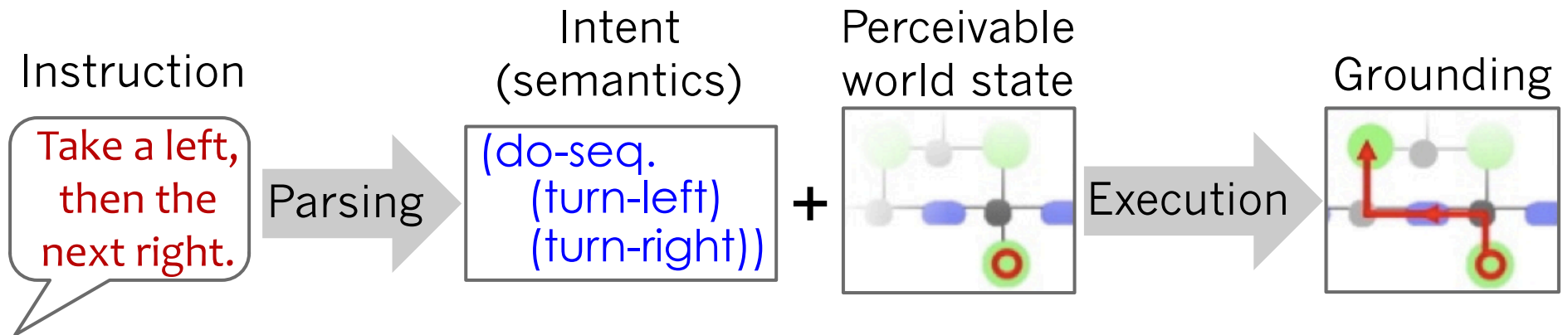
```
“Go left to the end of the hall.”  
(do-sequentially  
  (turn-left current-loc)  
  (do-until  
    (or  
      (not (exists forward-loc))  
      (room forward-loc))  
    (move-to forward-loc)))
```

```
“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; parser generates RCL
- ◆ Assumptions: robot can execute actions, recognize objects, and determine conditionals
- ◆ Primitives can encode complex activities

# Approach: Semantic Parsing

7



- ◆ **Parse from NL to a formal control language:** Robot Control Language, or RCL.
- ◆ **Train** semantic parsing model
  - ◆ → Distribution over RCL sequences for any NL sentence
- ◆ Application of learned system: parse new instructions, with simulated agent in unknown map

# Categorial Combinatory Grammars

8

- ◆ Capture **syntax** and **semantics** of language
- ◆ Parse sentences to expressions in  $\lambda$ -calculus
- ◆ Space of possible parses defined by:

lexical entries  $\begin{cases} \text{go to} \vdash S / NP : \lambda x. \text{move-to}(x) \\ \text{junction} \vdash N : \lambda x. \text{junction}(x) \end{cases}$

along with **combinatory rules**.

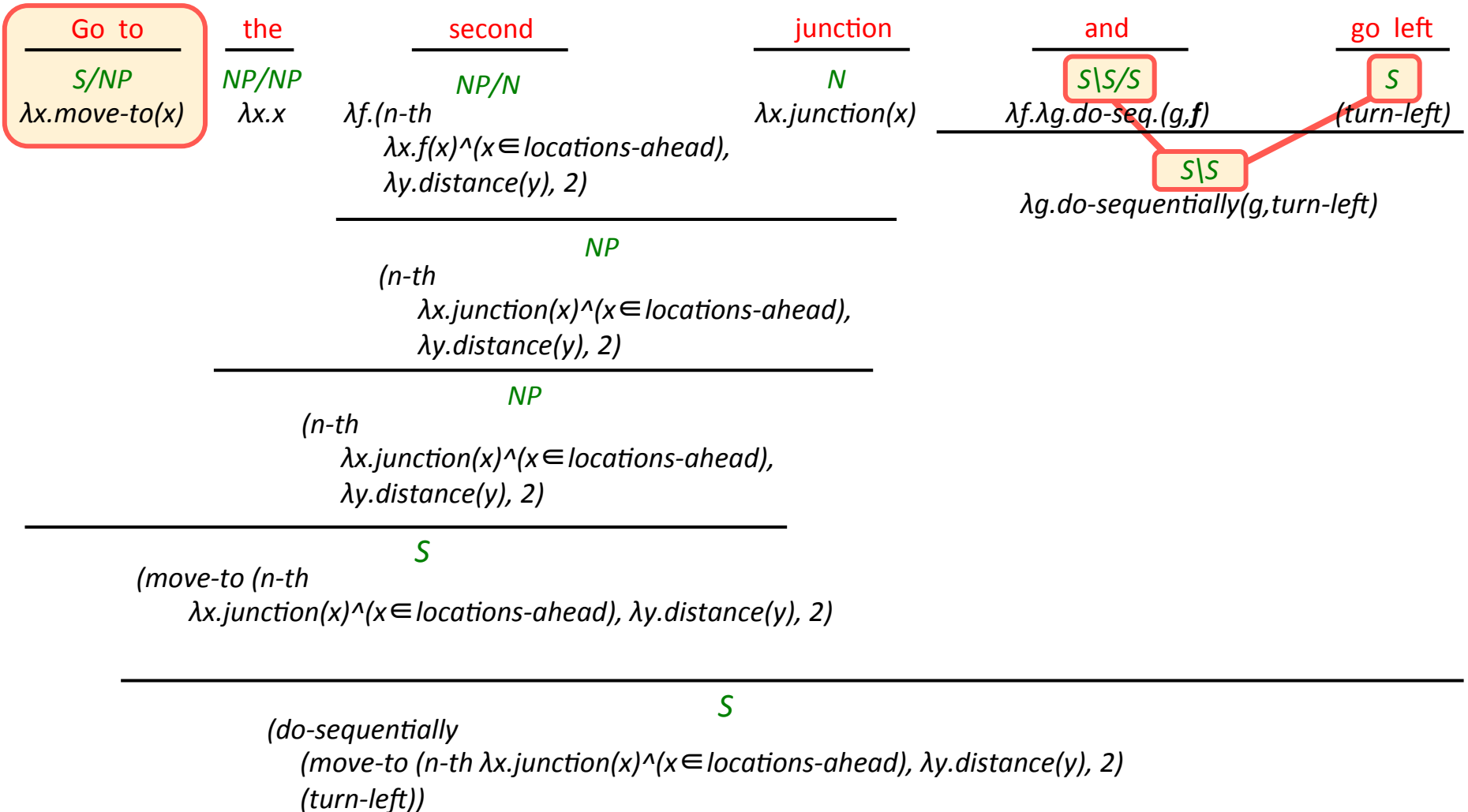
- ◆ **Probabilistic CCGs** define a log-linear model over:

$$\begin{array}{l}
 \text{sentence } x \\
 \text{parse } y \\
 \text{logical form } z
 \end{array}
 \quad
 p(y, z \mid x; \theta, \Lambda) = \frac{e^{\theta \cdot \phi(x, y, z)}}{\sum_{y', z'} e^{\theta \cdot \phi(x, y', z')}}$$



# Example CCG Parse

9



# Learning Probabilistic CCGs

10

- ◆ **Input:** Example pairs of sentences and logical forms
- ◆ **Output:** PCCG lexicon and feature weights
- ◆ **Structure learning:** Generate lexical items from examples
  - ◆ Via combination or splitting rules
- ◆ **Data driven updates:** add lexical items only when involved in generating most likely parse of formula
- ◆ Parameter estimation via gradient descent

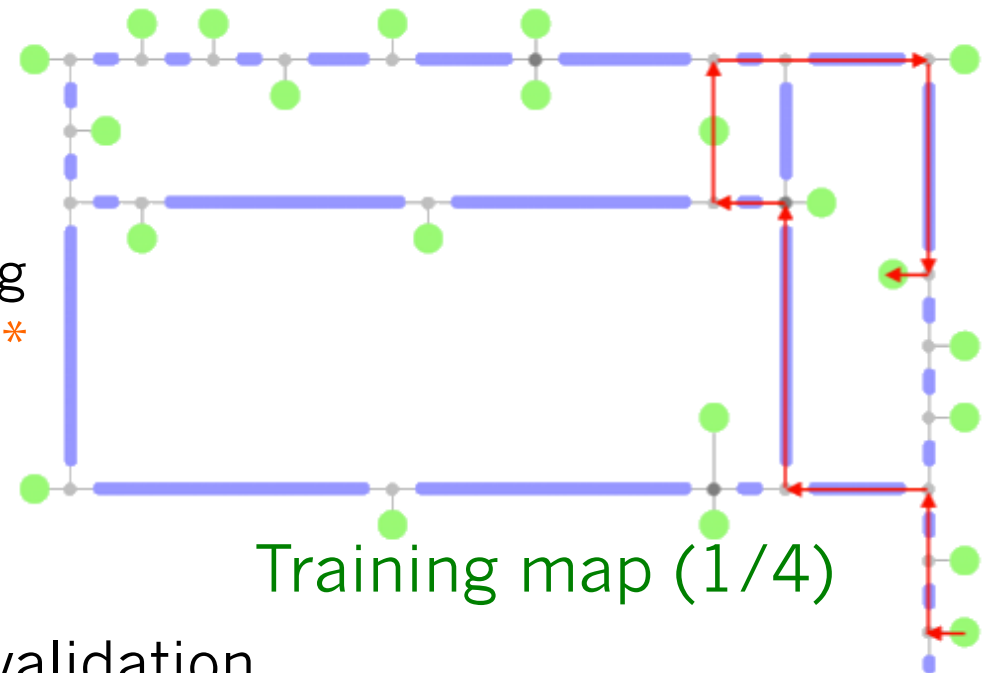
$$\frac{\partial \log(p(z_i | x_i; \theta, \Lambda))}{\partial \theta_j} = E_{p(y|x_i, z_i; \Theta, \Lambda)} [\phi_j(x_i, y, z_i)] - E_{p(y, z|x_i; \Theta, \Lambda)} [\phi_j(x_i, y, z)]$$

# Experimental Setup

11

## ◆ Training

- ◆ Route instructions:
  - ◆ 9 routes, 2 maps
  - ◆ Semantic labeling using Voronoi Random Fields\*
- ◆ Annotated in RCL



## ◆ Testing

- ◆ Parsing: 10-fold cross-validation
- ◆ Navigation
  - ◆ 1200 generated routes, 2 novel maps
  - ◆ Map discovery simultaneous with following RCL program

# Experiment: Parser

12

- ◆ Route instructions from non-expert users
  - ◆ **Segmented** and **annotated** in RCL
- ◆ Parser test: 10-fold cross-validation on parsing
- ◆ Compare produced parses against gold-standard RCL annotations
- ◆ Tests **exact** match only

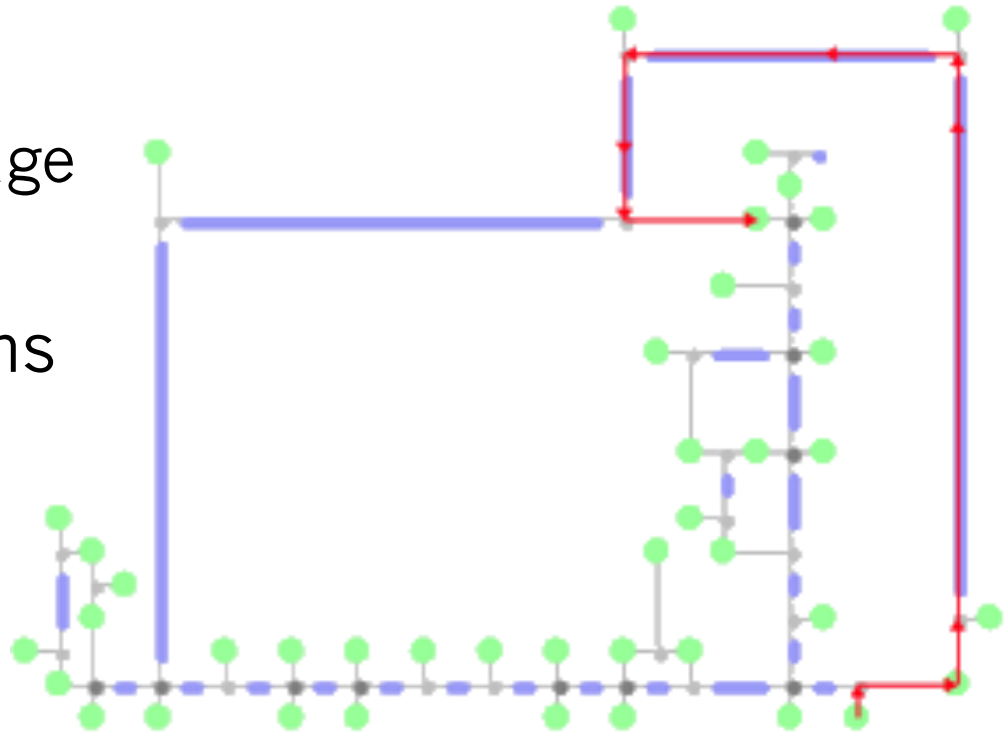
Precision	Recall	F1-measure
71.0%	72.6%	71.8%

- ◆ Evaluates performance on individual sentences, **not** testing full system against a map

# Experiment: Navigation

13

- ◆ Route Following with complex language
  - ◆ Novel route instructions, novel map
    - ◆ 418 sentences total
    - ◆ 25 participants
    - ◆ Complex language represented
  - ◆ Route instructions generated from 2 held-out participants

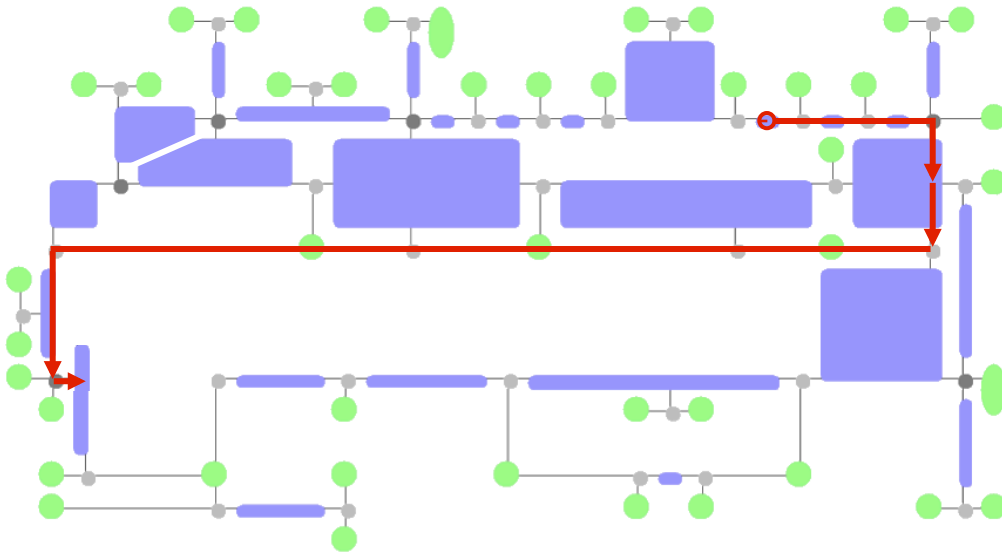




# Example Parse

15

- ◆ Go past two junctions and turn right, go forward to the 3-way intersection, take the first right, **go straight through the second junction then go left**, and turn left again.



```
(do-sequentially
  (do-sequentially
    (do-n-times 2
      (do-sequentially
        (do-until
          (junction current-loc)
          (move-to forward-loc))
        (move-to forward-loc)))
      (turn-right current-loc))
    (do-until
      (junction3 current-loc)
      (move-to forward-loc))
    (turn-right current-loc)
    (do-sequentially
      (do-n-times 2
        (do-sequentially
          (do-until
            (junction current-loc)
            (move-to forward-loc))
          (move-to forward-loc)))
        (turn-left current-loc))
      (turn-left current-loc))
  (turn-left current-loc))
```

# Conclusions

16

- ◆ **It is possible to combine advanced natural language processing with robotic perception and control.**
  - ◆ Parser can be **learned from data** to handle complex, procedural NL for robot instruction
  - ◆ Including counting, loops, conditionals, polysemies
- ◆ Future Work
  - ◆ Local error recovery; more/more varied training data
  - ◆ More interesting data
  - ◆ Generate ranked list of programs to execute
  - ◆ Analyze formal correctness of language; of produced programs
- ◆ Other work extends underpinnings of formal language –  
**ICML 2012**