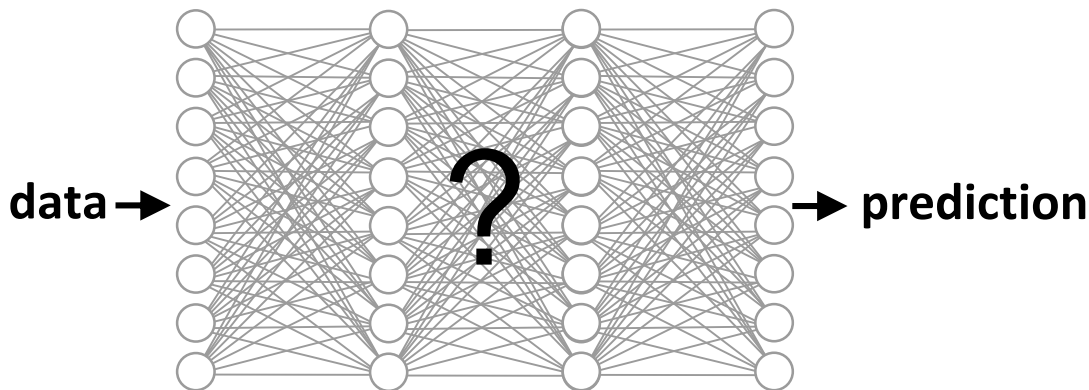


# Deeply Explainable AI

Darrell, Abbeel, Dragan, Klein, Griffiths, Canny, Saenko, Akata, Hoogs



# Current deep neural networks (DNNs) are “**black boxes**”



- Do not expose their **decision making** process
- Do not provide their **confidence** in their predictions
- Not clear whether they can be **trusted** and/or **corrected**

How do we make  
**Deep Learning**  
more **explainable** and  
**trustworthy?**

# Overview

Challenge Problems (Kate)

Explicit and Implicit Explanation Models (Kate)

Learning How to Explain (Zeynep)

Modeling the User (Trevor)

Program Schedule (Trevor)

# Challenge Problems

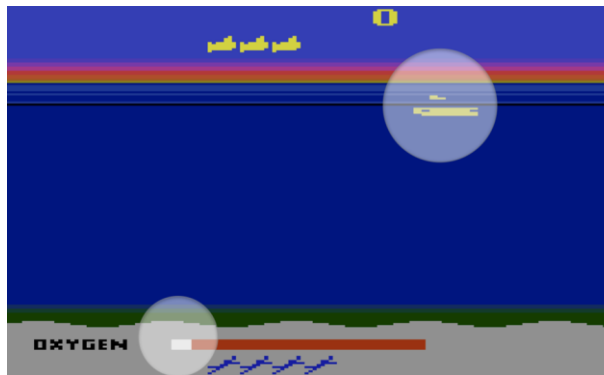
# Autonomy Challenge: control of autonomous vehicles

Can we explain the agent's behavior?

Demonstration-based control

Reinforcement-learned control

In simulated environments



**Textual Explanation:**  
Submarine is going up to surface to replenish oxygen

**Attention Explanation:**  
Agent decides to go up after looking at oxygen bar and current position



**Control Explanation:**  
A left banked turn is engaged to avoid crashing into canyon



**Route Explanation:**  
Avoid buildings, Avoid being seen by people

# Analytics Challenge:

## Multimedia Event Question Answering (MEQA)

Can we explain the system's answers?

Answering questions about images /  
video with associated audio / text

Interactive UI enables natural  
language dialogue with system

Data collection

Extend VQA, MovieQA datasets



**Q:** Can these people  
arrest someone?

**A:** Yes

**Explanation:** ...  
because they are  
Vancouver police



**Q:** What is he doing?

**A:** Juggling

**Explanation:**  
...because he has two  
balls in his hands  
while two are in the air

# Analytics Challenge:

## Multimedia Event Question Answering (MEQA)

Can we explain the system's answers?

Answering questions about images /  
video with associated audio / text

Interactive UI enables natural  
language dialogue with system

Data collection

Extend VQA, MovieQA datasets

Extend Berkeley DeepDrive dataset  
(100K hours of driving video,  
GPS/sensor IMU data) with text



**User:** “why did you turn left”?

*To avoid the traffic delay near  
the shopping mall.*

**User:** “why did we stop?”

*I can't tell, is this a shadow or a pothole?*

**User:** “how do you know how to drive through  
this intersection?”

*I learned from 236 prior driving exemplars  
transiting in this direction at this time of day in  
the past month.*





# Explicit Explanation Models

# Explicit Explanation Models

Explain higher-level reasoning in DNNs

Explainable decision path for multi-task, control and planning

Provide structure and intermediate state

**Q:** Can you park here?



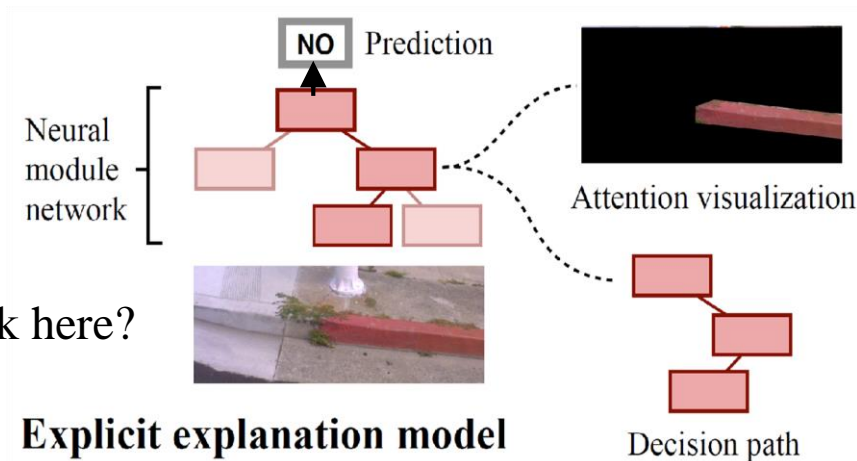
# Explicit Explanation Models

Explain higher-level reasoning in DNNs

Explainable decision path for multi-task, control and planning

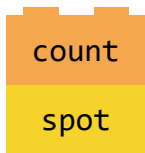
Provide structure and intermediate state

**Q:** Can you park here?

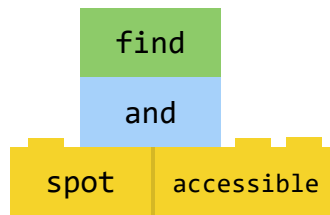


# Neural module networks

*How many parking spots are there?*

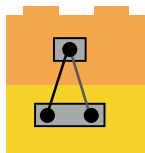


*Where is the accessible parking spot?*

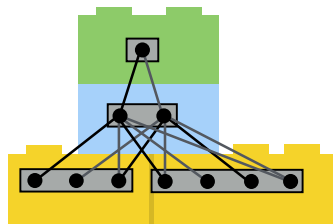


# Neural module networks

*How many parking spots are there?*



*Where is the accessible parking spot?*



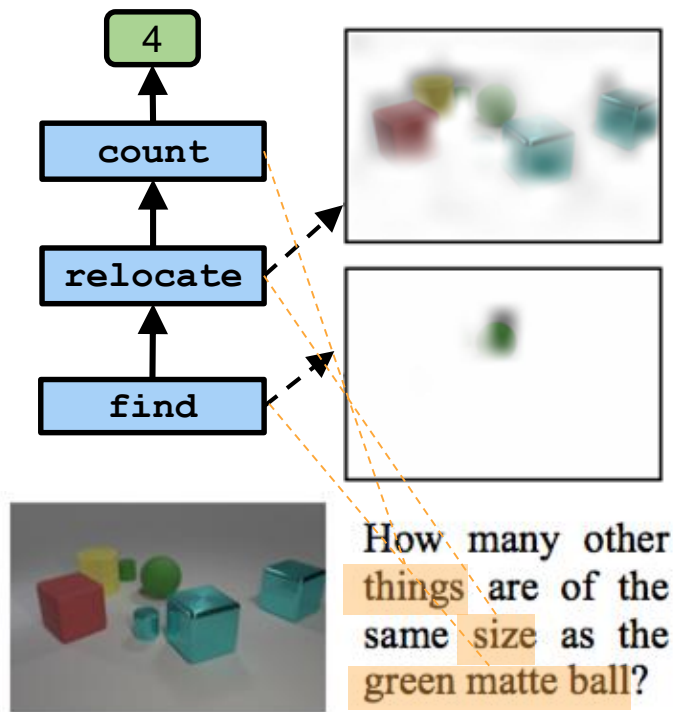
# Neural module networks

Explain answers via

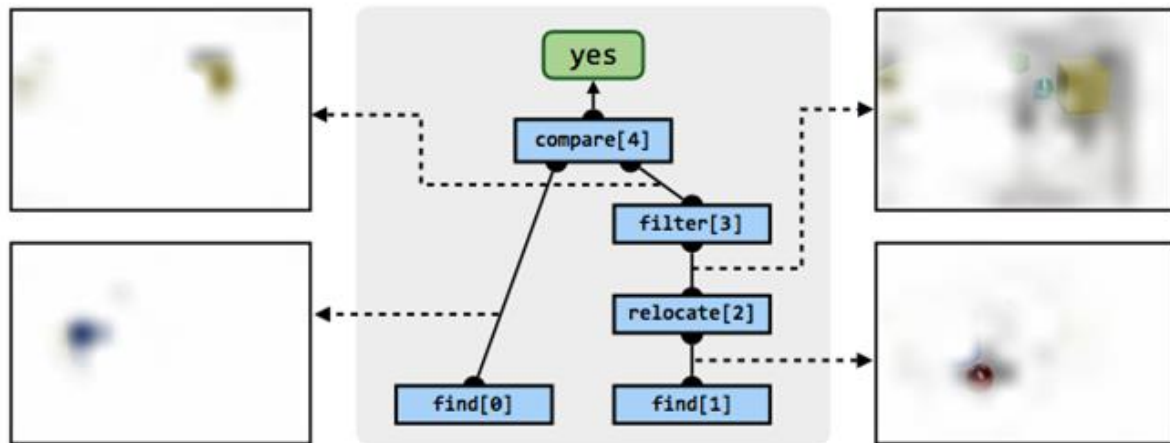
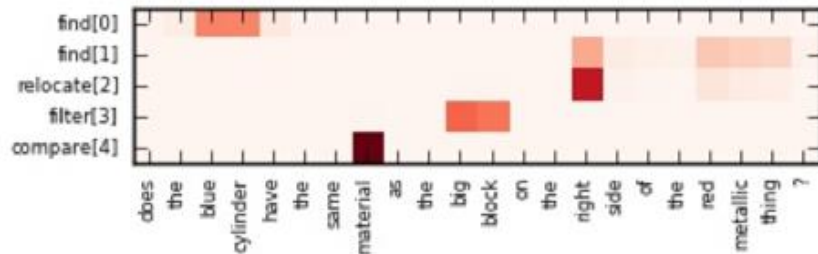
exposing modules used for prediction

showing the chain of reasoning

visualizing text and image attention



# Neural module networks



Does the blue cylinder have the same material as the big block on the right side of the red metallic thing?

# Learning Modular Neural Network Policies for Multi-Task and Multi-Robot Transfer

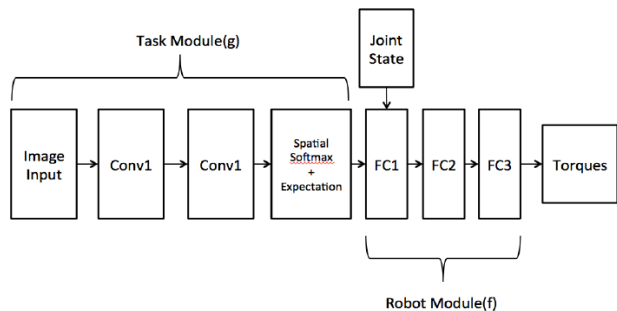
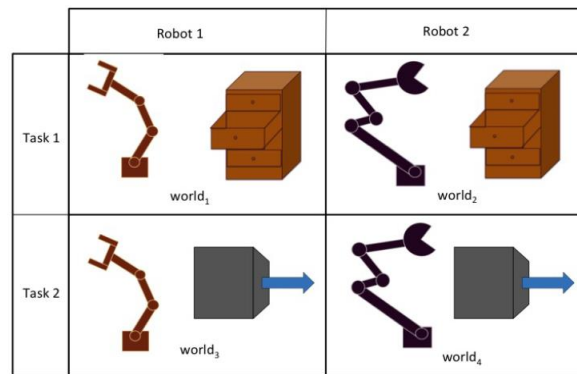
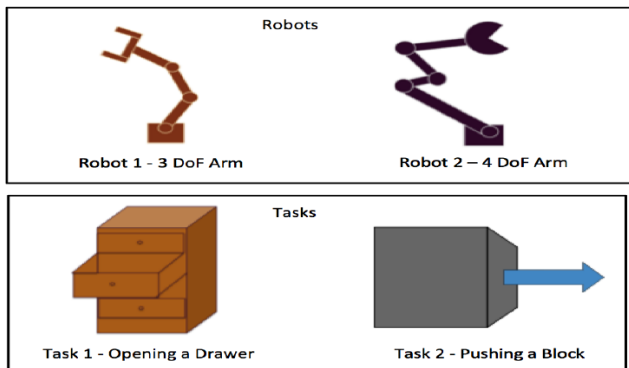
Coline Devin<sup>\*1</sup>

Abhishek Gupta<sup>\*1</sup>

Trevor Darrell<sup>1</sup>

Pieter Abbeel<sup>1</sup>

Sergey Levine<sup>1</sup>



## Available Modules

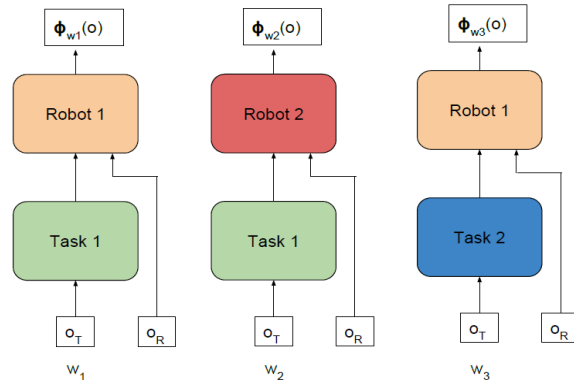
Robot Modules



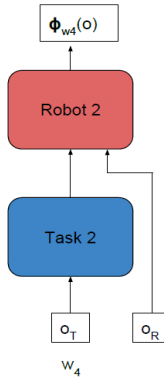
Task Modules



## Training Worlds



## Unseen World





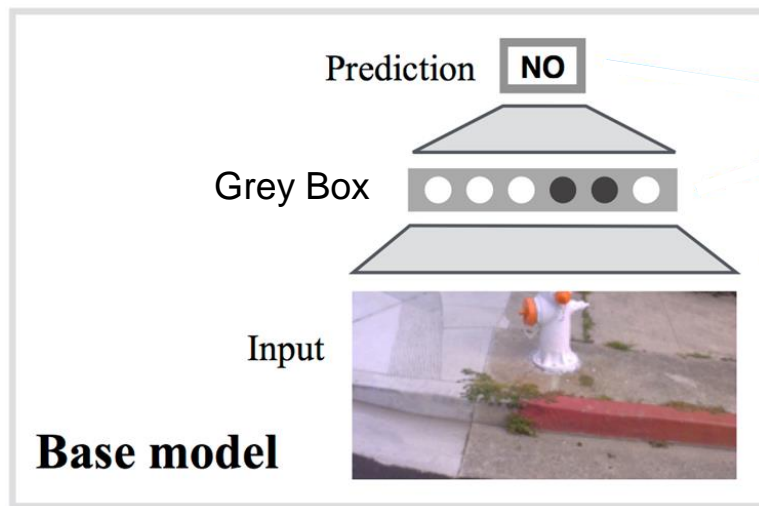
# Implicit Explanation Models

# Implicit Explanation Models

Recover visualizations or exemplars from black/grey box DNN

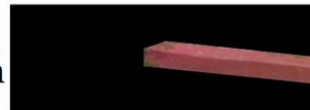
No explicit internal meaning representation is needed

**Q:** Can you park here?

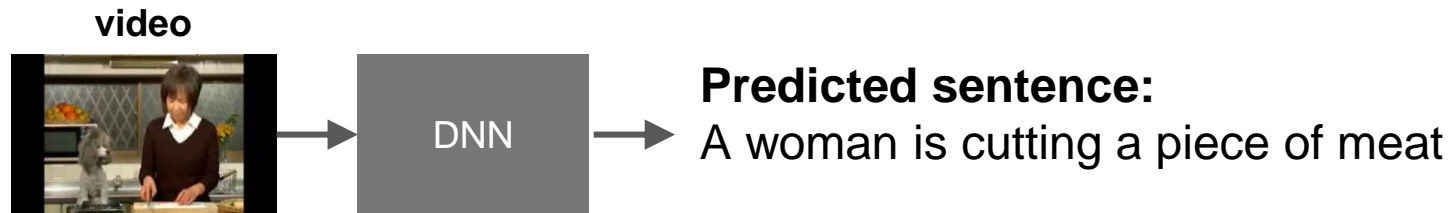


Model introspection

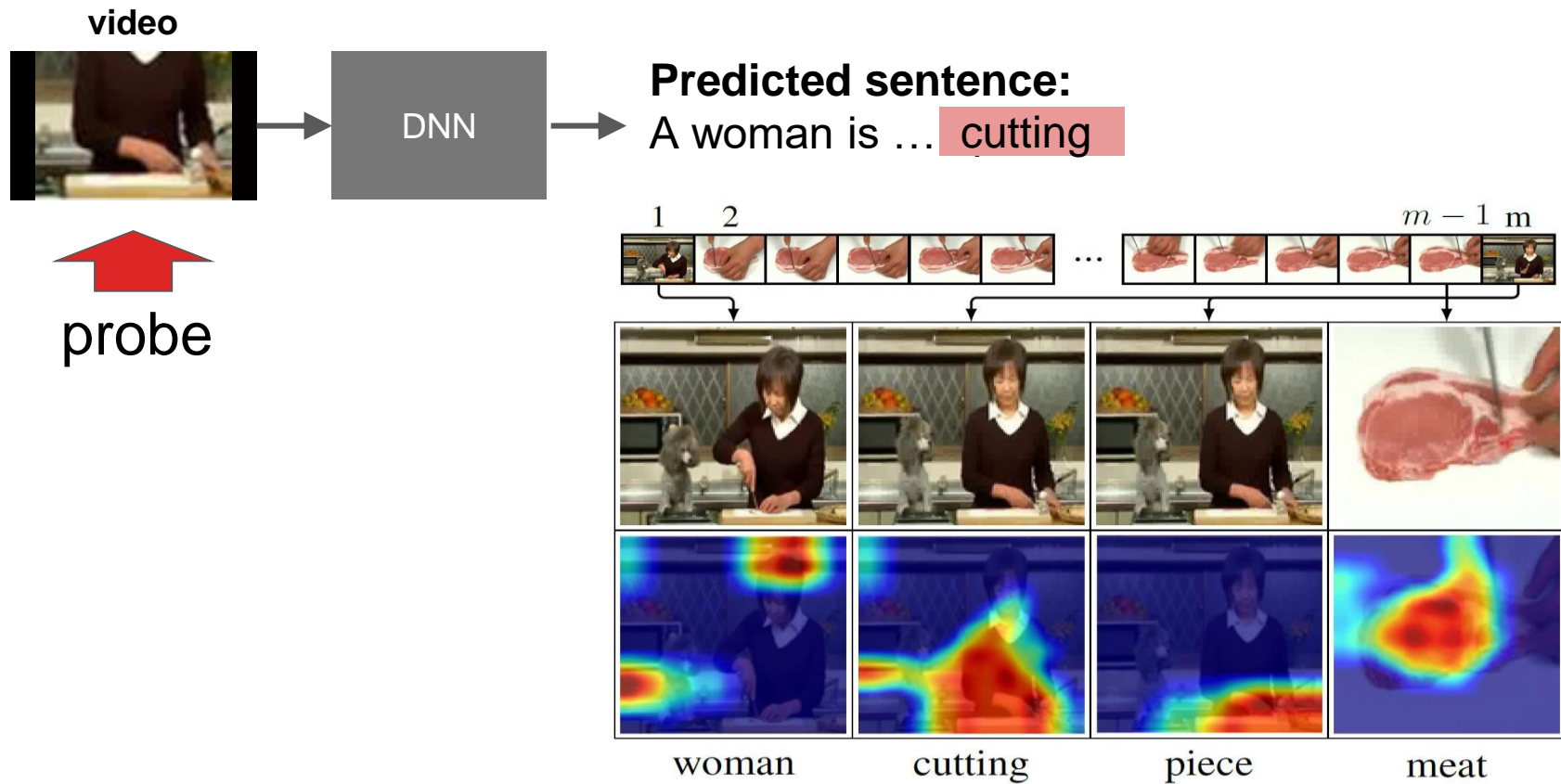
**Implicit explanation model**



# Implicit explanation via Top-down Saliency

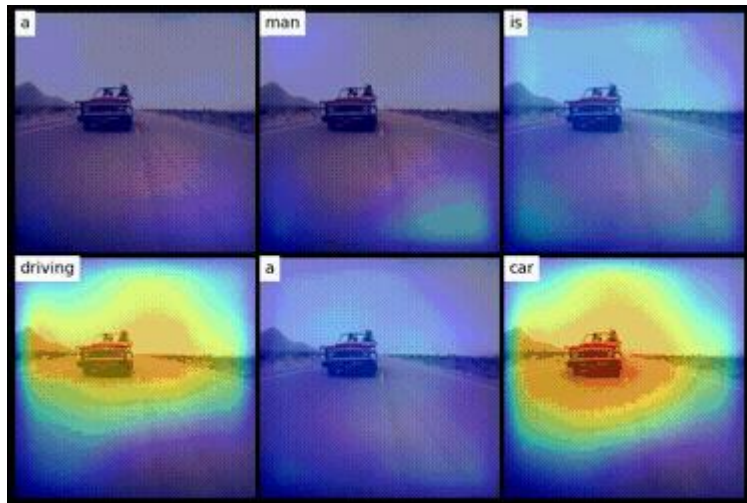


# Implicit explanation via Top-down Saliency

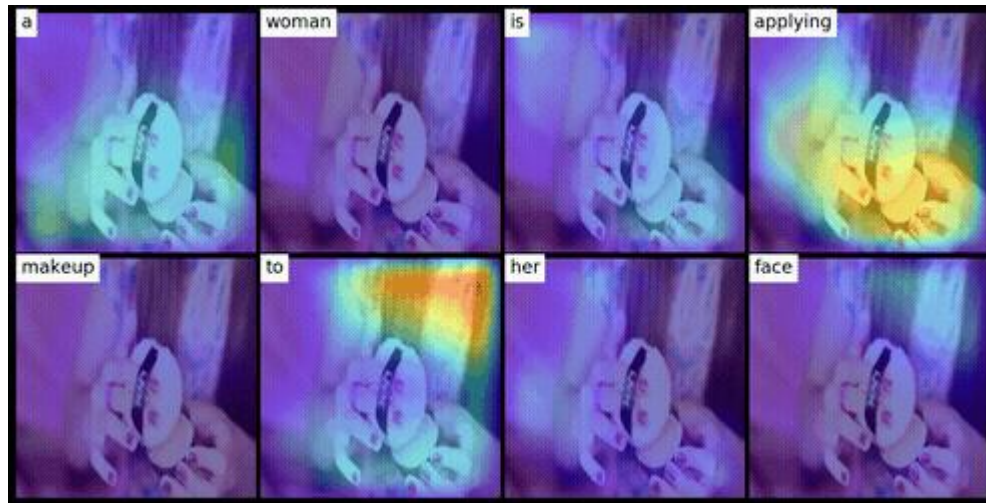


# Implicit Explanations via Top-down Saliency

**Prediction:** A man is driving a car



**Prediction:** A woman is applying makeup to her face

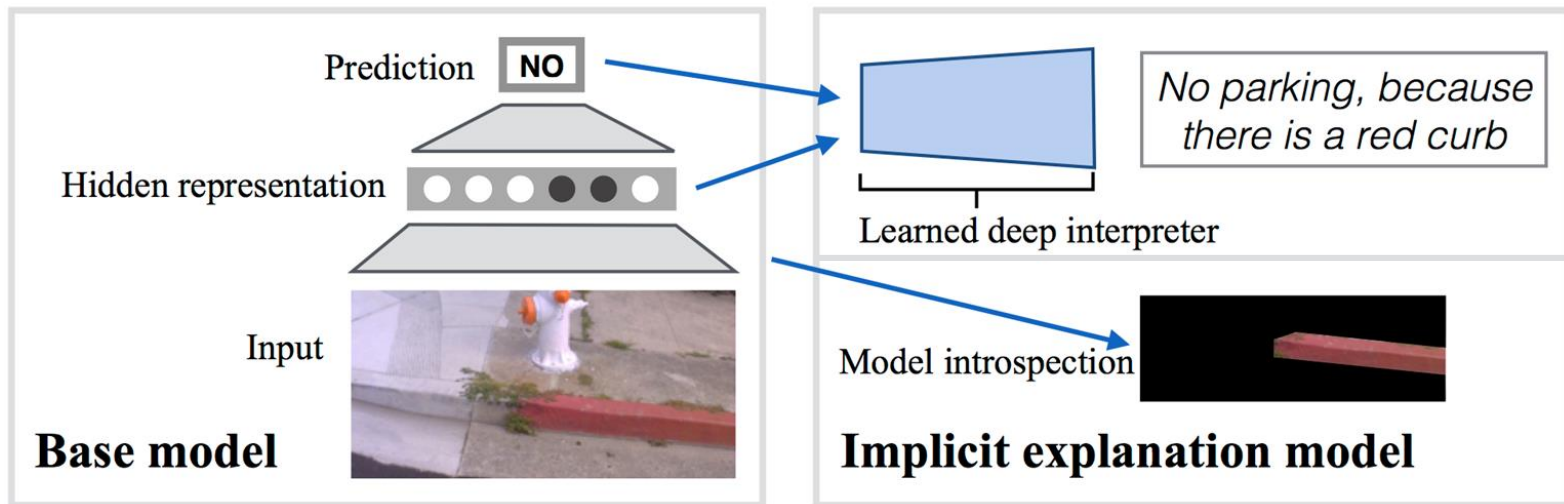


Learning how to explain  
(or, Talking to the User)

# Implicit Model for Textual Explanations

Translate DNN hidden state into  
visualizations and exemplars  
human-interpretable language

**Can you park here?**



# Textual Explanations





# Generating Textual Explanations

## Deep Finegrained Classifier



Compact  
Bilinear  
Classifier

Compact Bilinear  
Feature

Image Category:  
*Cardinal*

Input Sentence  
"A bright red bird with an  
orange beak."

Concat

$w_0$ : <SOS>

$w_1$ : A

$w_{T-1}$ : beak

LSTM

LSTM

LSTM

LSTM

LSTM

LSTM

LSTM

LSTM

$p(w_1|w_0, I, C)$

$p(w_2|w_{0:1}, I, C)$

$p(w_T|w_{0:T-1}, I, C)$

## Discriminative Loss

Sampled Sentence:  
"A red bird with black  
cheeks."

Sentence  
Classifier

Image Category:  
*Cardinal*

Reward  
Function

## Relevance Loss

$p(w_1|w_0, I, C)$   
 $p(w_2|w_{0:1}, I, C)$   
...  
 $p(w_T|w_{0:T-1}, I, C)$

Target Sentence  
"A bright red bird with an  
orange beak."

Cross Entropy  
Loss

# Definition vs. Explanation: Qualitative Results

*This is a **Downy Woodpecker** because...*



Definition: this bird has a white breast black wings and a red spot on its head.

Explanation: this is a black and white bird with a **red spot** on its crown.

*This is a **Downy Woodpecker** because...*



Definition: this bird has a white breast black wings and a red spot on its head.

Explanation: this is a white bird with a black wing and a **black and white** striped head.

# Definition vs Explanation: Failure Cases

**Correct class:** Laysan Albatross, **Predicted class:** Cactus Wren



Explanation: ...this is a **brown and white spotted** bird with a long pointed beak.

*Cactus Wren* Definition: This bird has a long thin beak with a **brown body** and black spotted feathers.

*Laysan Albatross* Definition: This bird has a **white head and breast** a grey back and an orange beak.

**Correct class:** Laysan Albatross, **Predicted class:** Laysan Albatross



Explanation: ...this bird has a **white head and breast** with a long hooked bill.

*Laysan Albatross* Definition: This bird has a **white head and breast** a grey back and an orange beak.

# How to Couple Visual and Textual Data?



CNN

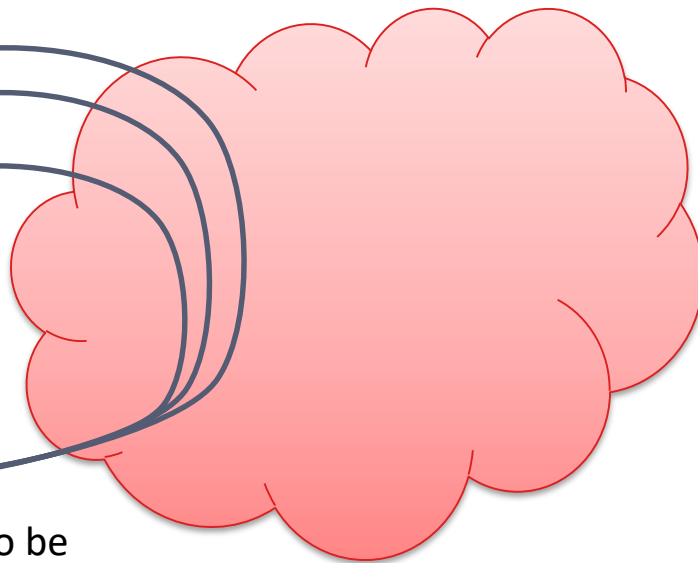
spoon  
plate  
bowl  
table  
food  
corn  
...

person

LSTM

Is?  
feast  
going to be  
...

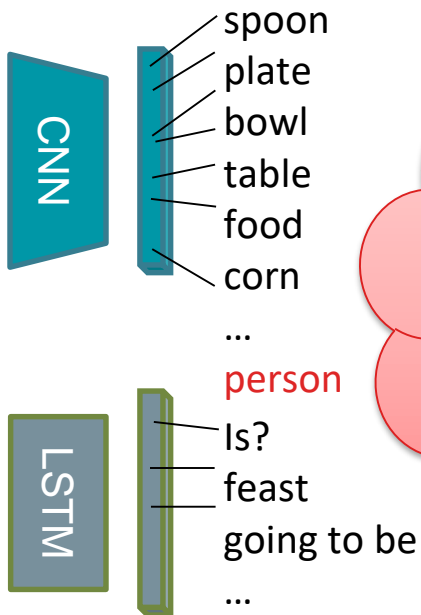
*Is this going to  
be a feast?*



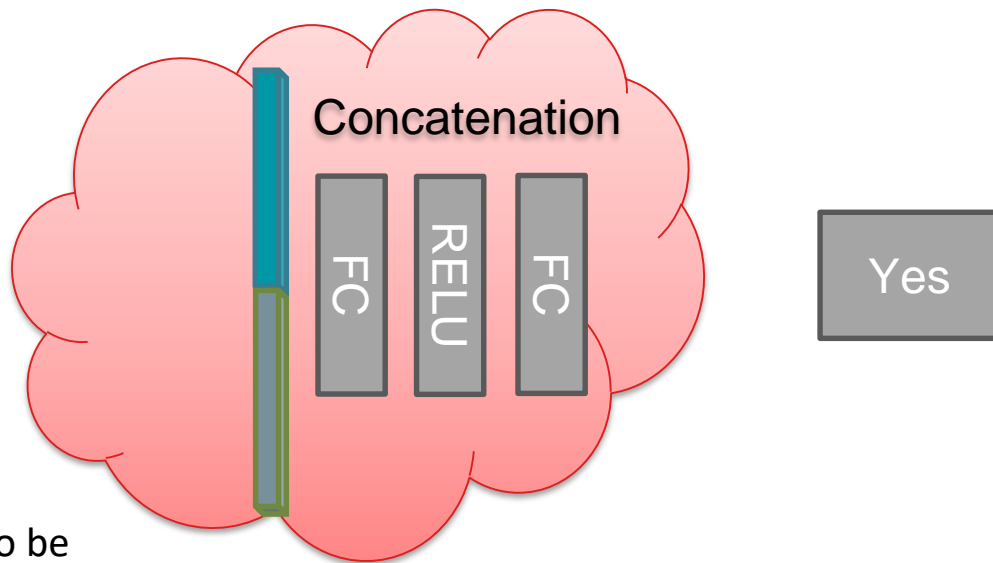
Yes

- ☐ All elements can interact
- ☐ Multiplicative interaction

# How to Couple Visual and Textual Data?

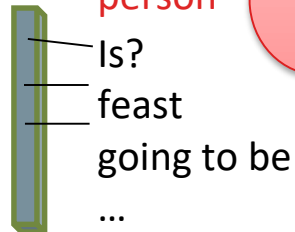
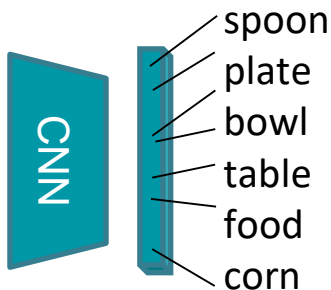


*Is this going to be a feast?*

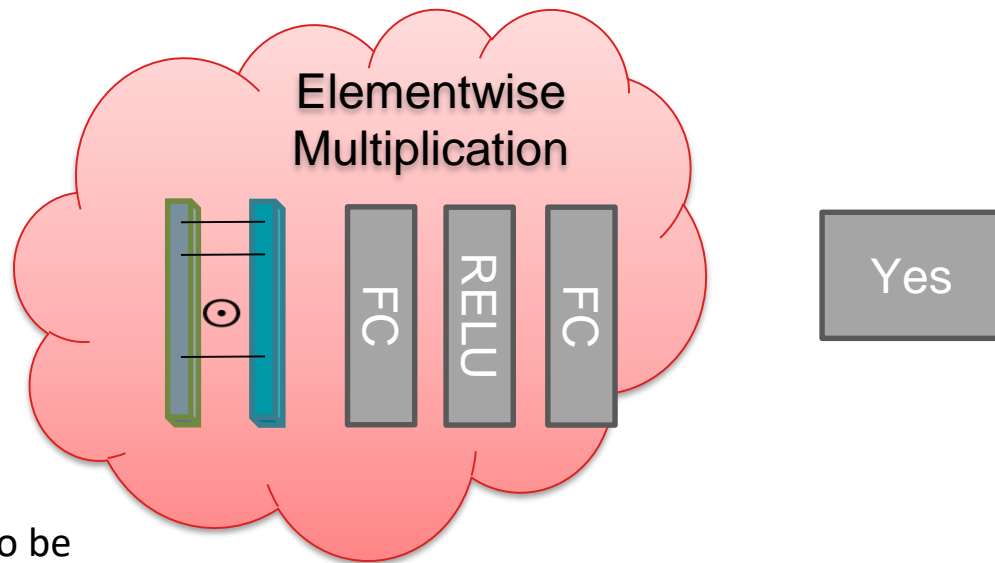


- ✓ **All elements can interact**
- ❓ **Multiplicative interaction**
  - **Difficult to learn output classification**

# How to Couple Visual and Textual Data?



*Is this going to  
be a feast?*



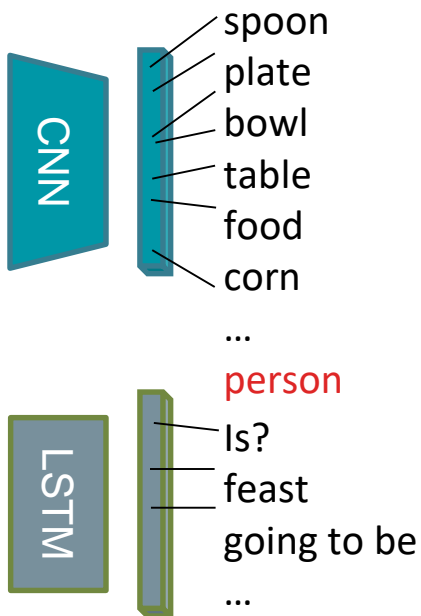
**❓ All elements can interact**

**✅ Multiplicative interaction**

- Difficult to learn input embedding

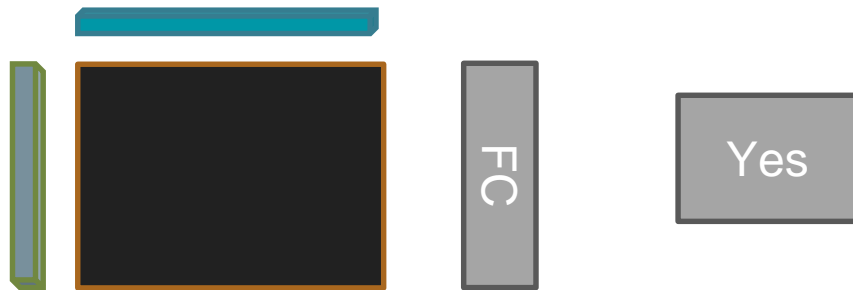


# How to Couple Visual and Textual Data?



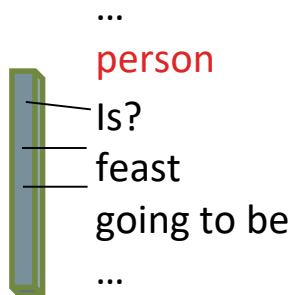
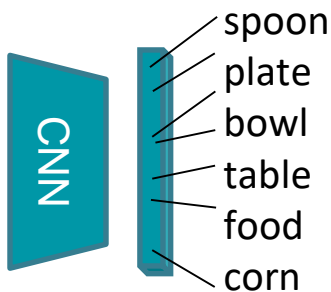
*Is this going to be a feast?*

Outer Product /  
Bilinear Pooling



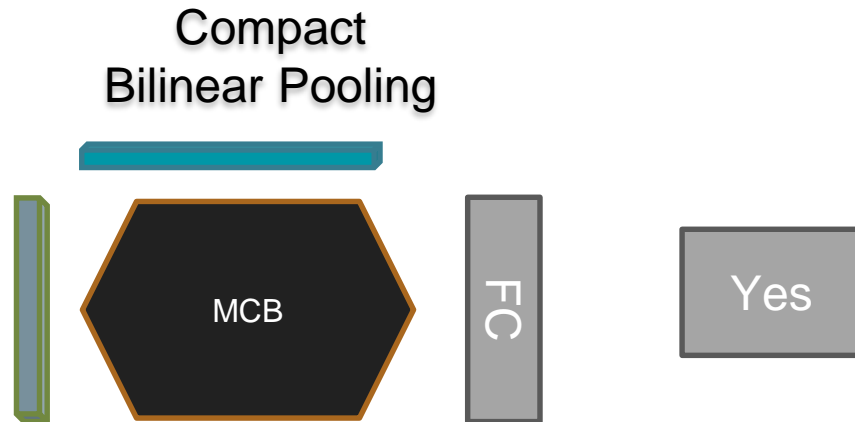
- ✓ All elements can interact
- ✓ Multiplicative interaction
- ⚠ High #activations & computation
- ⚠ High #parameters

# How to Couple Visual and Textual Data?



*Is this going to  
be a feast?*

LSTM



- ✓ All elements can interact
- ✓ Multiplicative interaction
- ✓ Low #activations & computation
- ✓ Low #parameters

[Zhang, Shelhamer, Gao, Darrell; ICLR workshop 2016]

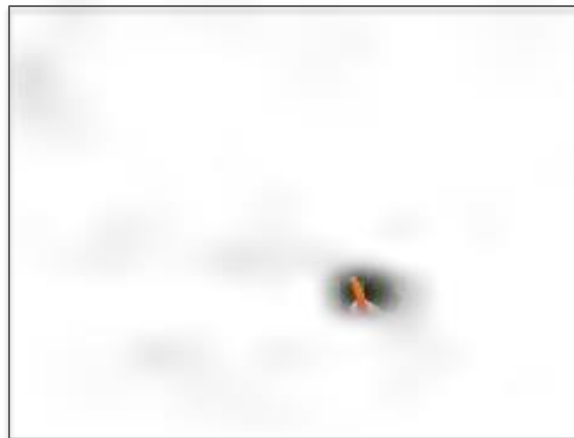
[Gao, Beijbom, Zhang, Darrell, CVPR 2016]



# Attention Visualizations

What is the woman **feeding** the giraffe?

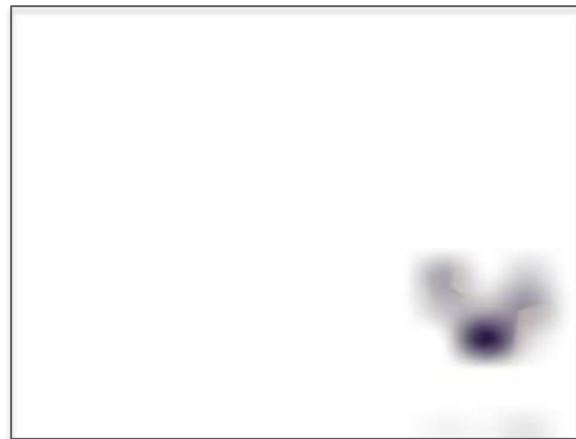
**Carrot**



# Attention Visualizations

What color is her **shirt**?

**Purple**



# Attention Visualizations

What is her **hairstyle** for the picture?

**Ponytail**

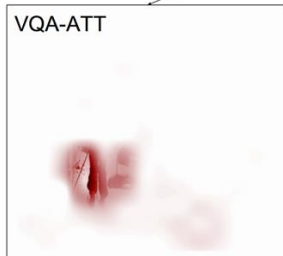


# Justifying Decisions and Pointing to the Evidence

*Q: What is the person doing?*

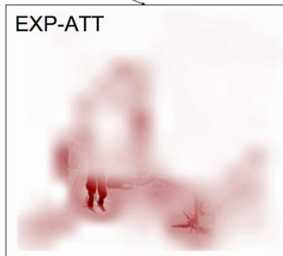


VQA-ATT



**A: Skiing**

EXP-ATT



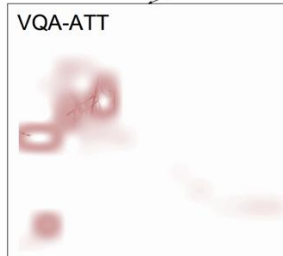
**Because:**

They are on **skis** and going down a **mountain**

*Q: What is the person doing?*

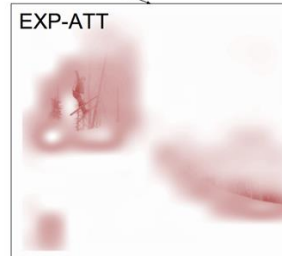


VQA-ATT



**A: Skiing**

EXP-ATT



**Because:**

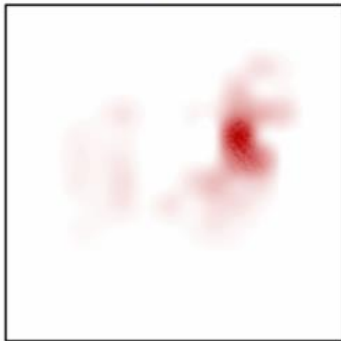
He is on a snowy **hill** wearing **skis** and **clothing** appropriate for skiing

# Attentive Explanations of VQA

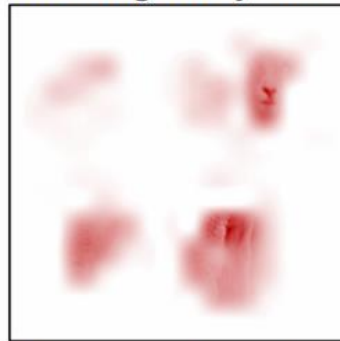
Q: What type of animal is this?



A: Sheep



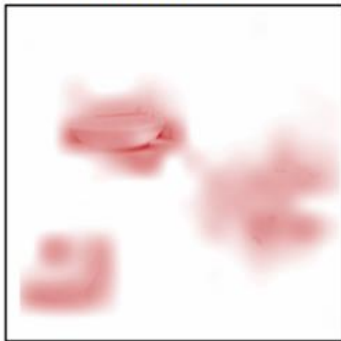
Because... it has four legs and long fluffy hair



Q: What room is this?



A: Bathroom

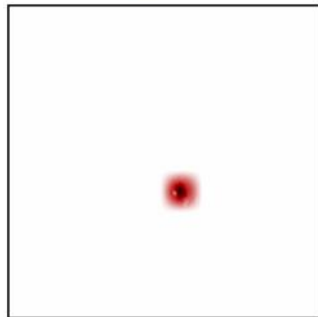


Because... there is a toilet and sink in the room

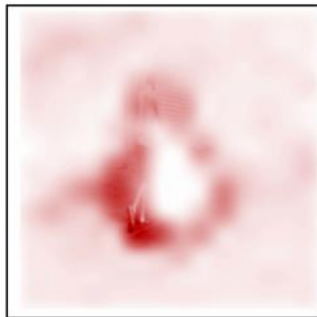


# Attentive Explanations of Activities

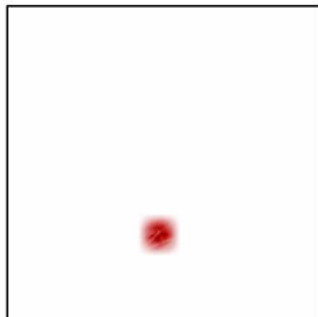
I can see that he is  
mowing lawn



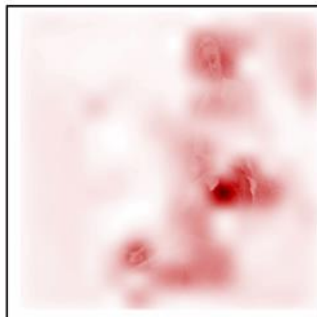
Because... he is pushing a  
lawn mower over a grassy  
lawn



I can see that he is  
mowing lawn



Because... he is kneeling in  
the grass next to a lawn  
mower

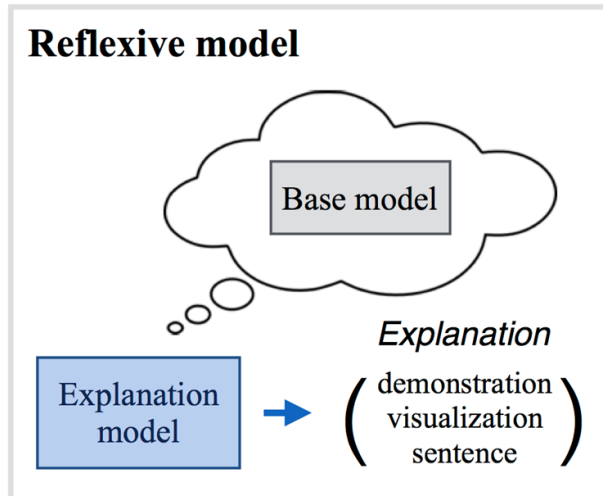


# Modeling the User

# Reflexive and Rational Models

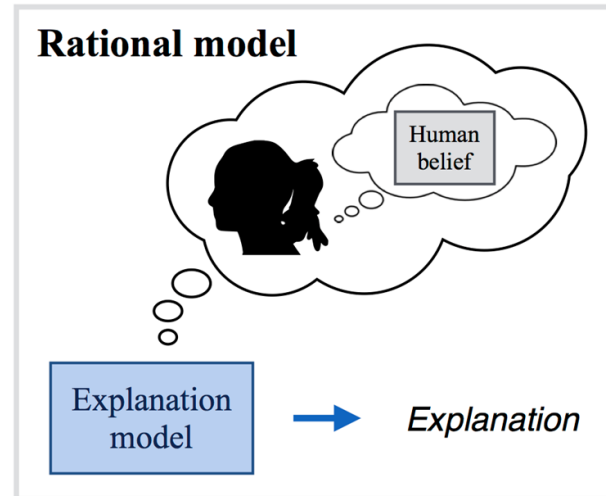
## Reflexive Agents

from examples by humans  
independent of user state.



## Rational Agents

how the system makes predictions  
why a mistake was made in a scenario



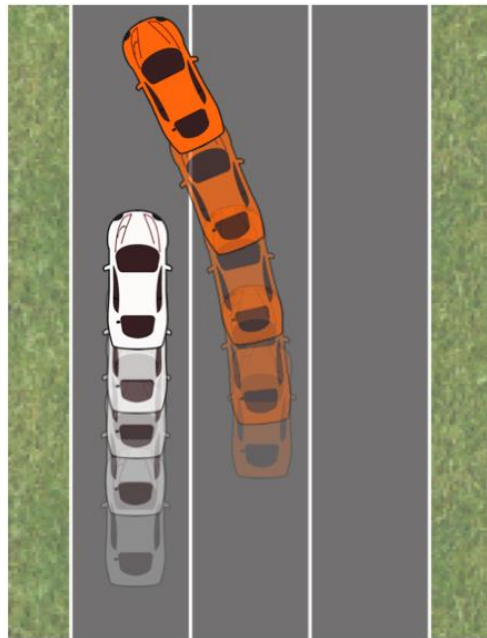


# Explainable Models for Dynamic Policies

## Dynamic policies

provide example state/action sequences

help user understand model behavior.



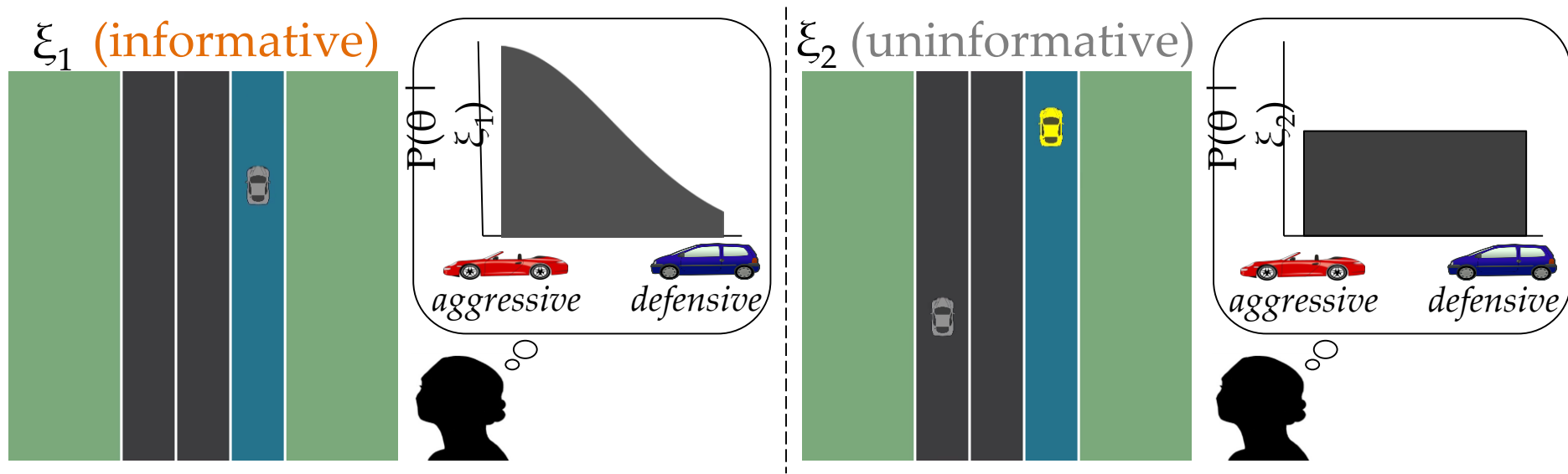
(a) Explain learned policy  
via example rollouts

How can we help users better anticipate  
what a robot will do?

Key insight:  
Users need to understand the *tradeoffs*  
that a robot makes

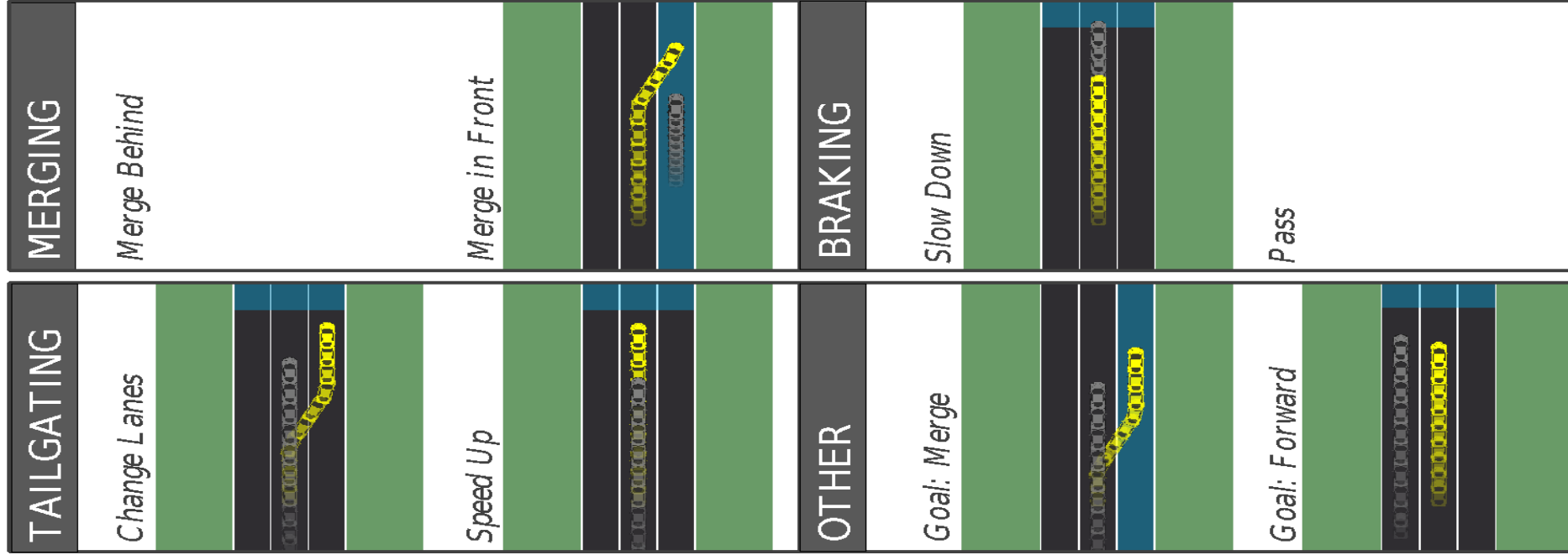
# Robots Inevitably Communicate via Their *Behavior*

Humans naturally reason about others' utility functions\*



\*Jara-Ettinger, J., et al. Trends in Cognitive Sciences (2016)

# Experiment Setup: Environments



# User Study: Noise Models

## Current Test Environment

Please watch all four video clips below. One video clip shows Carl driving, and the other three show imposter cars that look the same as Carl but drive differently. (Click on the video to start it.)



Which video clip do you think shows Carl driving?

- ☐ Video 1    ☐ Video 2    ☐ Video 3    ☐ Video 4

To what extent do you agree or disagree with the following statements?

The car in **Video 1** drives in a similar way as Carl.

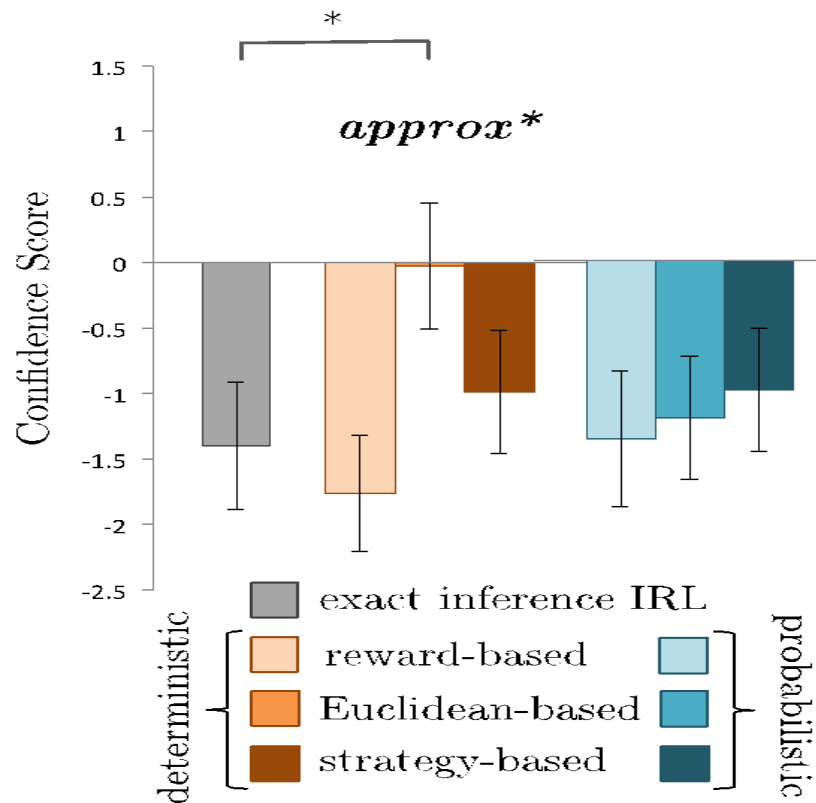
- ☐ Strongly disagree    ☐ Disagree    ☐ Somewhat disagree    ☐ Neither agree nor disagree    ☐ Somewhat agree    ☐ Agree    ☐ Strongly agree

The car in **Video 2** drives in a similar way as Carl.

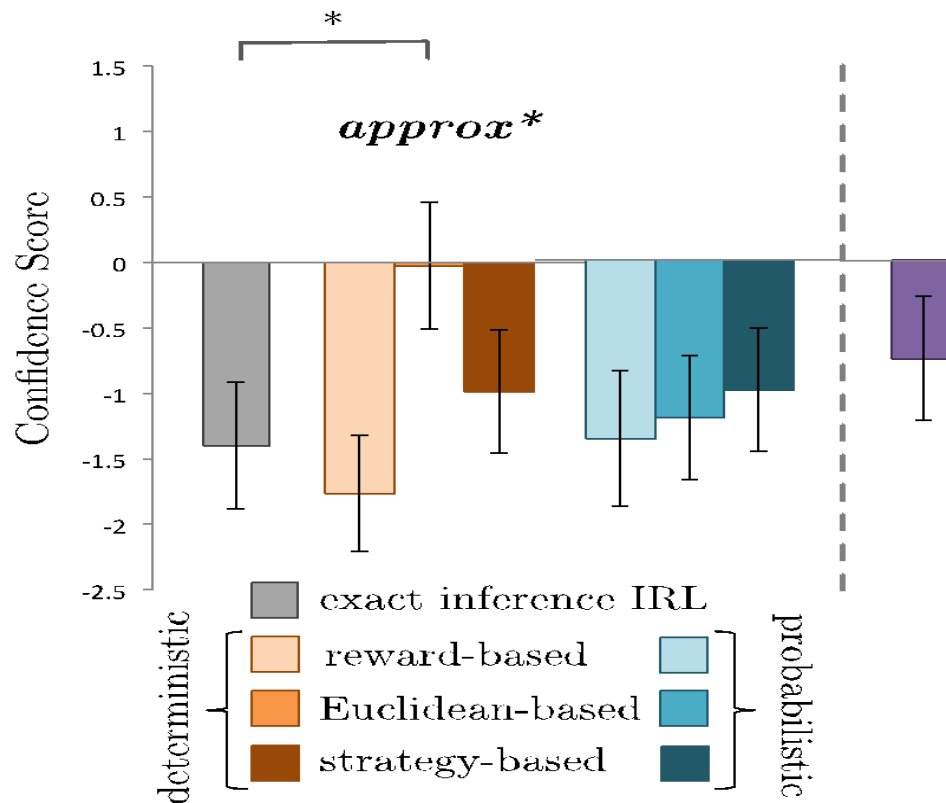
- ☐ Strongly disagree    ☐ Disagree    ☐ Somewhat disagree    ☐ Neither agree nor disagree    ☐ Somewhat agree    ☐ Agree    ☐ Strongly agree



# User Study: Noise Models



# User Study: Noise Models

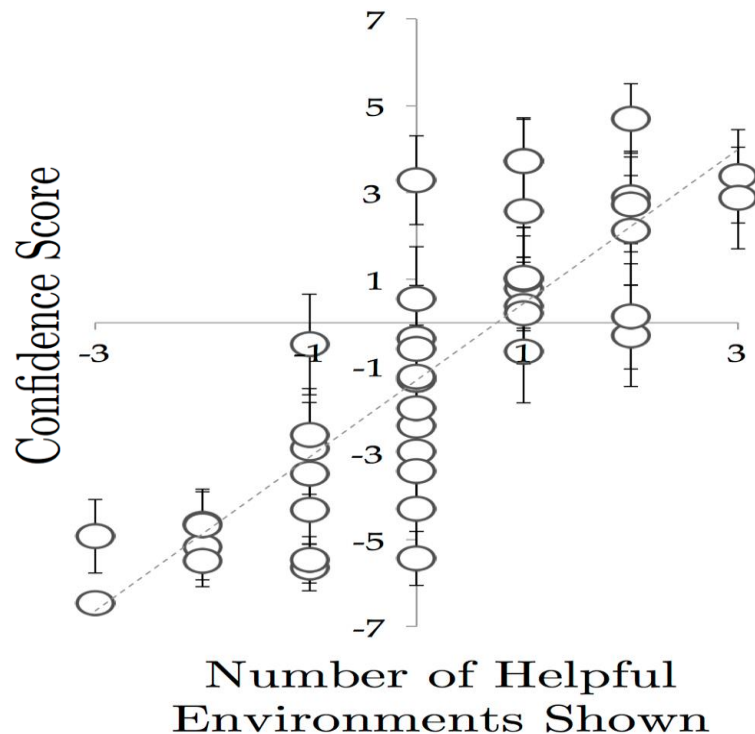


# Coverage Matters

Given  $x$  examples shown in strategy A and  $y$  from B:

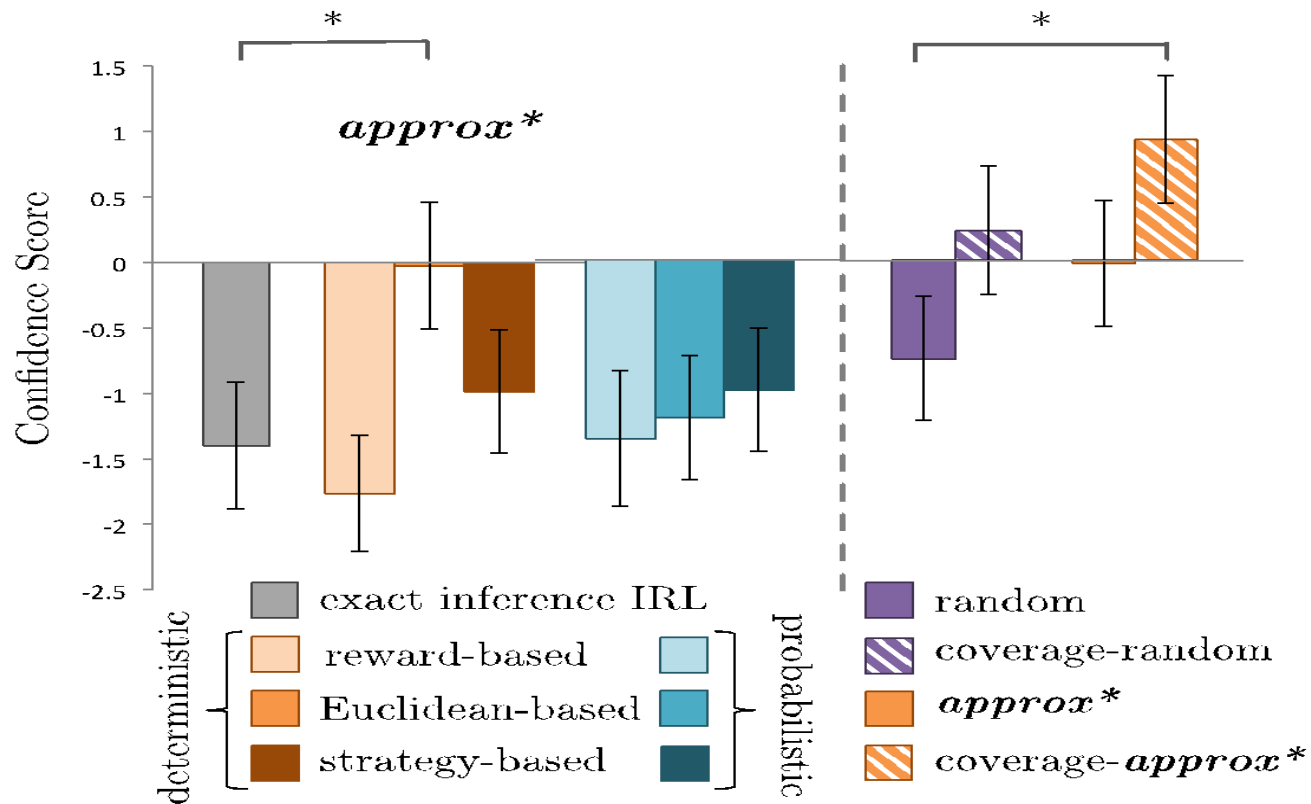
# *helpful environments* shown in A = 
$$\begin{cases} x, & \text{if } x > 0. \\ -y, & \text{otherwise.} \end{cases}$$

Pearson's  
 $r = 0.83$





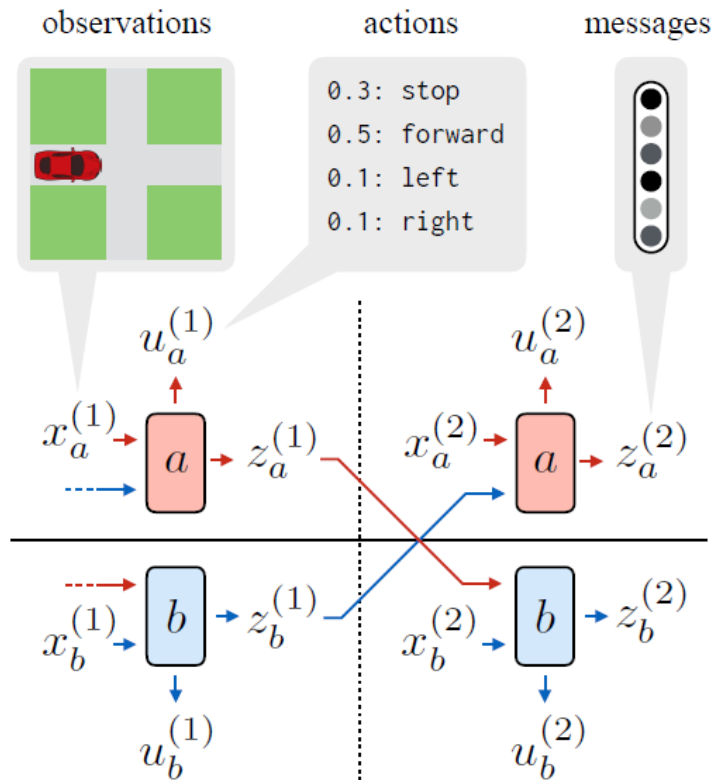
# User Study: Coverage



# What is Neuralese?

Idea: induce communication protocols for multiagent tasks

Example: navigating a contested intersection



# Translating Neuralese

- Idea: infer a mapping between neuralese and natural language
  - (1) Have neural system learn a code for a task
  - (2) Have humans do the same task, communicating in NL
  - (3) Compute mapping from neuralese to NL using *belief matching*
- Belief matching
  - Q: How do you know x in neuralese means the same as y in NL?
  - A: If they induce the optimally similar distributions over belief states

$$q(z, z') = \mathbb{E}[\mathcal{D}_{\text{KL}}(\beta(z, X_b) || \beta(z', X_b)) | z, z']$$

- (Applies more generally to knowing two messages mean the same thing)

# Examples

Task: Distinguish the starred image



*large bird, black wings, black crown*

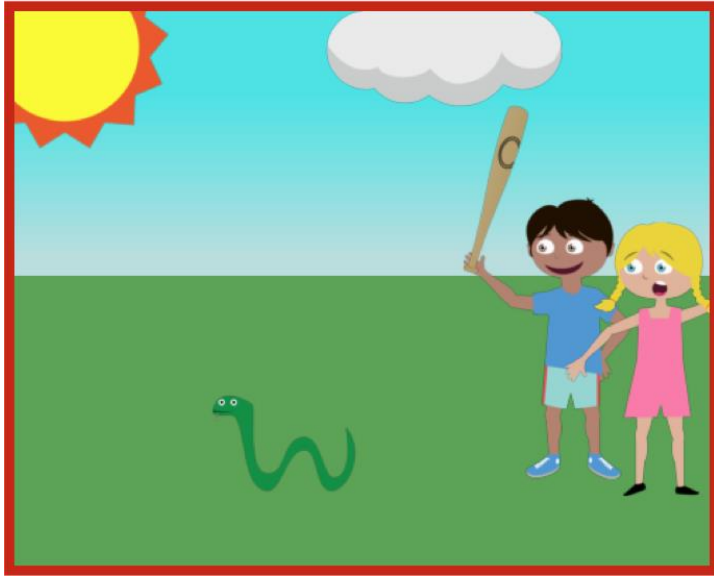
Task: Avoid collision with invisible car



*you first, following, going down*

# The reference game

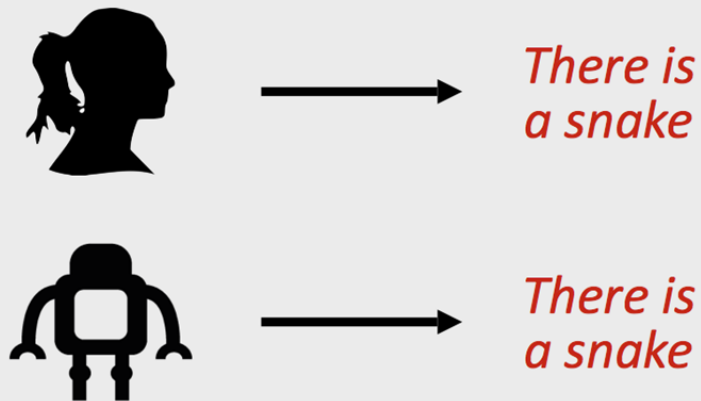
*Mike is holding a baseball bat*



# The reference game

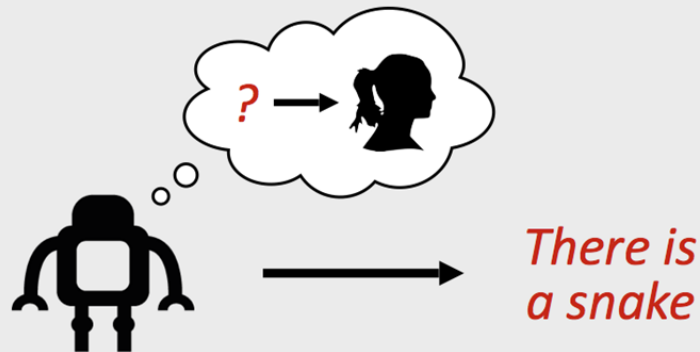
## DIRECT APPROACH:

Imitate successful human play

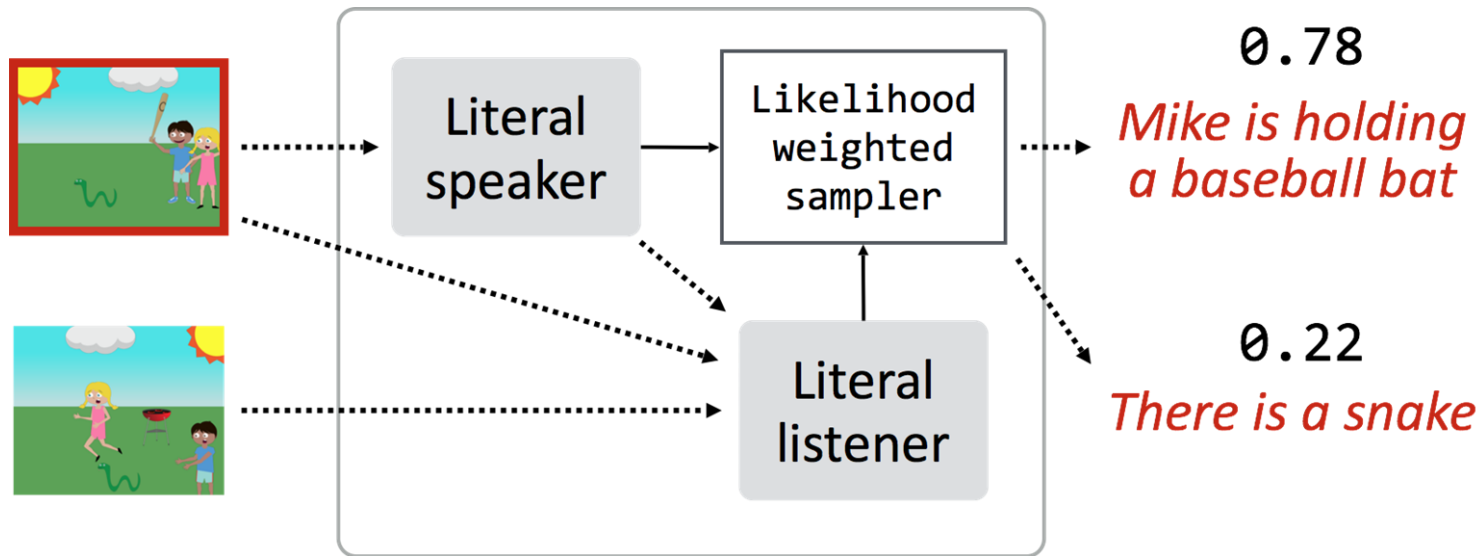


## DERIVED APPROACH:

Reason about listener beliefs



# The reference game



# Program Schedule



# Program Schedule

## Phase I

- algorithmic development and detailed specification

- definition, and initial evaluation of the above challenge areas and datasets

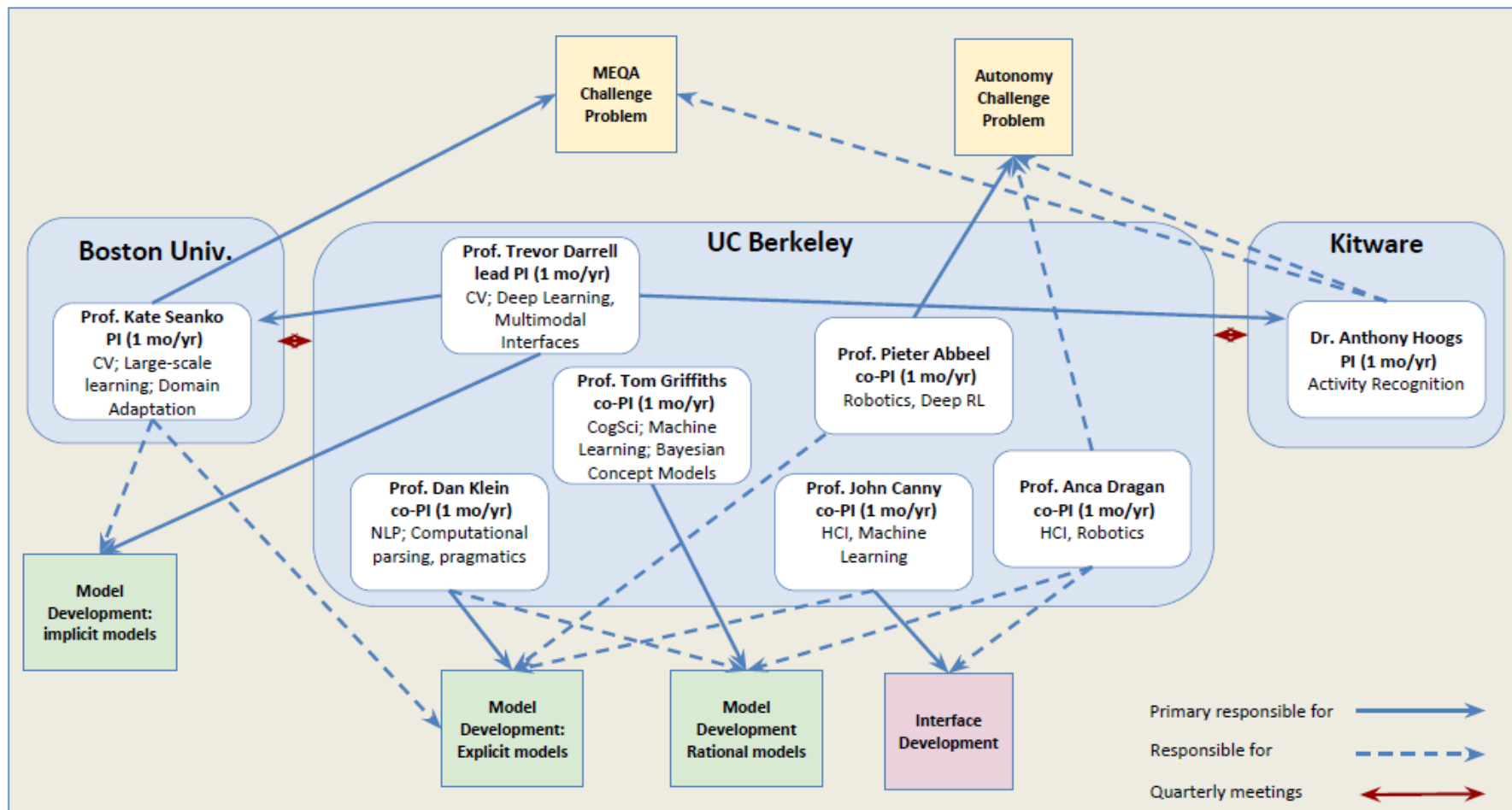
## Phase II

- performance improvements

- system demonstrations

- in-situ user-oriented evaluations

- preparation for transition



Thank you!