

# Natural Language Processing and Robotics

Jesse Thomason

[ <http://glamor.rocks/> ]



# Streetlight Effect

A policeman sees a drunk man searching for something under a streetlight and asks what the drunk has lost. He says he lost his keys and they both look under the streetlight together. After a few minutes the policeman asks if he is sure he lost them here, and the drunk replies, no, and that he lost them in the park. The policeman asks why he is searching here, and the drunk replies. "this is where the light is".



# GLAMOR Lab

- Grounding Language in Actions, Multimodal Observations, and Robots

## PhD Students



Bill Zhu  
2<sup>nd</sup> year PhD student (co-Robin Jia)  
📖 Compositional Semantics  
🚶 Vision-Language Navigation



Tejas Srinivasan  
2<sup>nd</sup> year PhD student  
📖 Vision-Language Alignment  
🚶 VL Continual Learning



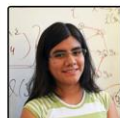
Anthony Liang  
2<sup>nd</sup> year PhD student  
🚶 Language-guided RL  
🚶 RL for Dialogue



Ting-Yun Chang  
2<sup>nd</sup> year PhD student (co-Robin Jia)  
📖 Pre-trained Model Understanding



Ishika Singh  
2<sup>nd</sup> year PhD student  
🚶 VLN and VL Interaction



Leticia Pinto-Alva  
2<sup>nd</sup> year PhD student  
📖 Vision-and-Language



Lee Kezar  
3<sup>rd</sup> year PhD student  
📖 VL for ASL



Abrar Anwar  
2<sup>nd</sup> year PhD student  
🚶 Speech and Gesture  
🚶 Assistive Robotics

## Undergraduate and Masters Students



Elle Szabo  
Senior Undergraduate @ USC  
Computer Science



Chu Fang  
Senior Undergraduate @ USC  
Computer Science (Games)  
URAP



Flora Jia  
Junior Undergraduate @ USC  
Computer Science and Applied Mathematics  
CURVE Fellowship



Junu Song  
Junior Undergraduate @ USC  
Computer Science  
CURVE Fellowship



Yuan Huang  
2<sup>nd</sup> year Masters Student @ USC  
Electrical and Computer Engineering  
Specializing in Machine Learning and Data Science

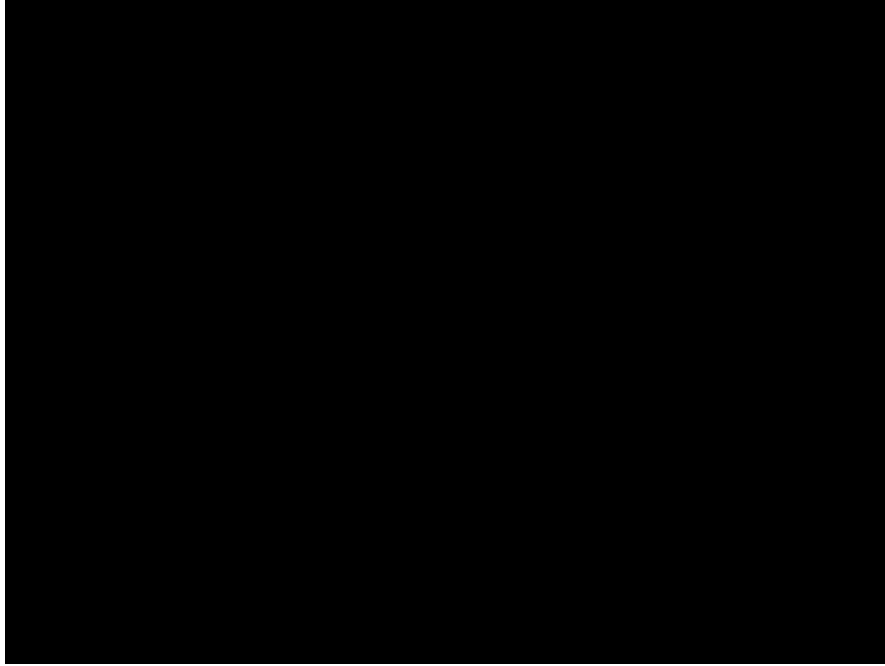


Leslie Moreno  
Sophomore Undergraduate @ USC  
Computer Engineering and Computer Science  
CURVE Fellowship (co-Maja Mataric)



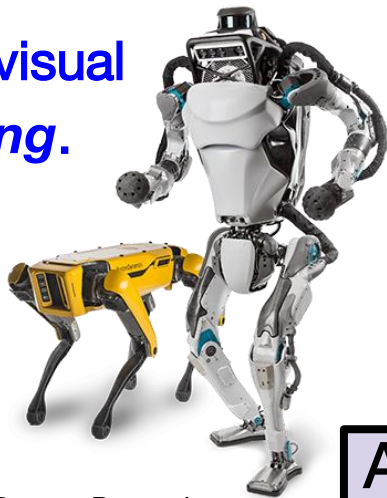
Julie Kim  
Sophomore Undergraduate @ USC  
Computer Engineering and Computer Science  
CURVE Fellowship (co-Maja Mataric)

# Why Would a Robot Need Language?



# Why Aren't Our Robots Language-Guided Already?

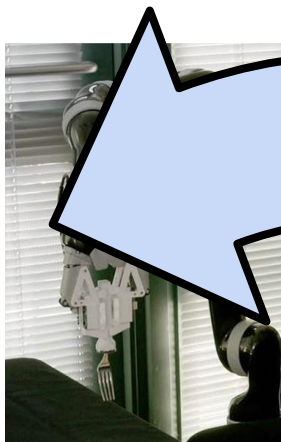
Require visual  
*grounding.*



Boston Dynamics



iRobot



???

No visual  
connections.



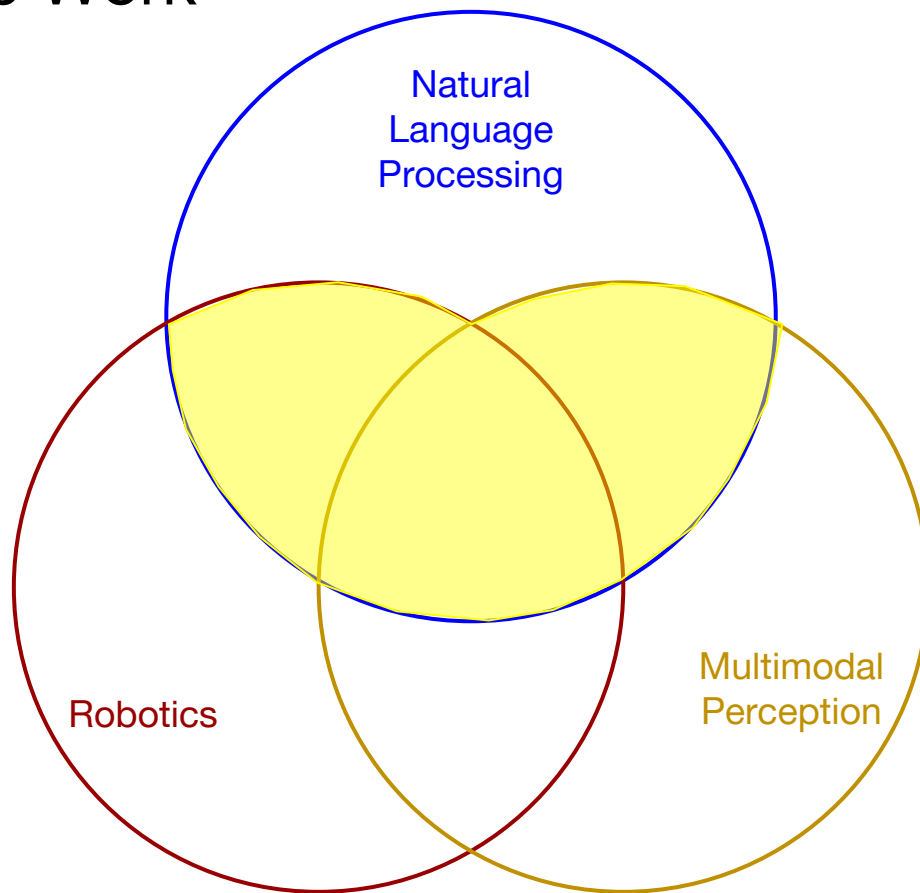
Amazon



Google

**Approach Summary:**  
Enable robots to perform  
language grounding via  
interaction with humans.

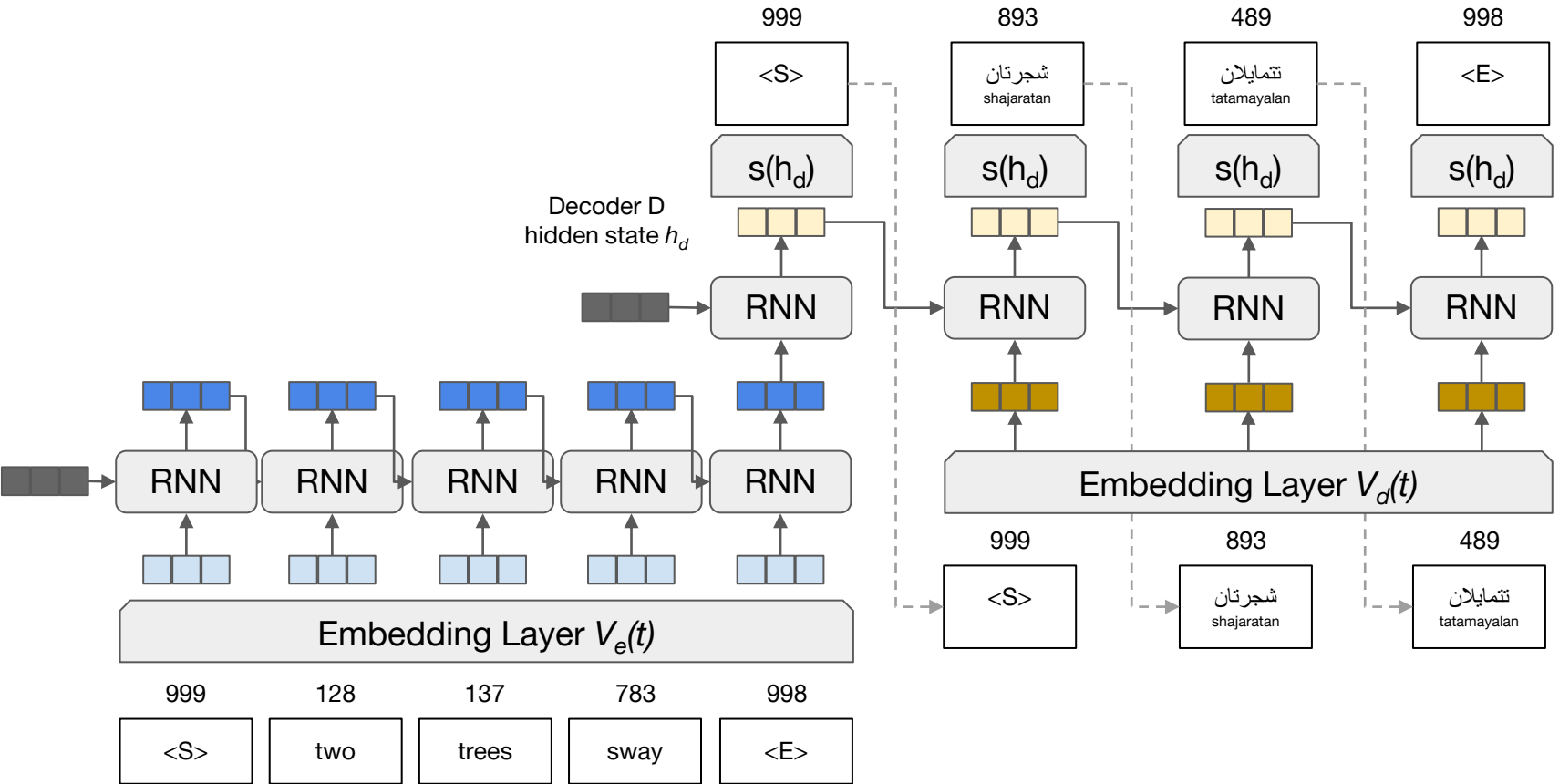
# GLAMOR Lab Work



# Natural Language Processing

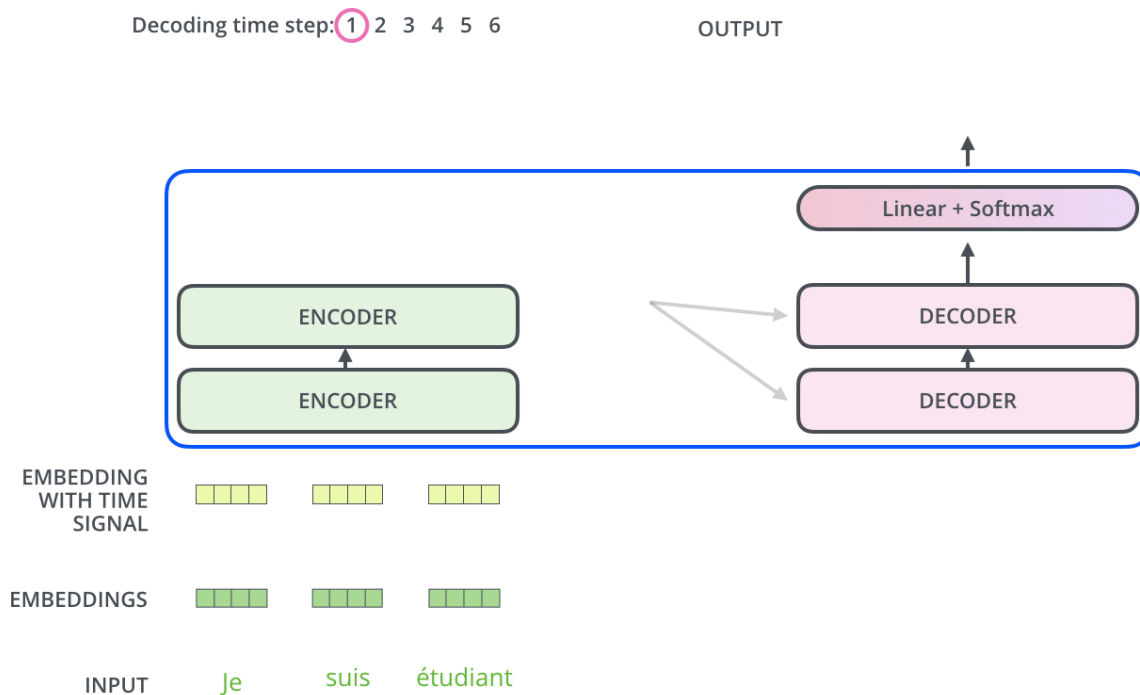
- Concerned with turning “in the wild” human language into something machine readable for downstream tasks like
  - Text classification
    - Sentiment analysis, fake news detection, hate speech detection, spam filters, ...
  - Textual language modeling
    - Text generation (e.g., GPT-3), autocomplete on your phone, spelling correction suggestions, machine translation, ...

# Machine Translation: Basic Idea

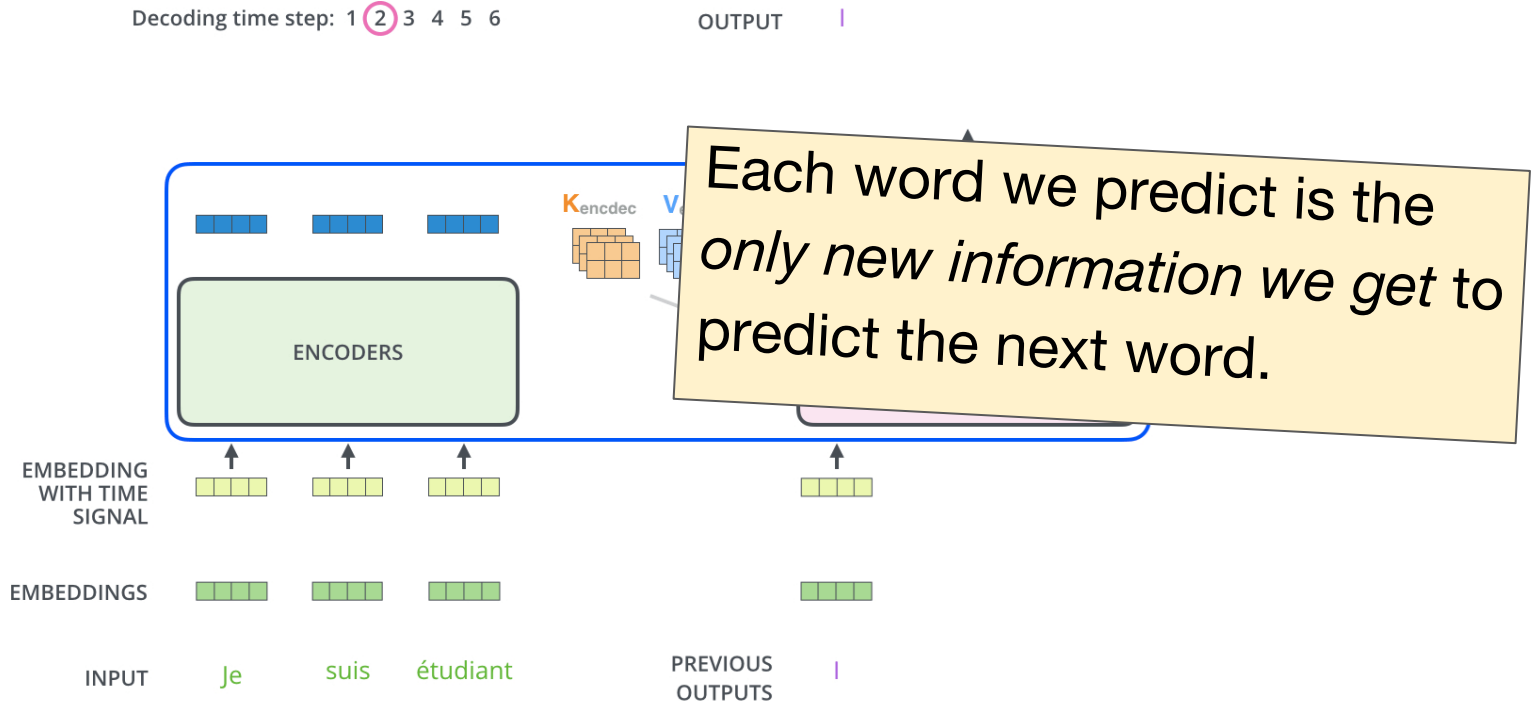




# Machine Translation: Transformer Models

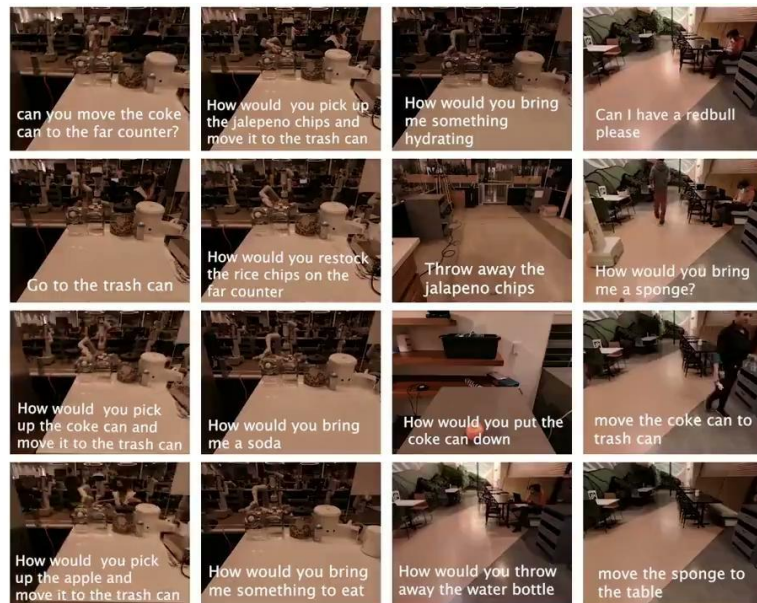


# Machine Translation: Transformer Models



# Machine Translation for Robots

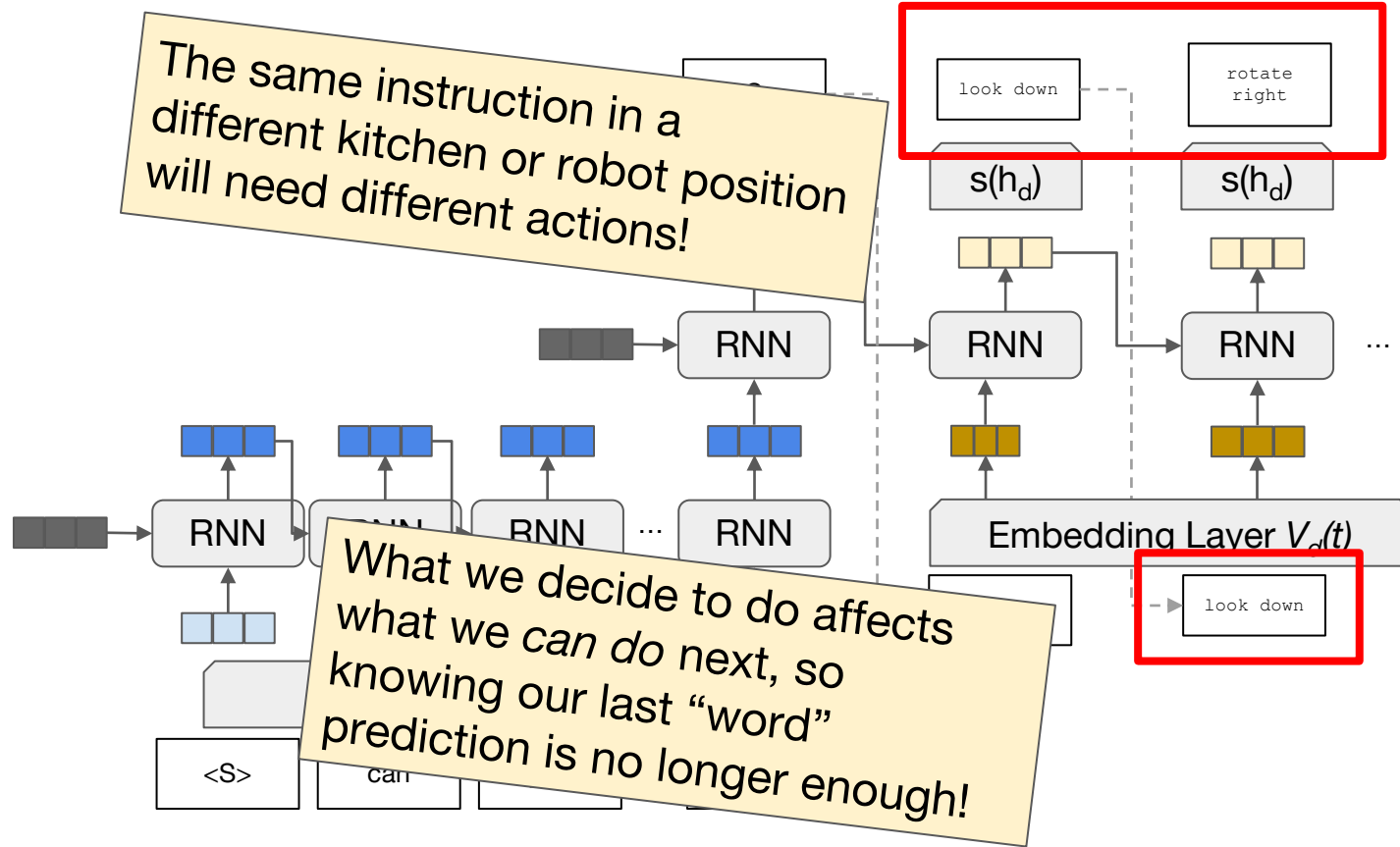
- Source language input:
  - “Can you move the coke can to the far counter?”
- Target language output:
  - look down, rotate right, move forward, move forward, rotate right, grasp coke can, rotate right, move forward, move forward, rotate right, place coke can



# Machine Translation for Robots

- Source language input:
  - “Can you move the coke can to the far counter?”
- Target language output:
  - look down, rotate right, move forward, move forward, rotate right, grasp coke can, rotate right, move forward, move forward, rotate right, place coke can
- How many input/output examples would you need to learn a *policy* that reliably translates English to sequences of robot actions?

# Machine Translation for Robots



# Streetlight Effect

A policeman sees a drunk man searching for something under a streetlight and asks what the drunk has lost. He says he lost his keys and they both look under the streetlight together. After a few minutes the policeman asks if he is sure he lost them here, and the drunk replies, no, and that he lost them in the park. The policeman asks why he is searching here, and the drunk replies, "this is where the light is".



# Experience Grounds Language - World Scopes


**Upshot:**

Most current L operates at the

**Upshot:**

Most L folks playing and vice versa are

**Upshot:**

Maybe this is AI Complete; but so is “language” 

**WS1**

Carefully  
Annotated  
Penn Treebank  
Brown Corpus  
WordNet

**WS2**

Unstructured,  
Unlimited Text  
Common Crawl  
Word2Vec  
ELMo, BERT, GPT\*

**WS3**

Text Paired with  
Sensory Data  
ImageNet, VQA,  
ViLBERT,  
Video Captioning

**WS4**

Language and an Embodied  
Agent  
VLN\*, IQA,  
NL+Games,  
RoboNLP

**WS5**

Social  
Embodiment  
Human-robot  
collaboration  
and dialog

**Claim:**

You can't learn language from the radio

**Claim:**

You can't learn language from a television.

**Claim:**

You can't learn language by yourself.

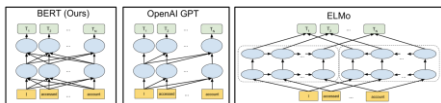
# Experience Grounds Language - World Scopes

Common Crawl



**WS1**

Carefully  
Annotated  
Penn Treebank  
Brown Corpus  
WordNet



**WS2**

Unstructured  
Unlimited Text  
Common Crawl  
Word2Vec  
ELMo, BERT, GPT\*

**WS3**

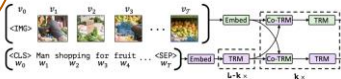
Text Paired with  
Sensory Data  
ImageNet, VQA,  
ViLBERT,  
Video Captioning



By Josh Seltzer  
the trail climbs steadily  
uphill most of the way.



By Daniel Hecker  
the stars in the night sky.



What color are her eyes?  
What is the mustache made of?



**WS4**

Language and an Embodied  
Agent  
VLN\*, IQA,  
NL+Games,  
RoboNLP

Doc:  
The Rebel Enclave consists of jacket  
spider, and wasp. Acolyte, blessed items  
are useful for poison monsters. Star  
Alliance contains bat, panther, and wolf.  
Goblin, jaguar, and lynx are on the same  
team - they are in the Order of the Forest.  
Gleaming and mysterious weapons beat  
cost monsters. Lightning monsters are  
weak against Grandmaster's and  
Soldier's weapons. Fire monsters are  
defeated by fanatical and shimmering  
weapons.  
Goal:  
Defeat the Order of the Forest



**WS5**

Social  
Embodiment  
Human-robot  
collaboration  
and dialog





# Images and Text for Retrieval

## Text Query

“A tropical bird  
perches in the jungle.”

## Candidate Images



# Images and Text for Retrieval

## Image Query



## Candidate Captions

“A rabbit sits in the palm of a hand.”

“Men are talking on a basketball court.”

“A tropical bird perches in the jungle.”

“Children play soccer in a field.”

“A white fox is looking at the camera.”

“Sports equipment is staged for a photo.”

# Formulating the Retrieval Problem

Images

$\mathcal{I}$

Captions

$\mathcal{L}$

Scoring Function

$$F : \mathcal{I} \times \mathcal{L} \rightarrow \mathbb{R}$$

# Formulating the Retrieval Problem

$$F : \mathcal{I} \times \mathcal{L} \rightarrow \mathbb{R}$$

$F(\# \text{ red pixels}, \# \text{ word "red"}) = 1 \text{ if } (a > 0, b > 0) \text{ else } 0$

$F(\text{RGB bins}, \text{color word counts}) =$   
sum of  $\#(\text{color word}, \text{bin}) / (\# \text{color word} + \# \text{bin})$  in training data

$F(\text{detected objects}, \text{word counts}) =$   
sum of  $\#(\text{object}, \text{word}) / (\# \text{object} + \# \text{word})$  in training data

$F(\text{image pixels}, \text{token sequence}) =$   
NN trained with contrastive matching loss

# Formulating the Retrieval Problem

Image Features

$$\psi : \mathcal{I} \rightarrow \mathbb{R}^n$$

Caption Features

$$\omega : \mathcal{L} \rightarrow \mathbb{R}^m$$

Scoring Function

$$F_{(\psi, \omega)} : \psi(\mathcal{I}) \times \omega(\mathcal{L}) \rightarrow \mathbb{R}$$

# Feature Extraction for Language

BOW

“A rabbit sits in the palm of a hand.”

a
bird
hand
in
men
on
rabbit
the

⋮

2
0
1
1
0
0
1
0

⋮

“A tropical bird perches in the jungle.”

1
1
0
1
0
0
0
1

⋮

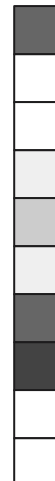
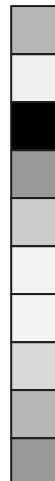
“Men are talking on a basketball court.”

1
0
0
0
1
1
0
0

⋮

# Feature Extraction for Vision

RGB  
kNN bins

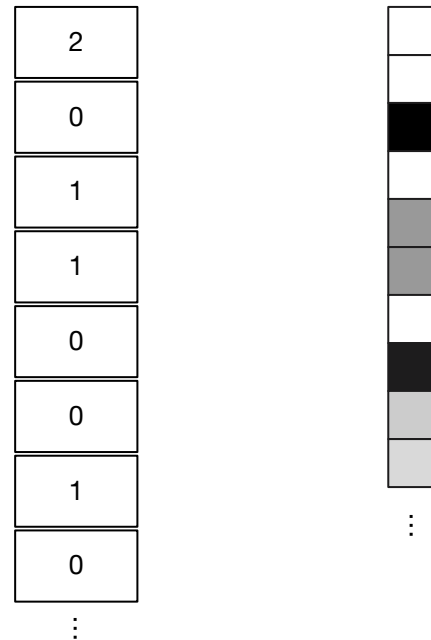


## A Simple Retrieval Solution

$$\psi : \mathcal{I} \rightarrow \mathbb{R}^n$$

$$\omega : \mathcal{L} \rightarrow \mathbb{R}^m$$

$$F_{(\psi, \omega)} : \psi(\mathcal{I}) \times \omega(\mathcal{L}) \rightarrow \mathbb{R}$$





# Contextual Information

*You shall know a word by the company it keeps* (Firth, J. R. 1957:11)

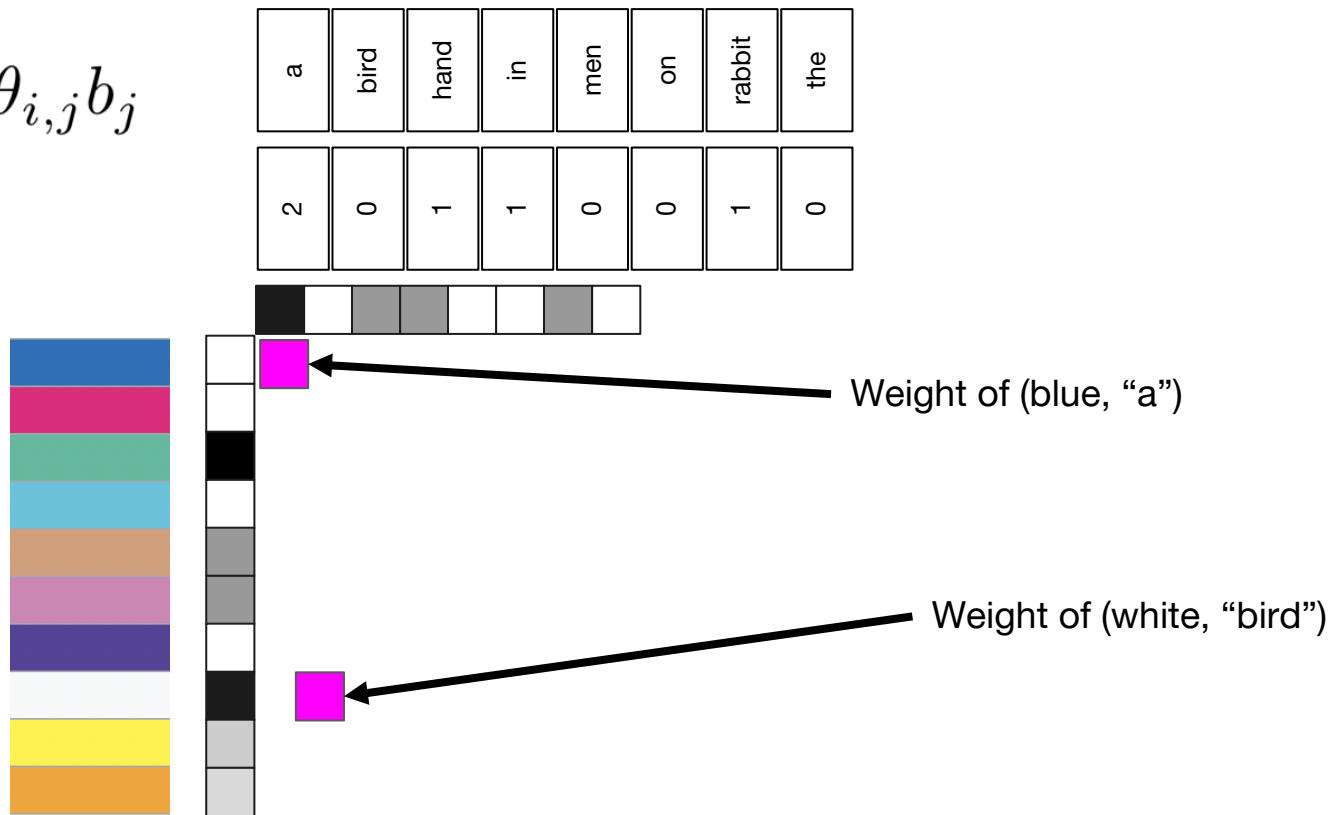
You shall know the **features of modality A** by the company they keep in the **features of modality B**. **And vice versa.**

## Formulating the Retrieval Problem as a Linear Model

$$F_{(\psi, \omega)} : \psi(\mathcal{I}) \times \omega(\mathcal{L}) \rightarrow \mathbb{R}$$

# Formulating the Retrieval Problem as a Linear Model

$$\sum_{i=1, j=1}^{n, m} a_i \theta_{i,j} b_j$$



## Pointwise Mutual Information

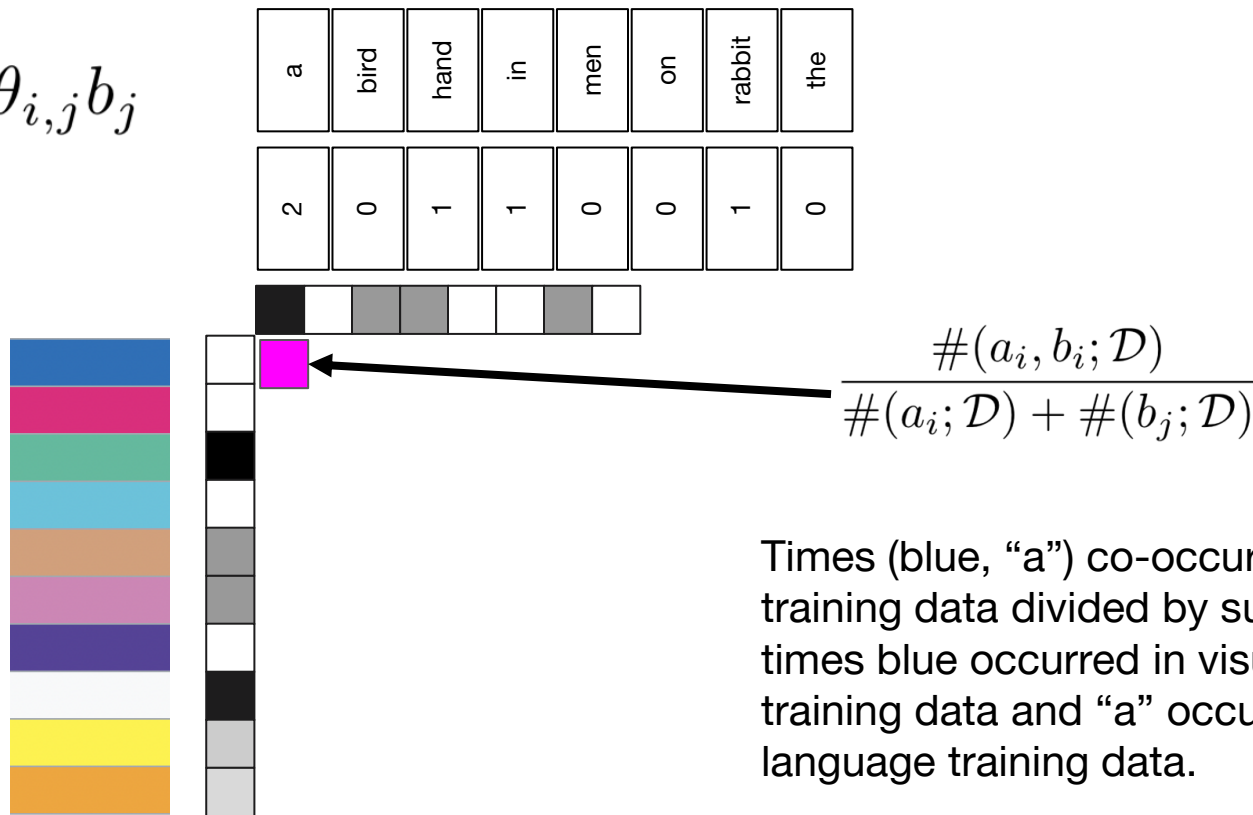
$$\text{pmi}(x; y) = \log \frac{p(x, y)}{p(x)p(y)}$$

$$F(A, B) = \text{pmi}(A; B); A \in \mathcal{I}, B \in \mathcal{L}$$

$$\begin{aligned}
F(\phi(A), \omega(B)) &= \text{pmi}(\vec{a}; \vec{b}) \\
&= \log \frac{p(\vec{a}, \vec{b})}{p(\vec{a})p(\vec{b})}
\end{aligned}$$

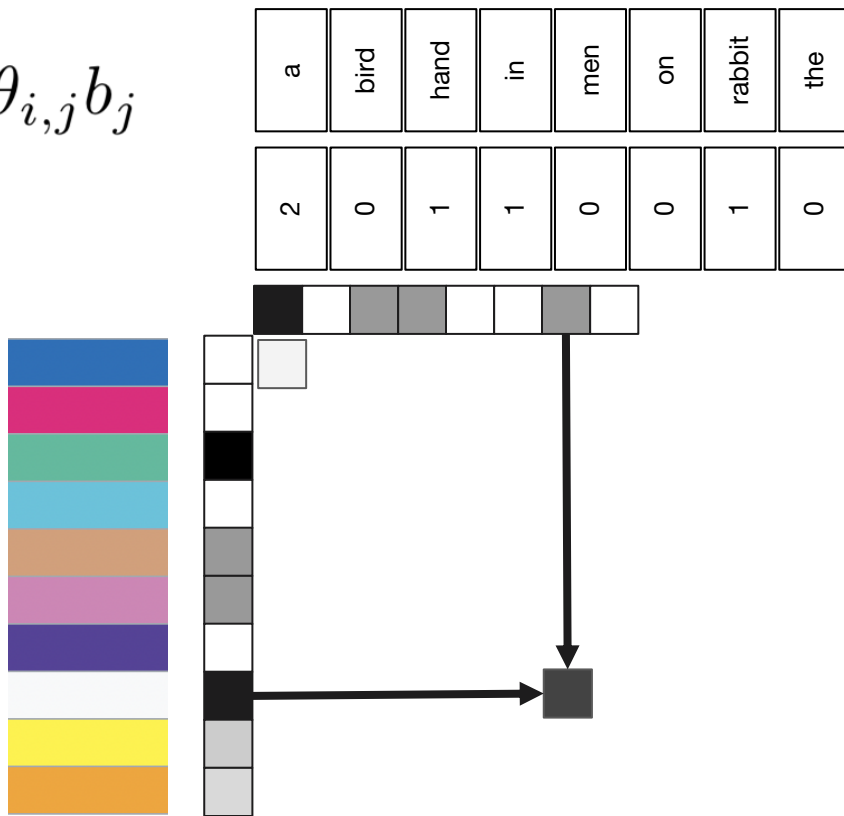
# Formulating the Retrieval Problem as a Linear Model

$$\sum_{i=1, j=1}^{n, m} a_i \theta_{i,j} b_j$$



# Formulating the Retrieval Problem as a Linear Model

$$\sum_{i=1, j=1}^{n, m} a_i \theta_{i,j} b_j$$



# Feature Extraction and Representation Learning

$$X \subset \mathcal{I}, Y \subset \mathcal{L}$$

$$a \in \psi(\mathcal{I}), b \in \omega(\mathcal{L})$$

$$\mathcal{H}(P(Y|a))$$

$$\mathcal{H}(P(X|b))$$



$b$  indicates "rabbit"

"A rabbit sits in the palm of a hand."

"A tropical bird perches in the jungle."

"Men are talking on a basketball court."

"Children play soccer in a field."

"Sports equipment is staged for a photo."

"A white fox is looking at the camera."

$a$  indicates detected furry ear



# Tokenization Considering Language and Vision Retrieval



“A tropical bird  
perches in the jungle.”



“Perching parakeet in  
a wire frame cage.”



“A snow owl lands on  
a wooden perch.”

$$\mathcal{H}(P(X|b))$$

# Tokenization Considering Language and Vision Retrieval



“A tropical bird  
perches in the jungle.”



“Perching parakeet in  
a wire frame cage.”



“A snow owl lands on  
a wooden perch.”

Stemming

a

tropic

bird

perch

in

the

jungl

.

perch

parakeet

in

a

wire

frame

cage

.

a

snow

owl

land

on

wood

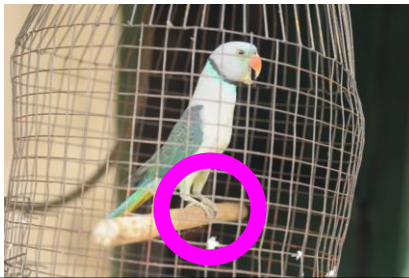
perch

.

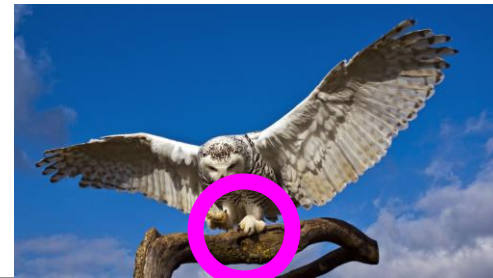
# Tokenization Considering Language and Vision Retrieval



“A tropical bird  
perches in the jungle.”



“Perching parakeet in  
a wire frame cage.”



“A snow owl lands on  
a wooden perch.”

Stemming

a

tropic

bird

perch

in

the

jungl

.

perch

parakeet

in

a

wire

frame

cage

.

a

snow

owl

land

on

wood

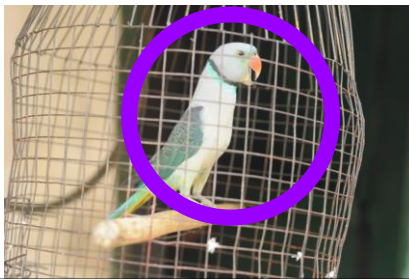
perch

.

# Tokenization Considering Language and Vision Retrieval



“A tropical bird  
perches in the jungle.”



“Perching parakeet in  
a wire frame cage.”



“A snow owl lands on  
a wooden perch.”

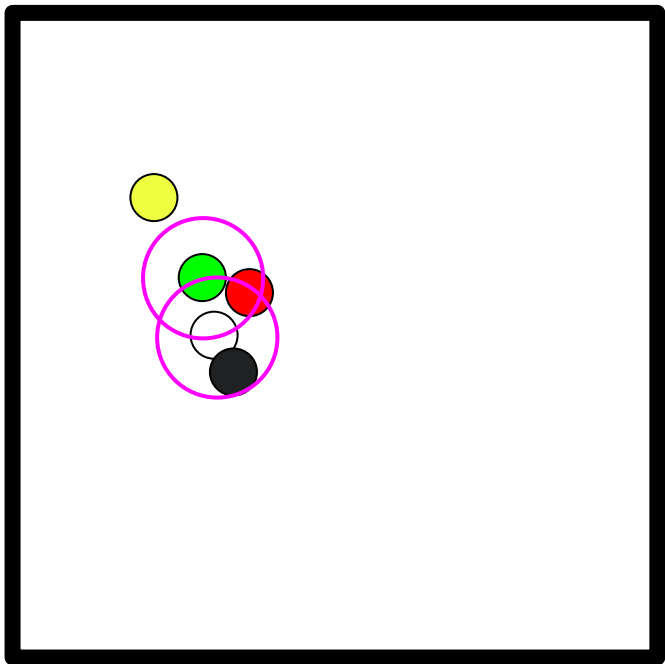
a	tropic	bird	perch	in	the	jungl	.
perch	parakeet	in	a	wire	frame	cage	.
a	snow	owl	land	on	wood	perch	.

# Pretrained Language Token Embeddings

- Cosine similarity of “bird”, “owl”, “parakeet” helps share information across training data
- Taking a guess: what are the nearest neighbors of “green” in word2vec embedding space?
  - Blue, white, red, yellow, black, grey, purple, pink, light, gray
- What can we learn for highly polysemous words like “play”?
  - “play guitar”, “play piano”, “play basketball”, “play tag”
- Text-based embeddings help us share information ***only to the extent that words used in a similar context share a visual representation***
  - which is true for, say, birds or trees, but not colors

# Pretrained Language Token Embeddings

$$\omega(\mathcal{L})$$



Training

“A white bunny rabbit  
held in a green field.”



Inference

“A black rabbit  
watches a red sunset.”



# Considering Language and Vision Retrieval



“A tropical bird  
perches in the jungle.”



“Perching parakeet in  
a wire frame cage.”



“A snow owl lands on  
a wooden perch.”

$$\mathcal{H}(P(Y|a))$$

# Image Feats. Considering Language and Vision Retrieval



“A tropical bird  
perches in the jungle.”



“Perching parakeet in  
a wire frame cage.”



“A snow owl lands on  
a wooden perch.”

Object  
Classification

toucan

cockatoo

great grey  
owl



# Image Feats. Considering Language and Vision Retrieval



“A tropical bird perches in the jungle.”



“Perching parakeet in a wire frame cage.”



“A snow owl lands on a wooden perch.”

Object  
Detection

bird

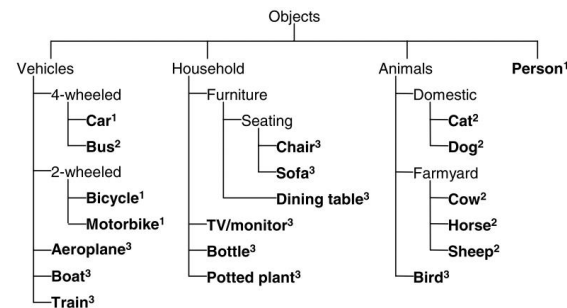
potted plant

bird

boat

bird

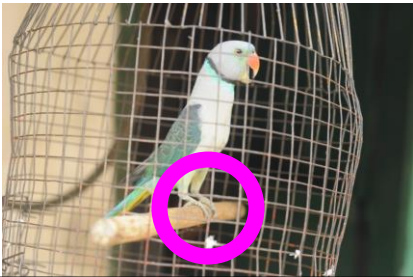
aeroplane



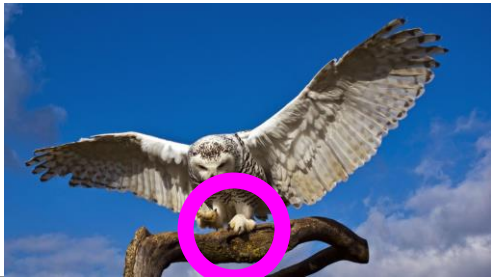
# Image Feats. Considering Language and Vision Retrieval



“A tropical bird  
perches in the jungle.”



“Perching parakeet in  
a wire frame cage.”



“A snow owl lands on  
a wooden perch.”

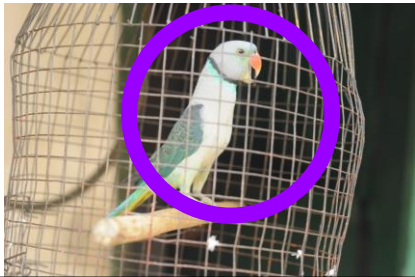
bird	potted plant
bird	boat
bird	aeroplane

a	tropic	bird	perch	in	the	jungl	.
perch	parakeet	in	a	wire	frame	cage	.
a	snow	owl	land	on	wood	perch	.

# Image Feats. Considering Language and Vision Retrieval



“A tropical bird  
perches in the jungle.”



“Perching parakeet in  
a wire frame cage.”



“A snow owl lands on  
a wooden perch.”

bird	potted plant
bird	boat
bird	aeroplane

a	tropic	bird	perch	in	the	jungl	.
perch	parakeet	in	a	wire	frame	cage	.
a	snow	owl	land	on	wood	perch	.

		"A rabbit sits in the palm of a hand."		"Children play soccer in a field."	
		"A white fox is looking at the camera."		"Men are talking on a basketball court."	

Image  
Embedder

$$\psi(\mathcal{I})$$

Caption  
Embedder

$$\omega(\mathcal{L})$$

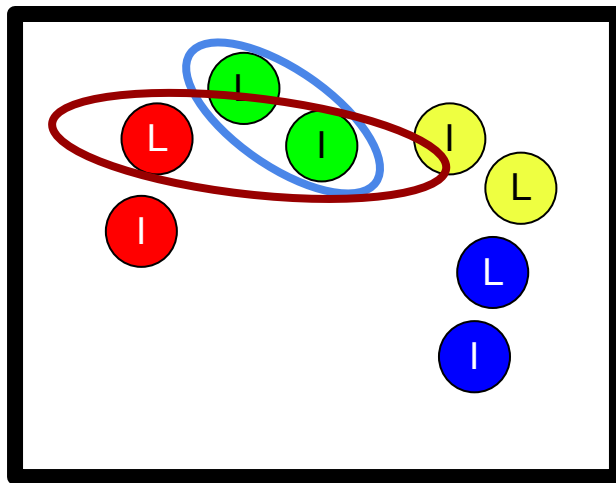
Pull matching

image

embeddings

together.

Learned  
Projection

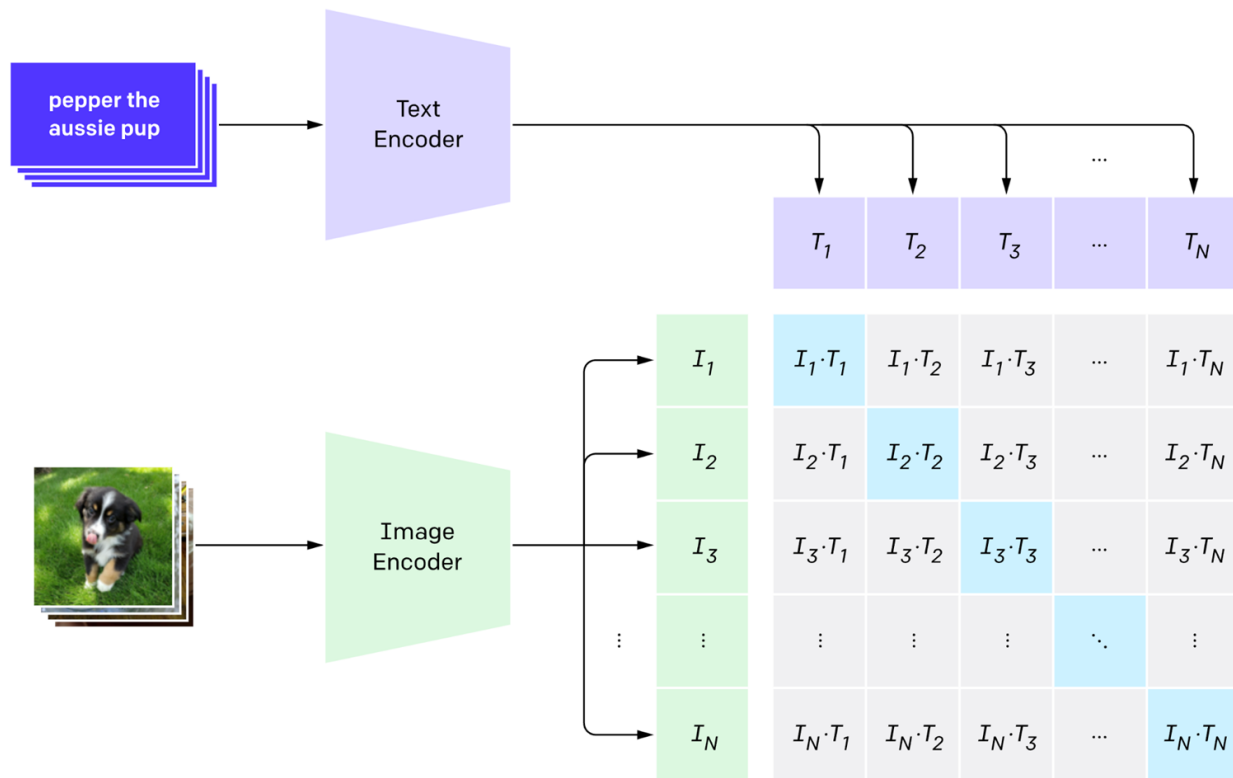


Push distractor

Learned  
Projection

embeddings  
apart.

# Contrastive Language–Image Pre-training (CLIP)



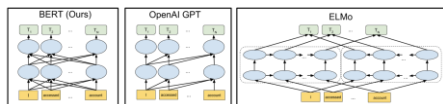
# Experience Grounds Language - World Scopes

Common Crawl



**WS1**

Carefully  
Annotated  
Penn Treebank  
Brown Corpus  
WordNet



**WS2**

Unstructured  
Unlimited Text  
Common Crawl  
Word2Vec  
ELMo, BERT, GPT\*

**WS3**

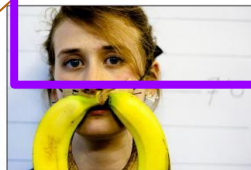
Text Paired with  
Sensory Data  
ImageNet, VQA,  
ViLBERT,  
Video Captioning



By Jiri Stehlik  
the trail climbs steadily  
uphill most of the way.



By Daniel Heisey  
the stars in the night sky.



What color are her eyes?  
What is the mustache made of?



**WS4**

Language and an Embodied  
Agent  
VLN\*, IQA,  
NL+Games,  
RoboNLP

Doc:  
The Rebel Enclave consists of jacket  
spider, and wasp. Acolyte, blessed items  
are useful for poison monsters. Star  
Alliance contains bat, panther, and wolf.  
Goblin, jaguar, and lynx are on the same  
team - they are in the Order of the Forest.  
Gleaming and mysterious weapons beat  
cold monsters. Lightning monsters are  
weak against Grandmaster's and  
Soldier's weapons. Fire monsters are  
defeated by fanatical and shimmering  
weapons.  
Goal:  
Defeat the Order of the Forest



**WS5**

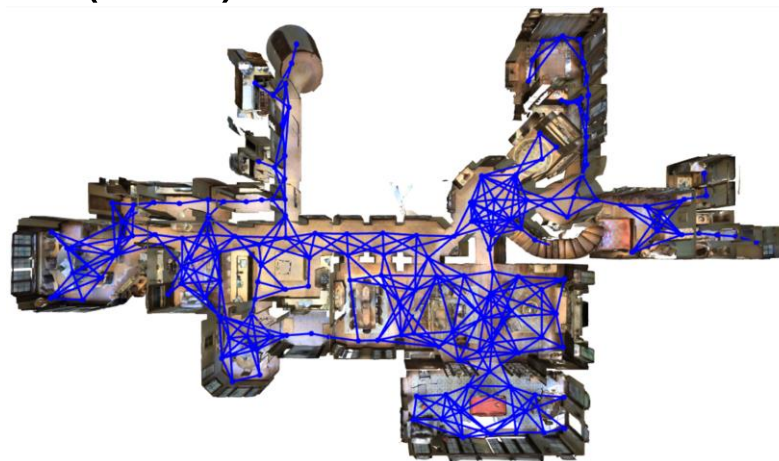
Social  
Embodiment  
Human-robot  
collaboration  
and dialog



# Vision-and-Language Navigation (VLN)

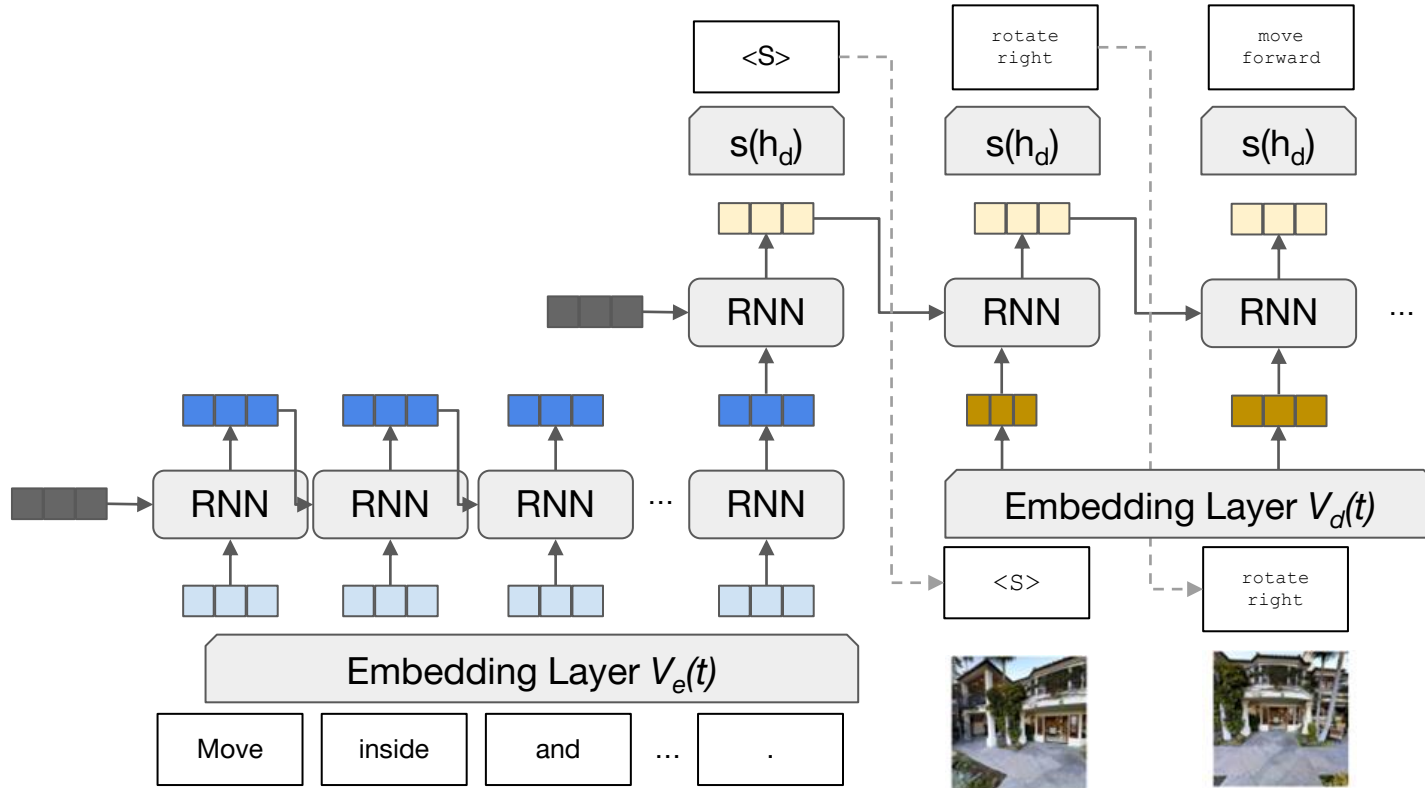


Leave the bedroom, and enter the kitchen. Walk forward, and take a left at the couch. Stop in front of the window.



- Source language input: instructions in English
- Target language output: sequence of actions like “forward”, “rotate left”, “rotate right”, and “stop”

# Machine Translation for Robots

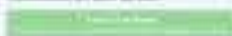




Get the map on the left from the page with the map in the center.



For to find the goal room, green the first (left) LUNA and drag to the center room, right click (blue cylinder to move). And your partner for keep being the stick. Your partner can only your movement on their screen.



Also, The goal room contains a ball.

Find the goal room

Find	Find
Find the ball in the room?	Find the ball in the room?
Find the ball in the room?	Find the ball in the room?

Find the ball in the room?

Find the ball in the room?

The ball is in the room. The ball is in the room. The ball is in the room.

Get the map on the left from the page with the map in the center.



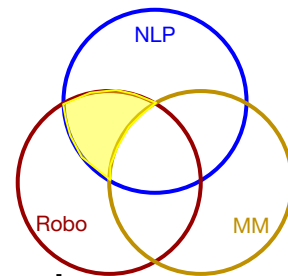
Your partner has moved for help in the first window. View the ball in the first window by clicking "Show Ball" button. And after you answer back to them.



Find the ball in the room?

## Activity: Cooperative Vision-and-Dialogue Navigation

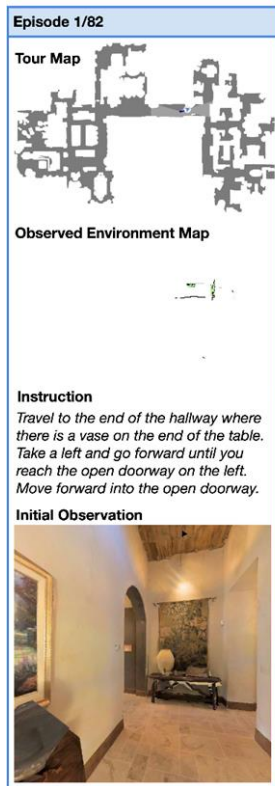
- If we ask a computer to do it, we have to be able to do it too!
- <https://cvdn.dev/>
  - “Two-player Demo”
  - Scroll to the bottom, click “Start Task”



# Language-guided Robots with Vision

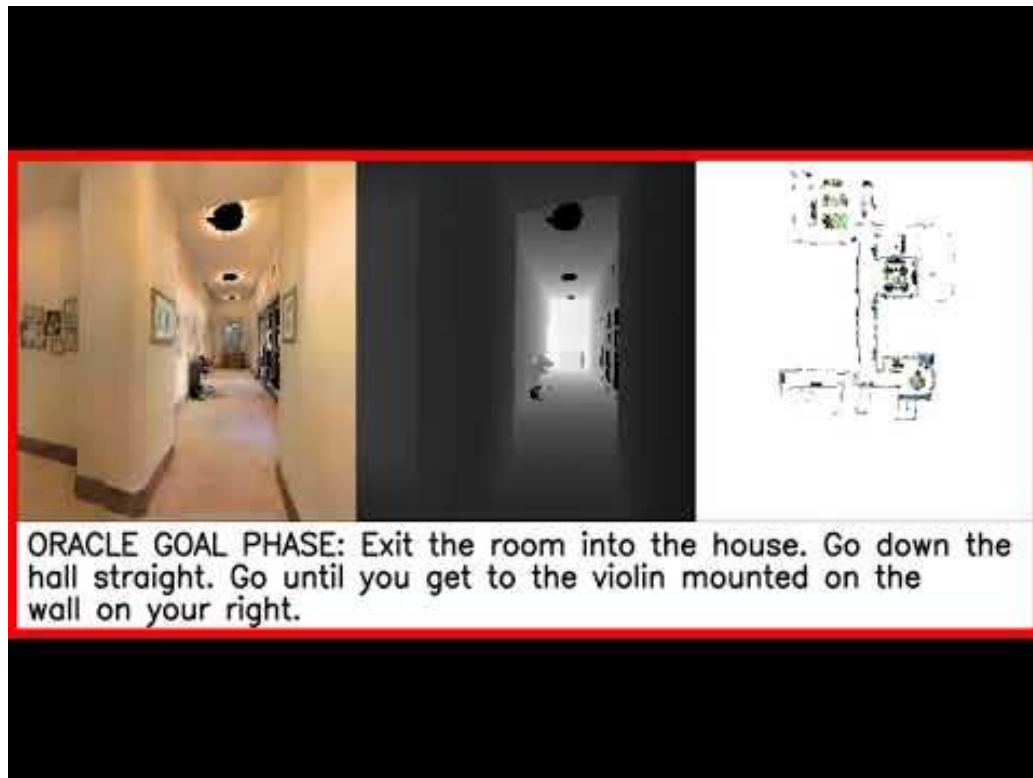
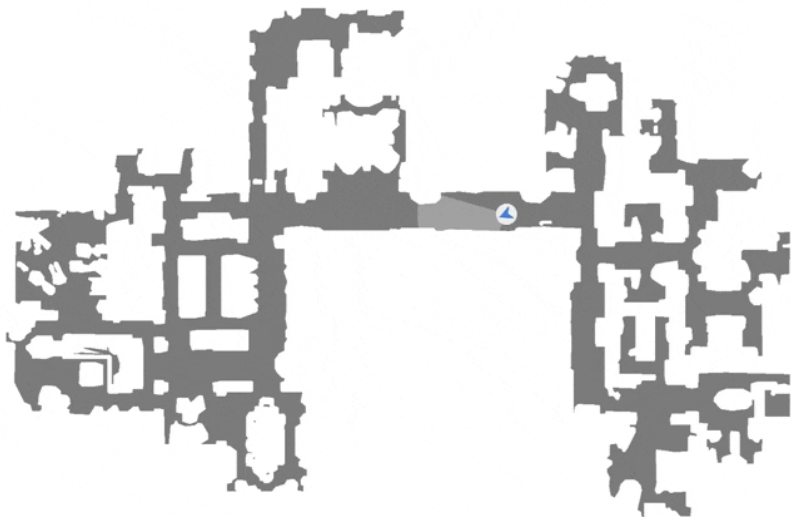
- We are exploring three directions in the space of translating language instructions to agent movement plans:
  - *Memory*: following instructions and building maps together.
  - *Physical Control*: overcoming the gap between simulation and real robot control.
  - *Planning*: translating underspecified language into coherent plans to get a job done.

# Memory: Iterative Vision-and-Language Navigation [in submission]

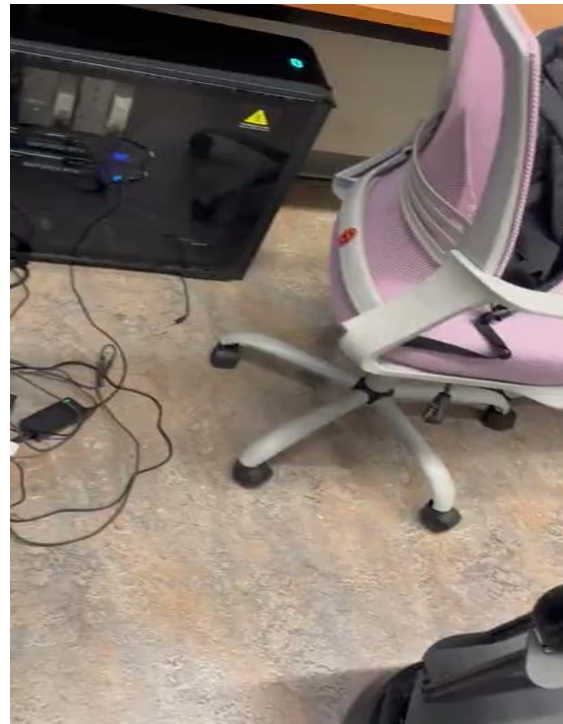
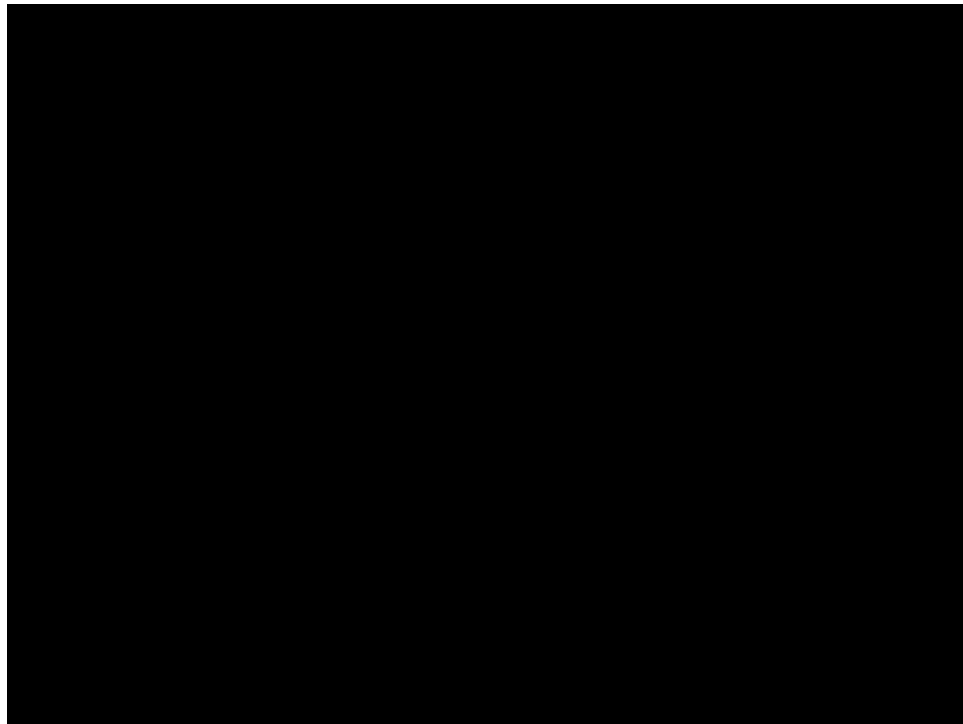


# Memory: Iterative Vision-and-Language Navigation

Phase: oracle\_start    Observed: 2%    Episode: 0/82



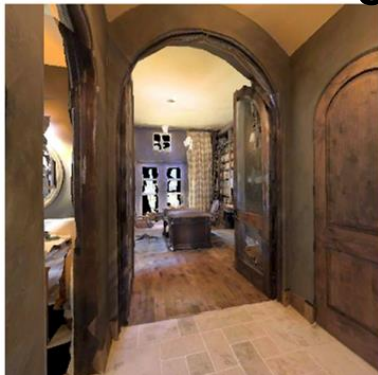
## *Physical Control: VLN With A Physical Robot*



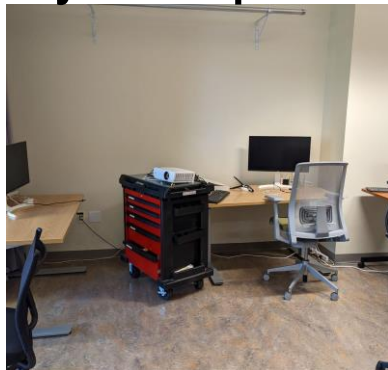
*“Go to the couch and stop.”*

# Physical Control: VLN With A Physical Robot

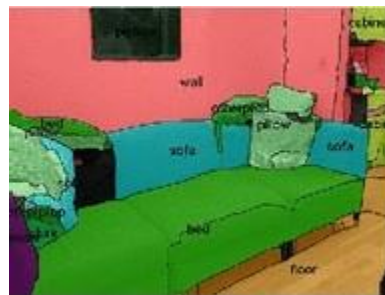
## Simulation Training



My Lab Space



- Use *mid-level* vision to align virtual and real world spaces

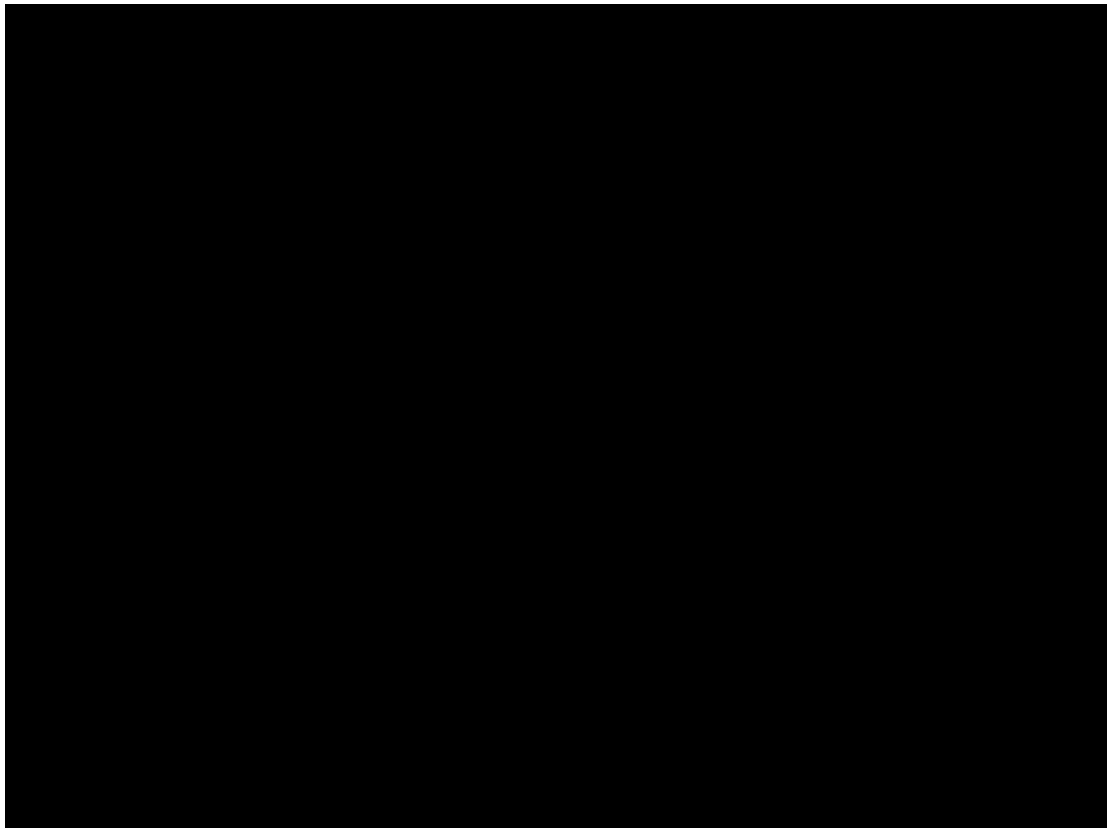


## *Planning: Prompting for Robot Control [in submission]*

- LLMs generate actions robots can't do with objects that aren't around
- Pythonic *prompts* can specify robot actions and world objects



## *Planning:* Prompting for Robot Control

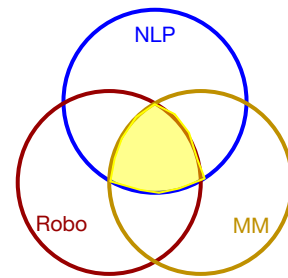


- Used to control an agent in a virtual environment
- Perform new tasks in a zero-shot setting

# *Planning:* Prompting for Robot Control



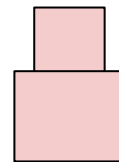
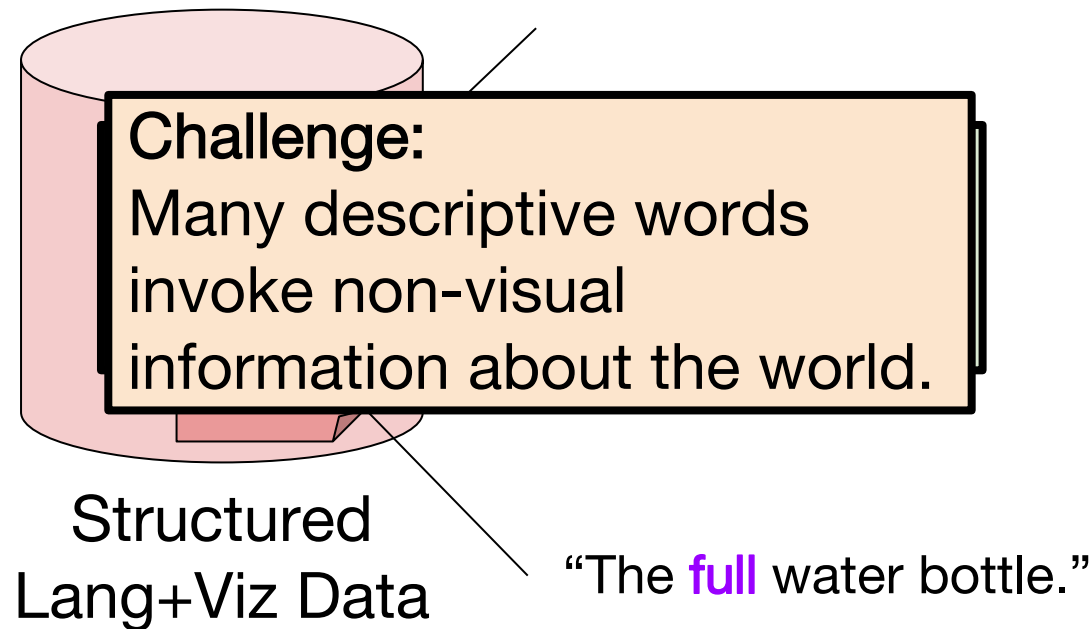
- Also enables generating pick-and-place robot plans!



# Language-guided Embodied, Multimodal Agents

- Sometimes robot *vision* is not enough.
- A robot is more than a mounted camera; it exists in the physical world and enables *physical sensory experience*
- We can combine language instructions and dialogue for learning with machine vision and perception for robots that learn from people and all their senses

# Language-guided Embodied, Multimodal Agents



Physical Robot  
Agent

# Sensory Perception Beyond Vision

## Look



***color***, ***shape***, and deep ***visual*** features extracted from images of the object.

## Push



## Press



## Grasp



## Lift / Lower



## Drop



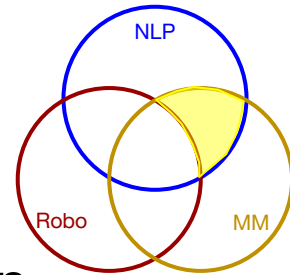
***haptic*** and ***audio*** features captured from arm and microphone.

# Grounding Multimodal Sensory Perception with Dialogue

## **Jointly Improving Parsing and Perception for Natural Language Commands through Human-Robot Dialog**

Jesse Thomason, Aishwarya Padmakumar, Jivko Sinapov,  
Nick Walker, Yuqian Jiang, Harel Yedidsion, Justin Hart,  
Peter Stone, Ray Mooney

# Multimodal NLP: Other Times Text is not Enough



- We are exploring the integration of language-and-vision representations with other sensory and linguistic signals for:
  - *Continual Learning*: creating systems that learn from multiple tasks that involve combinations of language and vision input [*to appear @ Neurips'22!*]
  - *Sign Recognition*: using linguistic priors to better identify ASL signs from videos [*in submission*]
  - *Dementia Detection*: using both language transcriptions and speech signals to detect disfluencies indicating early stage dementia.

# Natural Language Processing and Robotics

Jesse Thomason

[ <http://glamor.rocks/> ]

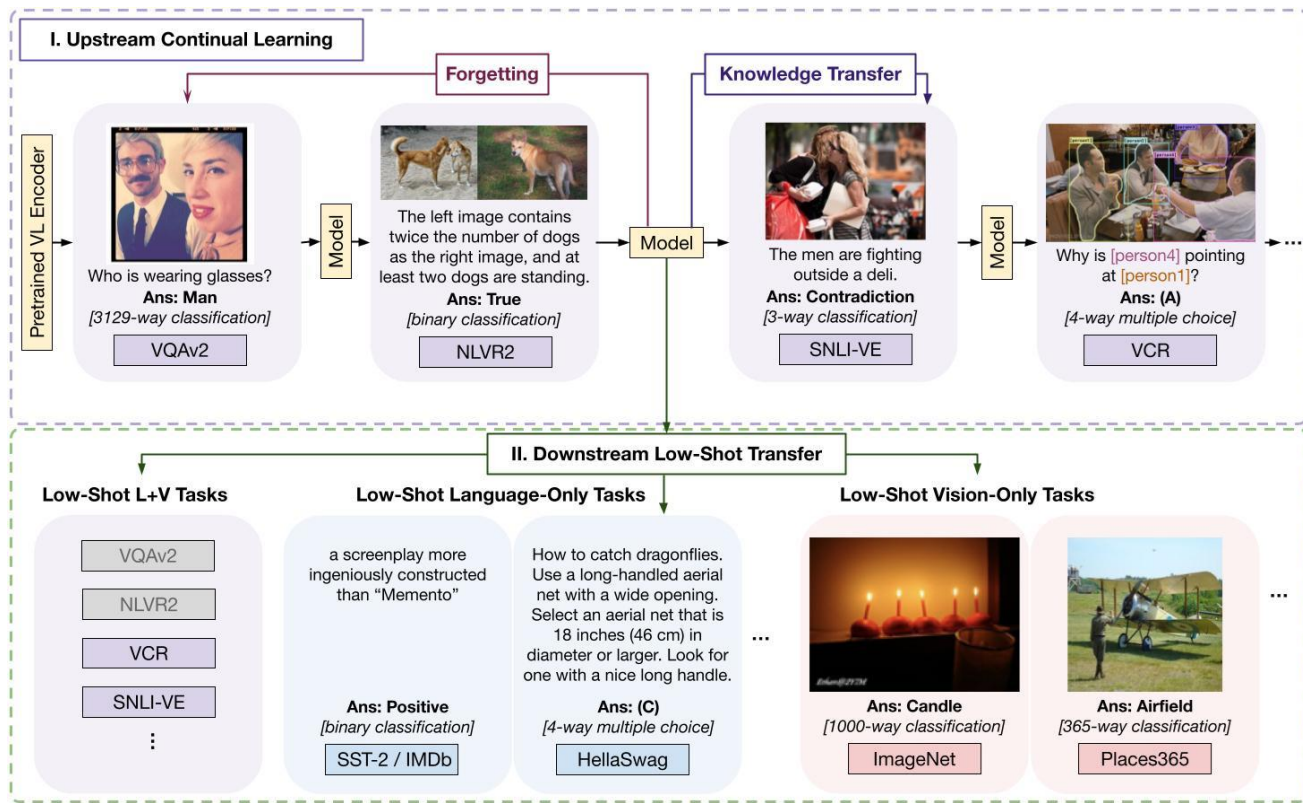
**USC** Viterbi  
School of Engineering

✧✧ GLAMOR

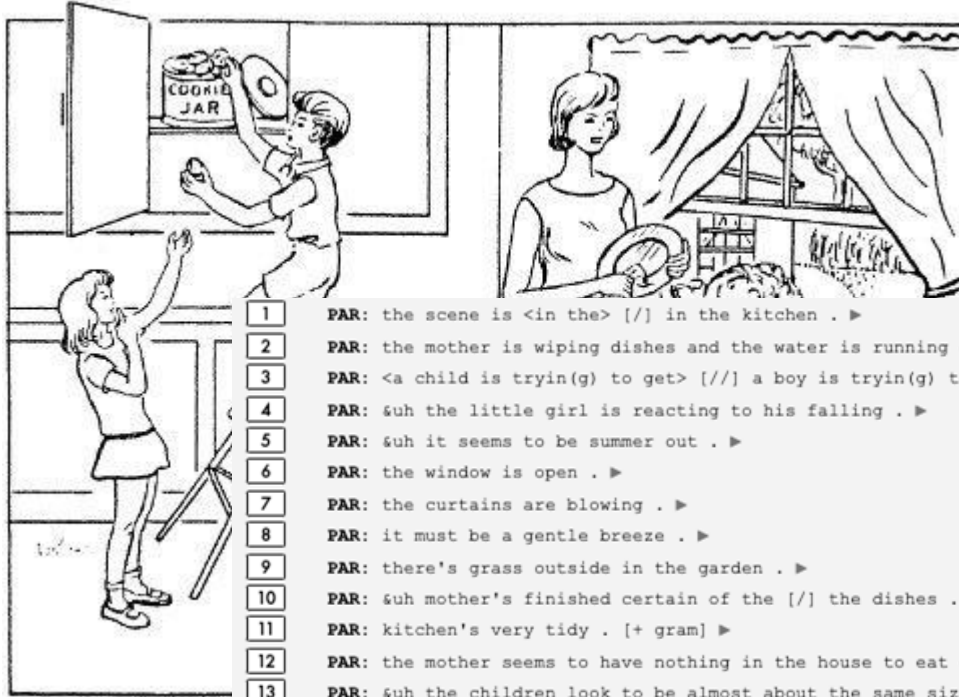


# Multimodal Continual Learning

- Explore *adapter fusion* and *split adapters* to learn future tasks, even when a modality disappears



# Early Detection of Dementia



- Predict based on both text *and* audio signals

1 PAR: the scene is <in the> [/] in the kitchen . ►  
2 PAR: the mother is wiping dishes and the water is running on the floor . ►  
3 PAR: <a child is tryin(g) to get> [/] a boy is tryin(g) to get cookies outta [: out of] a jar and he's about to tip over on a stool . ►  
4 PAR: suh the little girl is reacting to his falling . ►  
5 PAR: suh it seems to be summer out . ►  
6 PAR: the window is open . ►  
7 PAR: the curtains are blowing . ►  
8 PAR: it must be a gentle breeze . ►  
9 PAR: there's grass outside in the garden . ►  
10 PAR: suh mother's finished certain of the [/] the dishes . ►  
11 PAR: kitchen's very tidy . [+ gram] ►  
12 PAR: the mother seems to have nothing in the house to eat except cookies in the cookie jar . ►  
13 PAR: suh the children look to be almost about the same size . ►  
14 PAR: perhaps they're twins . ►  
15 PAR: they're dressed for summer warm weather . ►

# Isolated Sign Language Recognition



- Predict signed word from video frames
- Adding auxiliary predictions can help!
- Handshape
- Location

