

Toward Understanding Natural Language Directions

[Video](#)

Motivating Example

- **Robot** Someone is on the way to get you out of here. Are there any other people around who need help?
- **Victim** I saw someone in the main lobby.
- **Robot** Where is the main lobby?
- **Person** Exit this room and turn right. Go down the hallway past the elevators. The lobby is straight ahead.
- **Robot** Understood.

Data Corpus

- Data collection
 - 15 visitors wrote 10 sets of directions each (150 total)
 - Each visitor tries to follow someone else's directions to check quality
 - Best direction giver – 100% followable instructions
 - Worst direction giver – 30% followable instructions
 - Only landmarks shown on predetermined map could be used

With your back to the windows, walk straight through the door near the elevators. Continue to walk straight, going through one door until you come to an intersection just past a white board. Turn left, turn right, and enter the second door on your right (sign says “Administrative Assistant”).

Exploit the Structure of Language

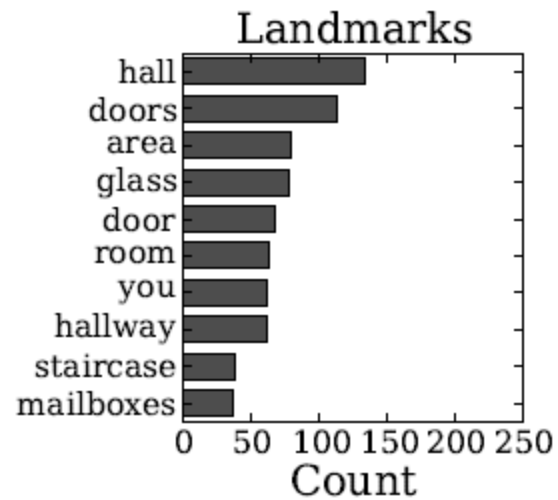
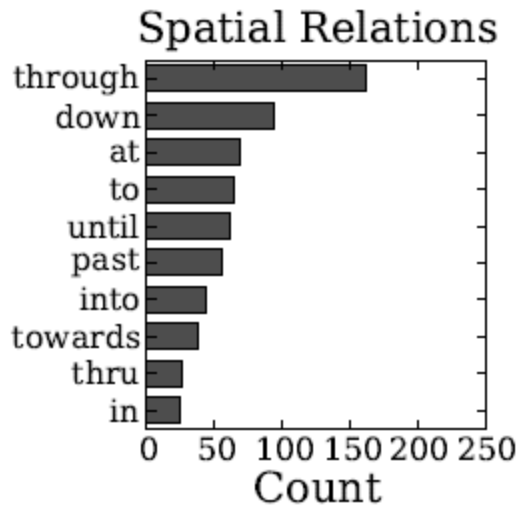
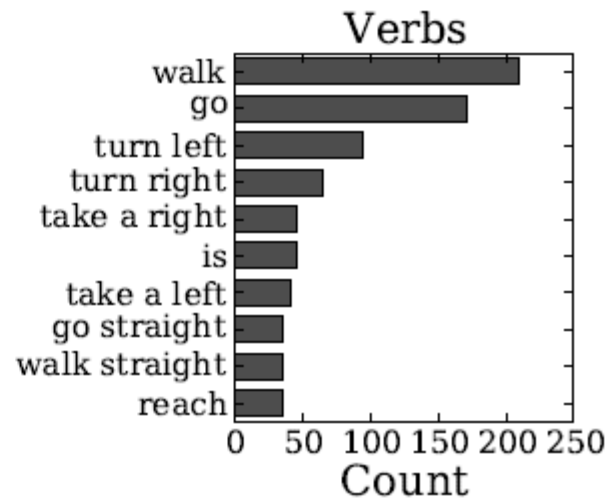
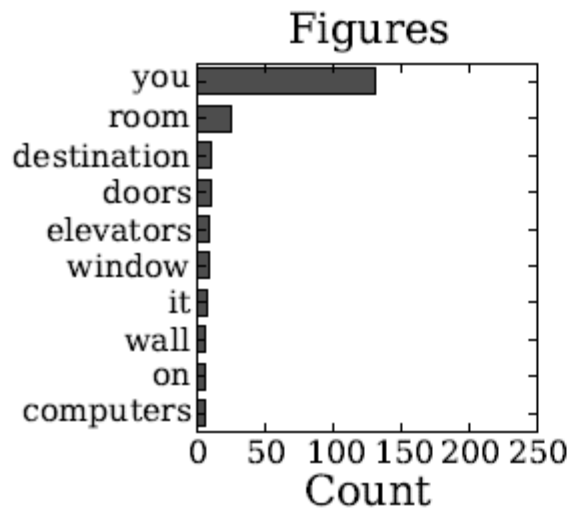
- Directions are:
 - Sequential
 - Contain references to landmarks
 - Contain spatial relations (though, past, etc)
 - Contain verbs

Spatial Description Clause (SDC)

- **figure** (the subject of the sentence)
 - **verb** (an action to take)
 - **landmark** (an object in the environment)
 - **spatial relation** (a geometric relation between the landmark and the figure)
- Any of these fields can be unlexicalized and therefore only specified implicitly.

“[you] Go down the hallway,”

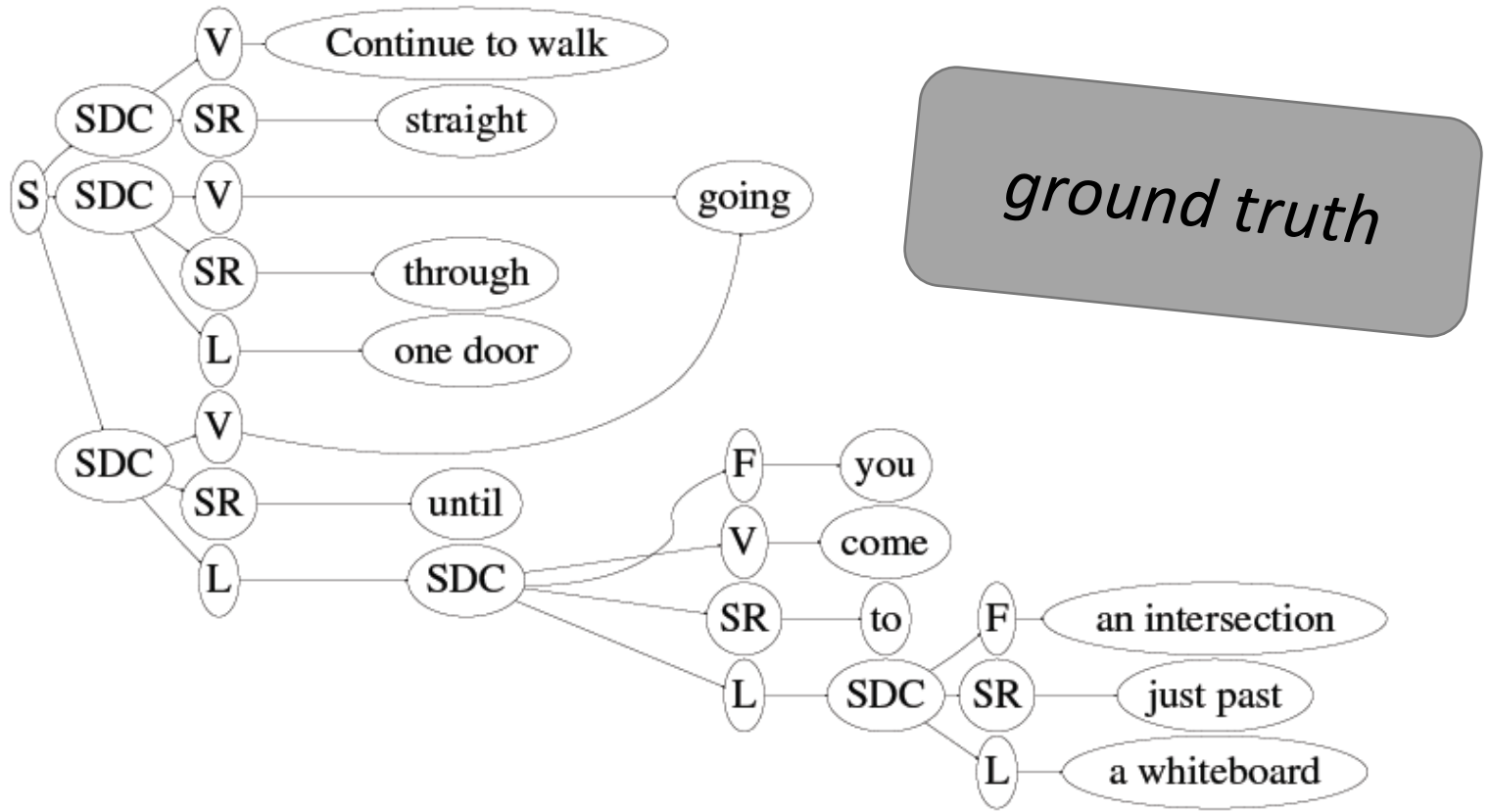
figure *verb* *spatial relation* *landmark*



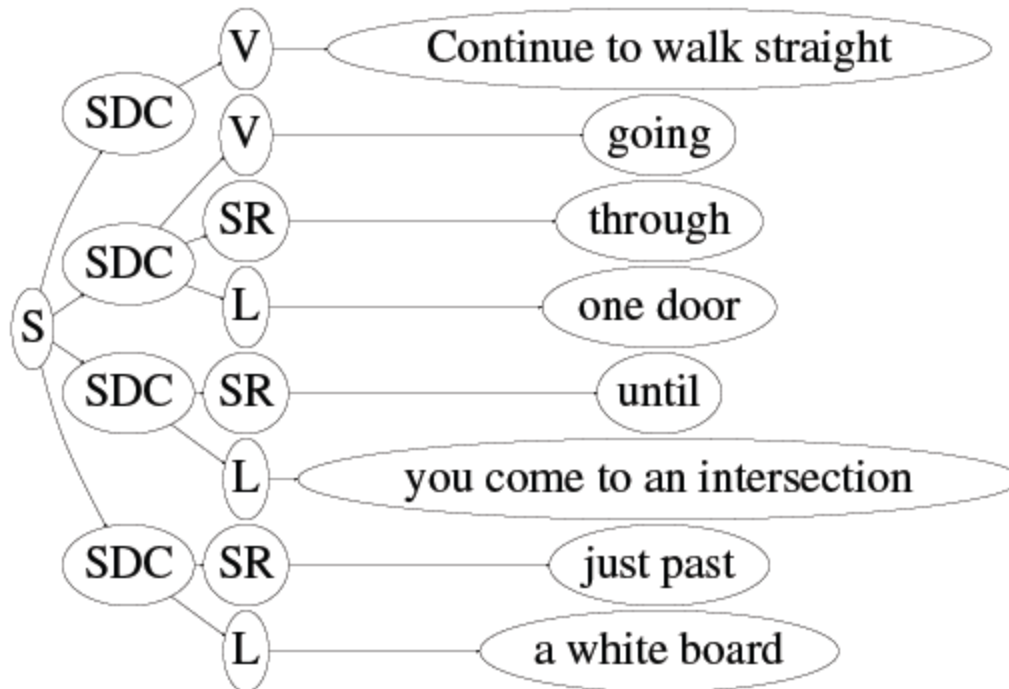
Most frequent words in each SDC field from the corpus if 150 directions (hand annotated).

Process

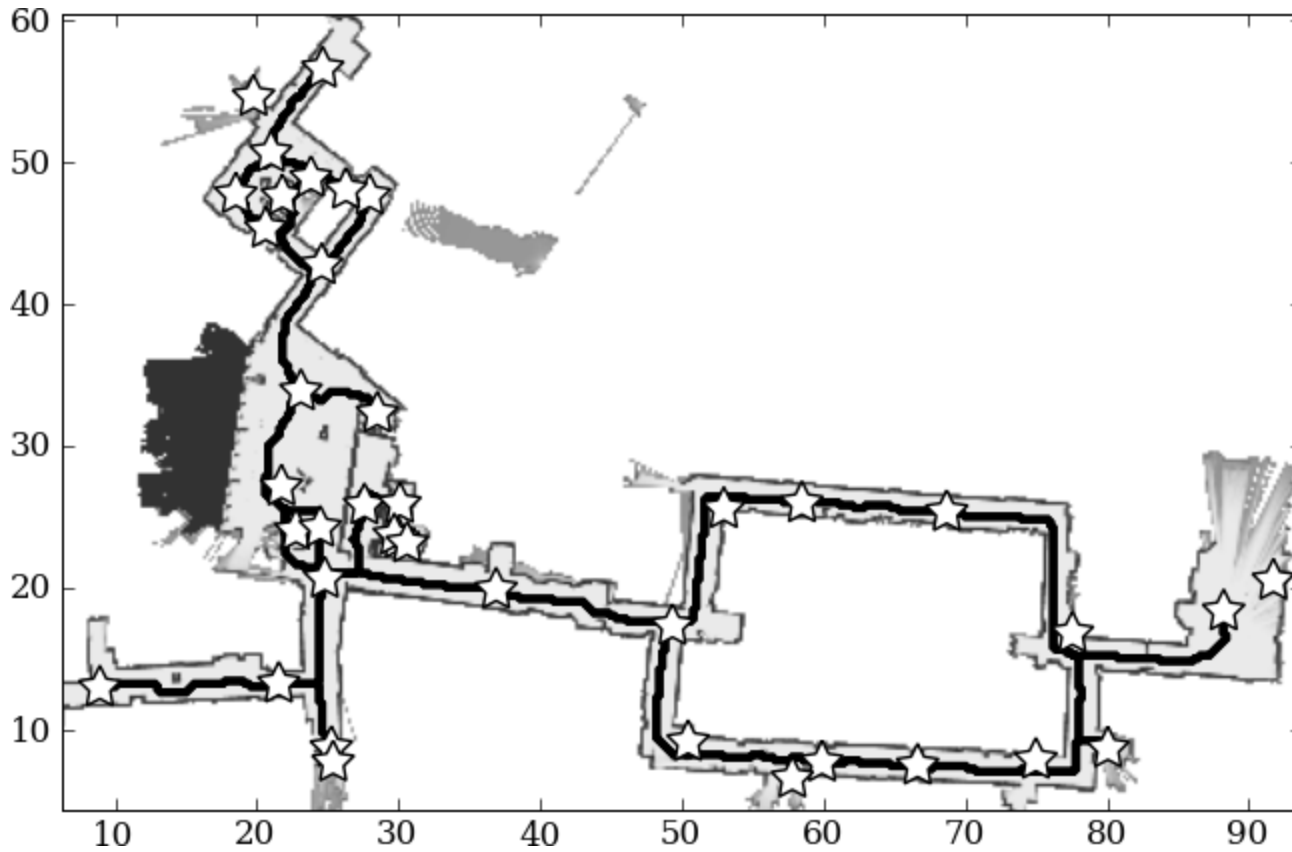
- Automatically extract SDCs from text (CRFs)
- Ground each part in the environment



automatic



Topological Map



$$\arg \max_P p(P, S|O) = p(S|P, O) \times p(P|O)$$

Conditional independence of three disjoint variables: once O is known, knowing S can no longer influence the probability of P .

They add an additional assumption that the path is independent of the objects.
Which leads to:

$$p(P, S|O) \approx p(\text{sdc}_1 \dots \text{sdc}_M | v_1 \dots v_{M+1}, O) \times p(v_1 \dots v_{M+1})$$

$$p(P, S|O) \approx p(\text{sdc}_1 \dots \text{sdc}_M | v_1 \dots v_{M+1}, O) \times p(v_1 \dots v_{M+1})$$

Additional simplifying assumptions (standard Markovian):

- 1) an SDC depends only on the current transition v_i, v_{i+1}
- 2) the next viewpoint v_{i+1} depends only on previous viewpoints.

$$p(P, S|O) = \left[\prod_{i=1}^M p(\text{sdc}_i | v_i, v_{i+1}, O) \right] \times \left[\prod_{i=1}^M p(v_{i+1} | v_i \dots v_1) \right] \times p(v_1)$$

Obtain the probabilities from their labeled training data:

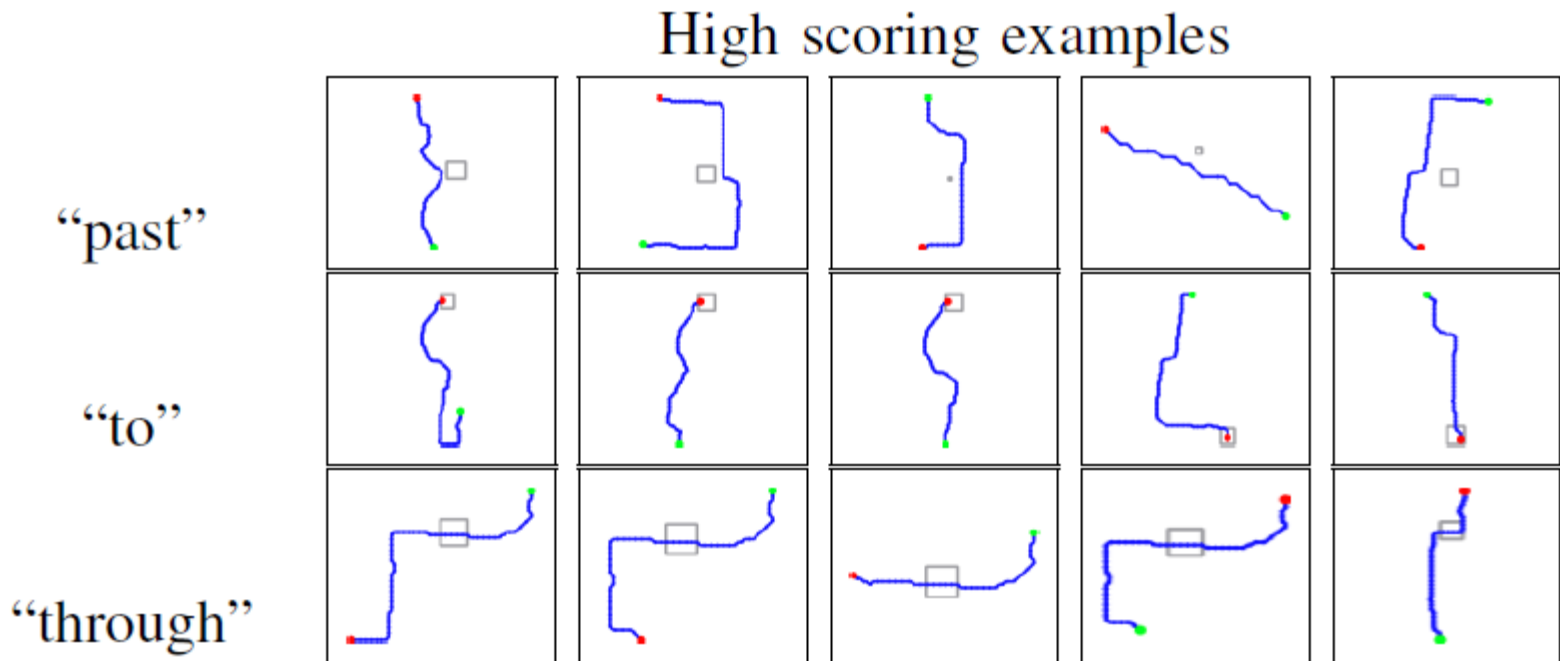
$$\begin{aligned} p(\text{sdc}_i | v_i, v_{i+1}, O) &= p(f_i, a_i, s_i, l_i | v_i, v_{i+1}, O) \\ &\approx p(f_i | v_i, v_{i+1}, o_1 \dots o_K) \times p(a_i | v_i, v_{i+1}) \times \\ &\quad p(s_i | l_i, v_i, v_{i+1}, o_1 \dots o_K) \times p(l_i | v_i, v_{i+1}, o_1 \dots o_K) \end{aligned}$$

Grounding the figures to physical landmarks

- “the door near the elevator”, “a beautiful view of the domes”
- Download >1M images from Flickr
- Used this dataset to model object co-occurrence
 - $P(kitchen|microwave, toaster)$

Grounding spatial relations

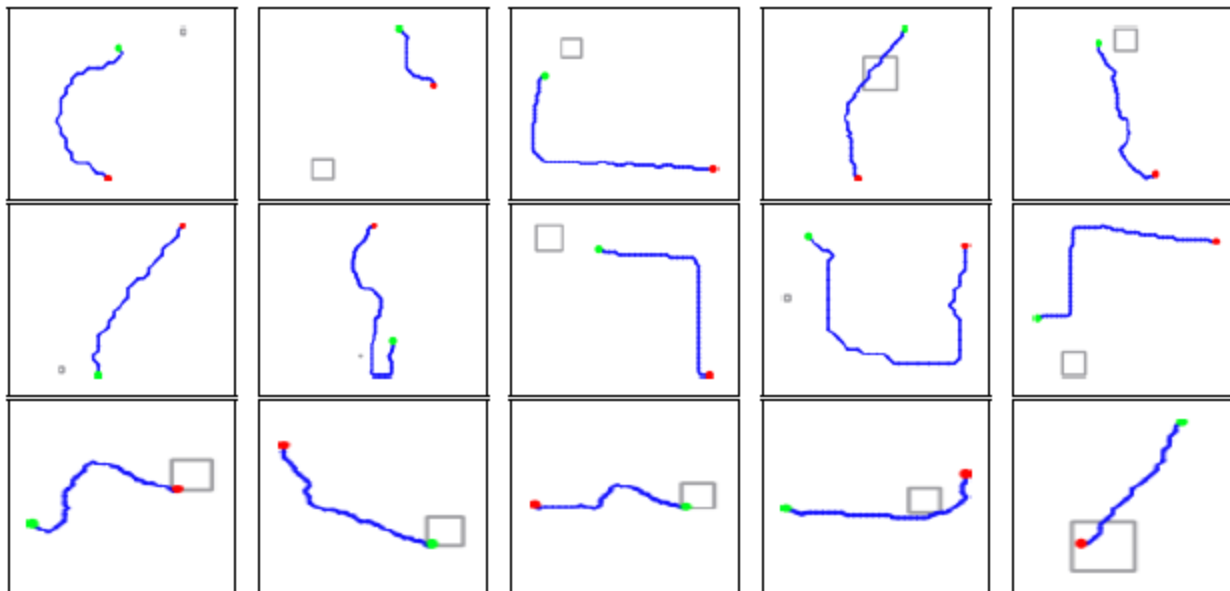
- Hand drawn training examples



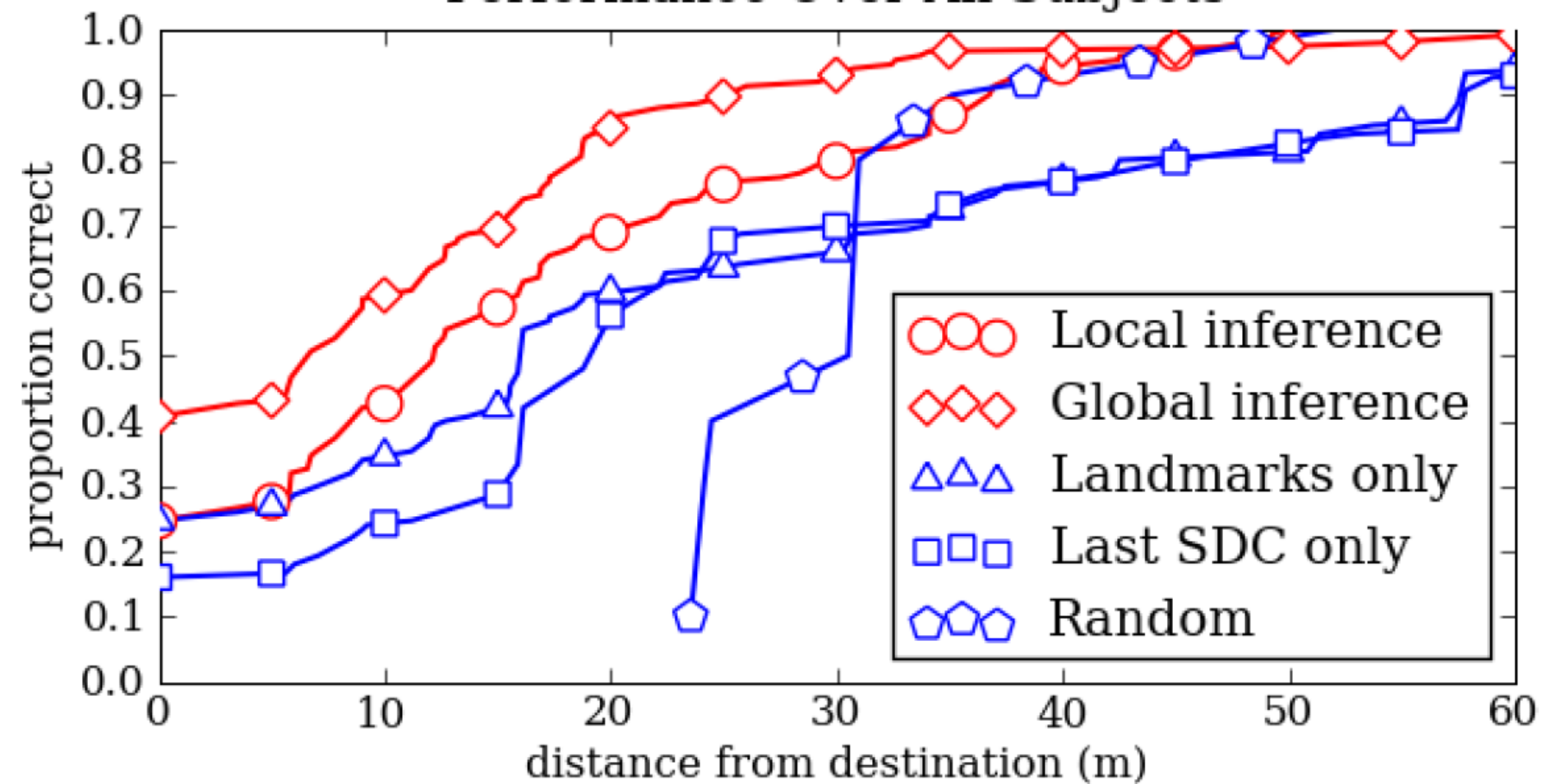
Grounding spatial relations

- Hand drawn training examples

Low scoring examples



Performance Over All Subjects



Evaluation

TABLE I
THE PERFORMANCE OF OUR MODELS AT 10 METERS.

Algorithm	% correct	
	Max Prob	Best Path
Global inference w/spatial relations	48.0%	59.3%
Global inference w/o spatial relations	48.0%	54.7%
Local inference w/ spatial relations	28.0%	42.0%
Local inference w/o spatial relations	26.7%	30.7%
Wei et al. [13]	34.0%	34.0%
Last SDC only	23.0%	24.0%
Random	0.0%	—

Understanding Natural Language Commands for Robotic Navigation and Mobile Manipulation

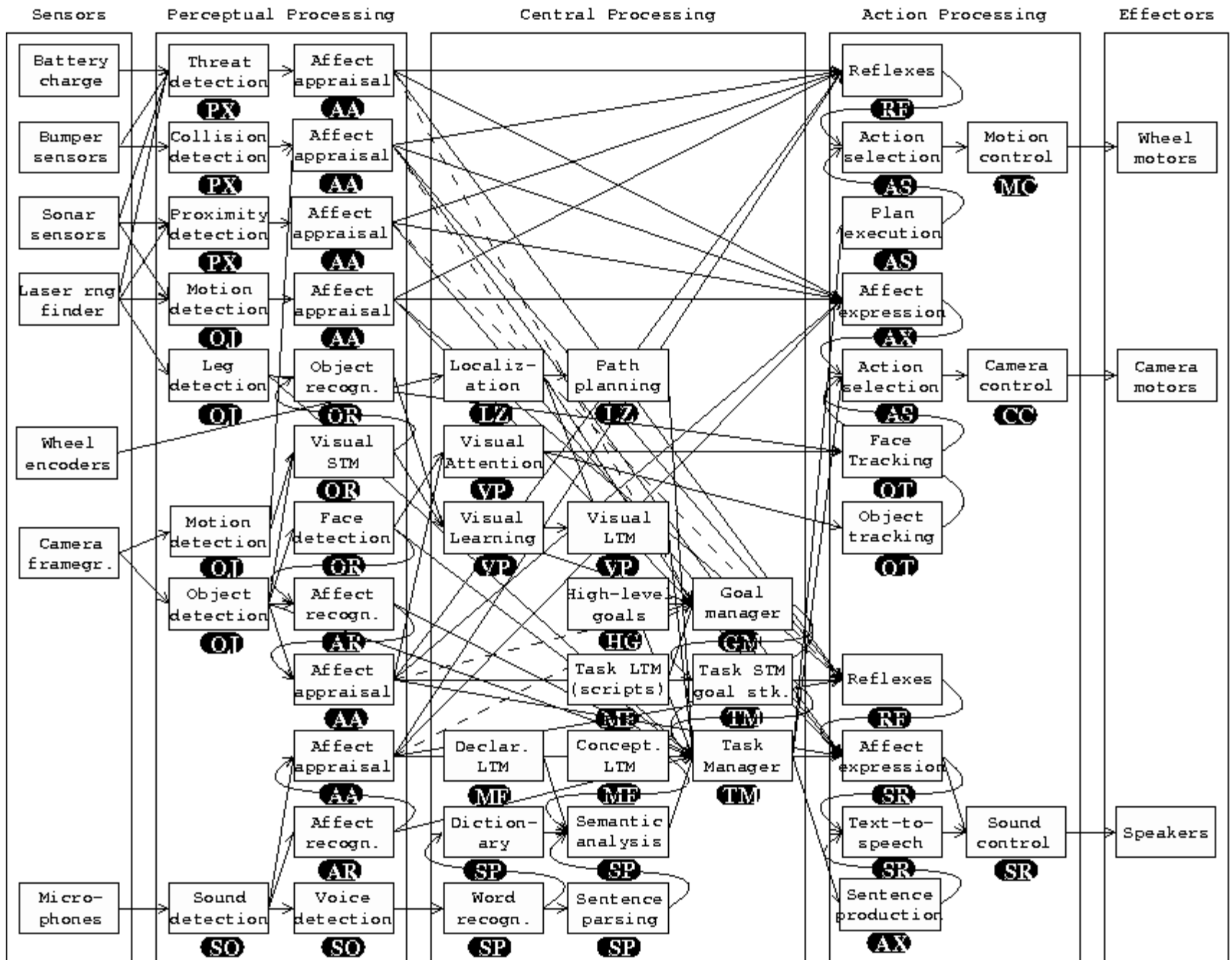


Commands from the corpus

- Go to the first crate on the left and pick it up.
 - Pick up the pallet of boxes in the middle and place them on the trailer to the left.
 - Go forward and drop the pallets to the right of the first set of tires.
 - Pick up the tire pallet off the truck and set it down
-

Other related projects

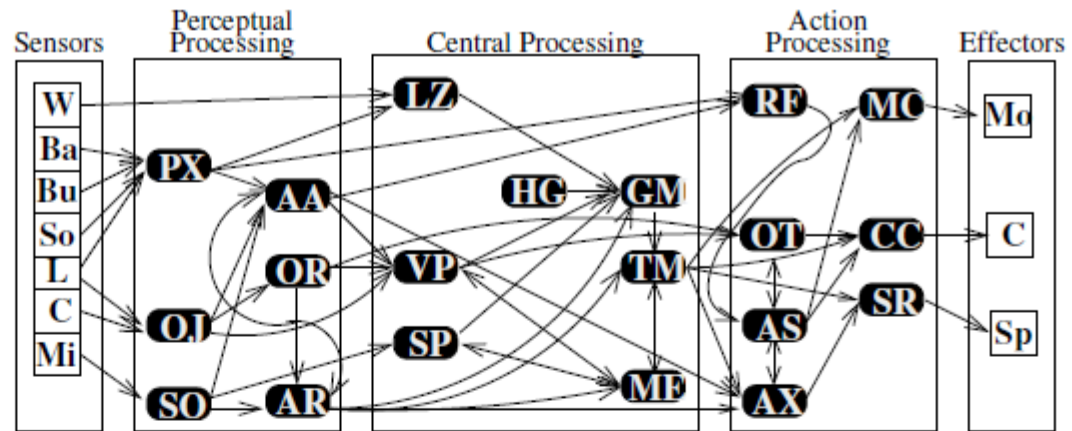
- <http://www.youtube.com/user/HRILaboratory?feature=watch>



- W** – Wheel Encoder
- Ba** – Battery
- Bu** – Bumper Device
- So** – Sonar Device
- L** – Laser Device
- Mo** – Motor Device
- C** – Camera Device
- Mi** – Microphone
- Sp** – Speakers

- – Architectural Link
- PX** – Proximity
- AA** – Affect Appraisal
- OR** – Object Recognition
- OJ** – Object Detection
- SO** – Sound Detection
- AR** – Affect Recognition
- LZ** – Localization
- HG** – High-level Goals
- GM** – Goal Manager
- VP** – Visual Processing
- TM** – Task Manager
- SP** – Speech Processing
- MF** – Memory
- RE** – Reflexes
- MC** – Motion Control
- OT** – Object Tracking
- CC** – Camera Control
- SR** – Speech Production
- AS** – Action Selection
- AX** – Affect Expression

- – ADEServer
- ↔ – Heartbeat Only
- ↔ – Data and Heartbeat
- ◊ – Data Wire
- ↔ – Network



Abstract Agent Architecture

ADE Components

