COMP 790.139 (Fall 2017) Natural Language Processing

Language+Robotics



THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

Mohit Bansal

NLP for Actions/Robotics

Task-based instructions, e.g., navigation, grasping, manipulation, skill learning



NLP for Actions/Robotics

Task-based instructions, e.g., navigation, grasping, manipulation, skill learning





Figure 1: This is an example of a route in our virtual world. The world consists of interconnecting hallways with varying floor tiles and paintings on the wall (butterfly, fish, or Eiffel Tower.) Letters indicate objects (e.g. 'C' is a chair) at a location.



Learning to interpret natural language navigation instructions from observations. Chen and Mooney. AAAI 2011.



Figure 2: Schematic diagram of a map environment and example of semantics of spatial phrases.

Weakly supervised learning of semantic parsers for mapping instructions to actions. Artzi and Zettlemoyer. TACL 2013.



Figure 4: A CCG parse showing adverbial phrases and topicalization.



Place your back against the wall of the "T" intersection. Go forward one segment to the intersection with the blue-tiled hall. This interesction [sic] contains a chair. Turn left. Go forward to the end of the hall. Turn left. Go forward one segment to the intersection with the wooden-floored hall. This intersection contains [sic] an easel. Turn right. Go forward two segments to the end of the hall. Turn left. Go forward one segment to the intersection containing the lamp. Turn right. Go forward one segment to the empty corner.

Figure 1: An example of a route instruction-path pair in one of the virtual worlds from MacMahon, Stankiewicz, and Kuipers (2006) with colors that indicate floor patterns and wall paintings, and letters that indicate different objects. Our method successfully infers the correct path for this instruction.



Figure 2: Our encoder-aligner-decoder model with multi-level alignment



Figure 4: Visualization of the alignment between words to actions in a map for a multi-sentence instruction.

Method	Single-sent	Multi-sent
Chen and Mooney (2011)	54.40	16.18
Chen (2012)	57.28	19.18
Kim and Mooney (2012)	57.22	20.17
Kim and Mooney (2013)	62.81	26.57
Artzi and Zettlemoyer (2013)	65.28	31.93
Artzi, Das, and Petrov (2014)	64.36	35.44
Andreas and Klein (2015)	59.60	—
Our model (vDev)	69.98	26.07
Our model (vTest)	71.05	30.34

 Table 1: Overall accuracy (state-of-the-art in bold)

	Full Model	High-level Aligner	No Aligner	Unidirectional	No Encoder
Single-sentence	69.98	68.09	68.05	67.44	61.63
Multi-sentence	26.07	24.79	25.04	24.50	16.67

Table 2: Model components ablations

Table 3: Accuracy as a function of distance from destination

Distance (d)	0	1	2	3
Single-sentence Multi-sentence	$71.73 \\ 26.07$	$86.62 \\ 42.88$	$92.86 \\ 59.54$	$95.74 \\ 72.08$



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

(a) Robotic forklift

(b) Sample commands

Figure 1: A target robotic platform for mobile manipulation and navigation (Teller et al., 2010), and sample commands from the domain, created by untrained human annotators. Our system can successfully follow these commands.

$$\begin{split} EVENT_1(r = \text{Put}, \\ l = OBJ_2(f = \text{the pallet}), \\ l2 = PLACE_3(r = \text{on}, \\ l = OBJ_4(f = \text{the truck}))) \end{split}$$

(a) SDC tree

$$\begin{split} EVENT_1(r = \text{Go} \\ l = PATH_2(r = \text{to}, \\ l = OBJ_3(f = OBJ_4(f = \text{the pallet}), \\ r = \text{on}, \\ l = OBJ_5(f = \text{the truck})))) \end{split}$$









Figure 3: (a) SDC tree for "Go to the pallet on the truck." (b) A different induced factor graph from Figure 2. Structural differences between the two models are highlighted in gray.

Understanding Natural Language Commands for Robotic Navigation and Mobile Manipulation. Tellex, Kollar, Dickerson, Walter, Banerjee, Teller, and Roy. AAAI 2011.



(a) Object groundings

(b) Pick up the pallet



Figure 4: A sequence of the actions that the forklift takes in response to the command, "Put the tire pallet on the truck." (a) The search grounds objects and places in the world based on their initial positions. (b) The forklift executes the first action, picking up the pallet. (c) The forklift puts the pallet on the trailer.



Fig. 1. An illustration of the robot trajectory $\mathbf{x}(t)$ generated from planning constraints that were inferred from the natural language instruction "move near the red box and the blue crate" using the Distributed Correspondence Graph (DCG) model. The dark gray, light gray, and white regions represent the goal states, admissible states, and inadmissible states respectively. The variables $o_1 \dots o_4$ identify the four objects in the environment model.



Fig. 2. A parse tree for the sentence "move near the red box and the blue crate". Part-of-speech tags in the parse tree are from the Penn Treebank [2].



Fig. 3. The factor graph resulting from the parse tree in Figure 2, used by the G^3 algorithm to infer the groundings of the instructions. Each linguistic is grounded to a object, location, or action through a factor that incorporates the grounding of its children. Black boxes, white spheres, and gray sphere are factors, known random variables, and unknown random variables respectively.



(a) "go to the blue (b) "move towards the (c) "travel to the orbox" green object" ange object"

Fig. 7. Images of labeled trajectories generated by constraint and environment sampling that form the training and test sets for constraint inference evaluation.

Utterance



Move the block that is currently located closest to the top left corner to the bottom left of the table, slightly higher than the block in the bottom right corner.

Error:

7.29 Block lengths



Move the block closest to the top left corner so it is above half a block length to the right of the blocks near the lower left corner of the table.

Error:

0.94 Block lengths

Table 6: Above are two commands and the worlds they applyto. Below we see the prediction error of our best model.



Figure 2: Our models all follow the above architecture. 1-Hot word vectors (orange) are fed as input to a Feed-Forward or Recurrent Neural Network for encoding. A semantic representation is extracted (green), which in conjunction with knowledge of the world (blue) is grounded to predict an action.

		MN Sc	NIST Pa	attern T	is with arget	label	ed b	locks	Ran So	dom Pa	attern T	is with arget	blanl	k blo	ocks
		Med	Mean	Med	Mean	S	R	D	Med	Mean	Med	Mean	S	R	D
	Human Performance Oracle	0.00	0.00	0.21 0.00	0.53 0.45	100 100 1	L00 1	L00	0.00	0.30	0.37 1.00	1.39 1.09	93 100 :	100	100
FFN	Discrete Predictions Continuous Predictions End-to-End	0.00 0.49 0.02	0.49 1.00 0.38	1.09 1.59 1.14	2.17 2.42 1.81	93	69	63	5.28 4.25 3.45	5.09 4.04 3.52	5.51 3.86 3.60	5.46 3.93 3.94	9	15	32
RNN	Discrete Predictions Continuous Predictions End-to-End	0.00 0.47 0.03	0.14 0.64 0.19	0.00 1.23 0.53	0.98 1.60 1.05	98	92	78	5.29 4.16 3.29	5.00 4.05 3.47	5.51 3.71 3.60	5.57 3.87 3.70	10	7	46
	Center Baseline Random Baseline	_ 6.37	- 6.49	3.46 6.12	3.43 6.21	100 5	5	11	_ 4.90	_ 4.97	4.09 5.51	4.06 5.44	100 10	11	12

Table 4: Model error when trained on only the subset of the data with decorated blocks or blank blocks. Where appropriate S, R, and D are the model's predictive accuracy at identifying the Source, Reference and Direction. All models are evaluated on the Median and Mean prediction error the source block and its final target location. Distances are presented in block-lengths.



Figure 1: An example of the configuration instruction understanding task (based on blank-labeled blocks). Our model is able to correctly predict the source block and the target position in this case.



Figure 2: Our overall model for the assembly instruction understanding task, showing instruction and world representation learning, language-to-block alignment modules, and source and target (expectation vs. sampling) loss functions.

Model	SOURCE			TARGET		
	Accuracy	Median	Mean	Median	Mean	
End-to-End FFN (Bisk, Yuret, and Marcu 2016)	9.0%	3.45	3.52	3.60	3.94	
End-to-End RNN (Bisk, Yuret, and Marcu 2016)	10.0%	3.29	3.47	3.60	3.70	
Our Expectation Model	56.1%	0.00	2.21	2.78	3.07	
Our Sampling Model	56.3%	0.00	2.18	3.12	3.18	
Our Expectation Model w/ Ensemble	56.6%	0.00	2.12	2.65	2.91	
Our Sampling Model w/ Ensemble	56.8%	0.00	2.11	2.71	2.90	

The box in the bottom right, slightly right of center, moves one space north of the tower.



Take the leftmost front block and place it on top of the stack of two blocks furthest to the back.

Move the highest block down to below and in front of the right stack of blocks.



Take the block from the last row and hide it behind the tower.



Positive Examples

Move the block closest to the bottom left corner so that it is on top of the block at the top of the backwards L .



Slide the block left of the two in the top right over and on top of the block in front of the tower.

The box next to the Tetris structure moves two spaces left and one half up.



In the 3-piece-long line, the middle box takes a second story from the middle box in the top row.



Negative Examples



Fig. 2 The human interaction with the BakeBot system for recipe execution. First the person provides the plain-text recipe and the measured ingredients. Then BakeBot infers a sequence of baking primitives to execute that correspond to following the recipe. If BakeBot encounters an unsupported baking primitive, it asks its human partner for help executing the instruction. The end result is baked cookies.

Interpreting and Executing Recipes with a Cooking Robot. Bollini, Tellex, Thompson, Roy, Rus. ISER 2012.



Fig. 3 Architecture of the BakeBot system. The NL system processes the plain text recipe, producing a high-level plan which is sent to the robot. For each instruction in the high-level plan, the motion planner assembles a motion plan and executes it on the PR2 robot.

http://projects.csail.mit.edu/video/research/robo/bakebot_final.mp4

Interpreting and Executing Recipes with a Cooking Robot. Bollini, Tellex, Thompson, Roy, Rus. ISER 2012.

Recipe Text	Inferred Action Sequence
Afghan Biscuits 200g (7 oz) butter 75g (3 oz) sugar 175g (6 oz) flour 25g (1 oz) cocoa powder 50g cornflakes (or crushed weetbix)	
Soften butter. Add sugar and beat to a cream. Add flour and cocoa. Add cornflakes last. Put spoonfuls on a greased oven tray. Bake about 15 minutes at 180°C (350°F).	<pre>pour(butter, bowl); mix(bowl) pour(sugar, bowl); mix(bowl) pour(flour, bowl); pour(cocoa, bowl) pour(cornflakes, bowl); mix(bowl) scrape() preheat(350); bake(pan, 20)</pre>

Fig. 4 Text from a recipe in our dataset, paired with the inferred action sequence for the robot.



Interpreting and Executing Recipes with a Cooking Robot. Bollini, Tellex, Thompson, Roy, Rus. ISER 2012.

Recipes: Tell Me Dave (<u>http://tellmedave.cs.cornell.edu/</u>)



Fig. 1. Natural Language Instructions to sequence of instructions for a given new environment. Our approach takes description in natural language and sequences together robotic instructions that are appropriate for a given environment and task. Note that the NL instructions are often ambiguous, and are incomplete, and need to be grounded into the environment.

Recipes: Tell Me Dave (<u>http://tellmedave.cs.cornell.edu/</u>)



Take some coffee in a cup.

Add ice cream of your choice.

Finally, add raspberry syrup to the mixture.

Fig. 4. **Robot Experiment.** Given the language instruction for making the dessert 'Affogato': '*Take some coffee in a cup. Add icecream of your choice. Finally, add raspberry syrup to the mixture.*', our algorithm outputs a sequence that the PR2 executes to make the dessert. (Please see the video.)

Recipes: RoboBarista (<u>http://robobarista.cs.cornell.edu/</u>)





RoboBarista: <u>http://robobarista.cs.cornell.edu/</u>



Fig. 10. **Examples of transferred trajectories** being executed on PR2. On the left, PR2 is able to rotate the 'knob' to turn the lamp on. In the third snapshot, using two transferred trajectories, PR2 is able to hold the cup below the 'nozzle' and press the 'lever' of 'coffee dispenser'. In the last example, PR2 is frothing milk by pulling down on the lever, and is able to prepare a cup of latte with many transferred trajectories.

Navigation Instruction Generation



Fig. 4. Participants' field of view in the virtual world used for the human navigation experiments.

Output: route instruction

"turn to face the grass hallway. walk forward twice. face the easel. move until you see black floor to your right. face the stool. move to the stool"

Fig. 1. An example route instruction that our framework generates for the shown map and path.

Navigation Instruction Generation



Fig. 2. Our method generates natural language instructions for a given map and path.



Fig. 3. Our encoder-aligner-decoder model for surface realization.

Navigation Instruction Generation



Fig. 8. Examples of paths from the SAIL corpus that ten participants (five for each map) followed according to instructions generated by humans and by our method. Paths in red are those traversed according to human-generated instructions, while paths in green were executed according to our instructions. Circles with an "S" and "G" denote the start and goal locations, respectively.

Navigation Dialogue



Fig. 1. A user gives a tour to a robotic wheelchair designed to assist residents in a long-term care facility. (Left) The guide provides an ambiguous description of the kitchen's location. (Right) When the robot is near one of the likely locations, it asks the guide a question to resolve the ambiguity.

Manipulation Dialogue



(a) Unmerged grounding graphs for three dialog acts. The noun phrases "the pallet," "one" and "the one near the truck" refer to the same grounding in the external world but initially have separate variables in the grounding graphs.



(b) The grounding graph after merging γ_2 , γ_3 and γ_5 based on linguistic coreference.

Figure 2. Grounding graphs for a three-turn dialog, before and after merging based on coreference. The robot merges the three shaded variables.

Clarifying Commands with Information-Theoretic Human-Robot Dialog. Deits et al. JHRI 2013.

Manipulation Dialogue



Figure 3. System diagram. Grayed-out blocks show components developed in previous work and are therefore not discussed in detail in this paper; black blocks show the question-asking feedback system new to this paper.

Manipulation Dialogue

Command:	Move your pallet further right.
Question:	What do the words 'your pallet' refer to?
Answer:	Your pallet refers to the pallet you are currently carrying.
Command:	Move closer to it.
Question:	What does the word 'it' refer to?
Answer:	It refers to the empty truck trailer.
Command:	Take the pallet and place it on the one to the left.
Question:	What do the words 'the one' refer to?
Answer:	The one refers to the empty trailer.
Command: Question: Answer:	Place the pallet just to the right of the other pallet. What do the words 'the pallet' refer to? The wooden crate that the merchandise sits on top of.

Figure 4. Sample commands, questions, and answers from the corpus.

Learning New Actions via Dialogue



Figure 1: An example setup and dialogue. Objects are marked with labels only for the illustration purpose.

Back to the Blocks World: Learning New Actions through Situated Human-Robot Dialogue. She, Yang, Cheng, Jia, Chai, Xi. SigDial 2014.

TASK TO COMPLETE

DIRECTORY

Dave Daniel wants the item in slot 5.

ROBOT	How can I help?
YOU	please bring the item in slot 5 to dave daniel
ROBOT	To whom should I bring something?
YOU	dave daniel
ROBOT	What action did you want me to take involving dave daniel?
YOU	5
ROBOT	What should I bring to dave daniel?
YOU	a calander
ROBOT	I'm sorry, but I couldn't pinpoint what you meant by that.
ROBOT	What should I bring to dave daniel?
YOU	a day planner
ROBOT	I'm sorry, but I couldn't pinpoint what you meant by that.
ROBOT	What should I bring to dave daniel?
YOU	calendar
ROBOT	You want me to bring calendar to dave daniel?
YOU	yes
ROBOT	I thought so
ROBOT	Happy to help
YOU	
Final Task	

People:

Alice Ashcraft; Secretary Francis ("Frannie") Foster Robert ("Bob") Brown Carol Clark, PhD Dave Daniel, PhD George Green; Intern Evelyn ("Eve") Eckhart Mallory Morgan; Director Peggy Parker, PhD Walter Ward; Supervisor Items available to robot:



Figure 3: The Mechanical Turk interface for the delivery task. This abridged conversation is from a Turker in training batch 0, when the system had access to only the seed lexicon. Because of this conversation, the agent learned that "calender" and "planner" mean "calendar" during retraining.



Figure 5: Left: Robot platform (Segbot) used in experiments. **Right:** Segbot architecture, implemented using Robot Operating System (ROS).

Learning to Interpret Natural Language Commands through Human-Robot Dialog. Thomason, Zhang, Mooney and Stone. IJCAI 2015.