





Embodied Language Grounding

Katerina Fragkiadaki

Carnegie Mellon University



Localize the referrents it mentions

Modeling Relationships in Referential Expressions with Compositional Modular Networks, Ronghang et al.

Generate the image it describes

Probabilistic Neural Programmed Networks for Scene Generation, Deng et al.



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."

Caption an image

Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Li, Fei fei



- After wading barefoot in the lake, Erik used his shirt to dry his feet.
- After wading barefoot in the lake, Erik used his glasses to dry his feet.

To act upon it, and infer its affordability

Embodied Cognition



- Words and phrases are indexes to objects in the world or to prototypical symbols of those objects
- We derive affordance from those objects
- The derived affordances constrain the way ideas can be coherently combined

To act upon it, and infer its affordability



- The bowl inside the cube
- The cube inside the bowl

To act upon it, and infer its affordability

Simulation Semantics



- Words and phrases are indexes to objects in the world or to prototypical symbols of those objects
- We derive affordance from those objects
- The derived affordances constrain the way ideas can be coherently combined

To act upon it, and infer its affordability

Simulation Semantics



Many philosophers agree and support such simulation semantics, no computational model for language grounding

Visually grounded interaction and language, NIPS workshop 2018

Language grounding to visual cues

Remains disconnected from affordances

2D boxes or 2D CNN features do not have any affordability attached Are themselves ungrounded

- The bowl inside the cube
- The cube inside the bowl

Reward learning using natural language

Given a NL utterance, learn a visual detector

"Can is to the right of the mug"



Reward learning using natural language

Given a NL utterance, learn a visual detector

"Can is to the right of the mug"



Use the learned visual detector to guide policy learning for achieving the NL described goal

It did not really worked, the reward detector could not effectively generalize across camera placements

3. Language Grounding from Narrated Demonstrations CVPR'18

"Can is to the right of the mug"



reward detector

Learned reward detector



can is to the right of the book



Learned policy



3. Language Grounding from Narrated Demonstrations CVPR'18



Prior work used manually coded rewards

We ground rewards to the sensory input of the agent

Affordandable visual representations

We seek visual cues that obey basic common sense constraints and basic affordability reeasoning:

- Objects have 3D extent
- Objects do not interpenetrate in 3D
- Objects come in regular sizes
- Objects persist over time

Grounding language on such representations would be able to support affordability inference

Affordandable visual representations

We seek visual cues that obey basic common sense constraints and basic affordability reeasoning:

- Objects have 3D extent
- Objects do not interpenetrate in 3D
- Objects come in regular sizes
- Objects persist over time

Grounding language on such representations would be able to support affordability inference





egomotion-stabilized update





3D-to-2D mapping





geometry-aware RNN



2D RNN [1]



[1] Neural scene representation and rendering DeepMind, Science, 2018



Truly novel scenes

[1] Neural scene representation and rendering DeepMind, Science, 2018

Geometry-Aware Recurrent Networks (GRNNs)



3D object detection



3D object detection

of input views





3D object detection Results - 3D object detection

of input views



3D object detection





- ask too much: high level of 3D details may be impossible to obtain
- ask too little: information about semantics of the objects is not prese



1. We consider an embodied agent that can see a scene from multiple viewpoints



1. We consider an embodied agent that can see a scene from multiple viewpoints



1. Our agent learns to map an RGB image to a set of 3D feature maps by training GRNNs to predict views



1. Our agent maps noun phrases to object-centric 3D feature maps



1. Our agent maps noun phrases to object-centric 3D feature maps



1. Our agent maps spatial expressions to relative 3D offsets



1. Our agent populates a 3D canvas with the predicted object tensors adn their relative offsets



1. Our agent populates a 3D canvas with the predicted object tensors adn their relative offsets



1. Our agent populates a 3D canvas with the predicted object feature maps and their relative spatial offsets



1. The generated canvas when projected should match the RGB image views

Red Rubber Cylinder to the left front of Blue Rubber Cube to the left front of Green Rubber Cylinder to right front of Blue Rubber Cube

Red Rubber Cube to the left front of the Blue Rubber Sphere to the right front of Cyan Metal Cylinder



Neural render

Blender render



Natural language utterance

Neural render



Purple Cylinder to the left behind

of Brown Cube to the left front of

Purple Cylinder to the left behind of Cyan Cube to the left front of Cyan Cube





Blender render

Natural language utterance Cyan Cube to the left behind of Gray Sphere to the left front of Blue Cube

Neural render

Blender render



Red Sphere to the left behind of Cyan Cylinder to the left front of Red Sphere



Natural language utterance

"red cylinder to the right behind of green cube"

Neural render

Blender render



"pink cylinder to the left front of red cylinder"



Scene alteration

Natural language utterance

"blue sphere to the right behind of green cube"

Neural render

Blender render



"green cube to the left front of cyan cylinder"



Affordability Inference

Natural language utterance

"blue sphere to the right behind of green cube"

Neural render

Blender render



"green cube to the left front of cyan cylinder"



Grounding Language on 3D visual feature representations

- Objects have regular sizes: object size is disentangled from the camera viewpoint
- Objects have 3D extent
- Objects do not interpenetrate in 3D: during iterative scene generation we can detect 3D intersection and continue sampling valid configurations
- Objects persist over time

Next steps

- Grounding action descriptions
- Use intuitive physics and dynamics beyond static spatial constraints



Thank you







Mihir

Fish Tung

Jyed Javed

Adam Harley

Max Sieb

- Embodied language grounding, arxiv
- Reward Learning from Narrated Demonstrations, Tung et al., CVPR 2018