# Lecture 9.2

Embodiment

# Analysis of baselines

## Questions to Ask

- How well do they do on overall task performance

- How well do they do on intrinsic metrics? Can simple models at least memorize the dataset?

- Error distributions:
  - Are baselines doing better on certain classes of questions/problems than others?
  - What aspects of the data analysis that you did previously are hurting your models now? Is it OOV? Is it unknown objects in BBoxes?

## Things to Write

- Several copied / trained rows in the results table

- A second table or two for intrinsic metrics

- A table/plot or two of error types

- Ideally, a few qualitative examples

- Most important: Insights! Did any of the models do better or worse than expected on certain aspects of the data? Does this align with your intuition from data analysis?

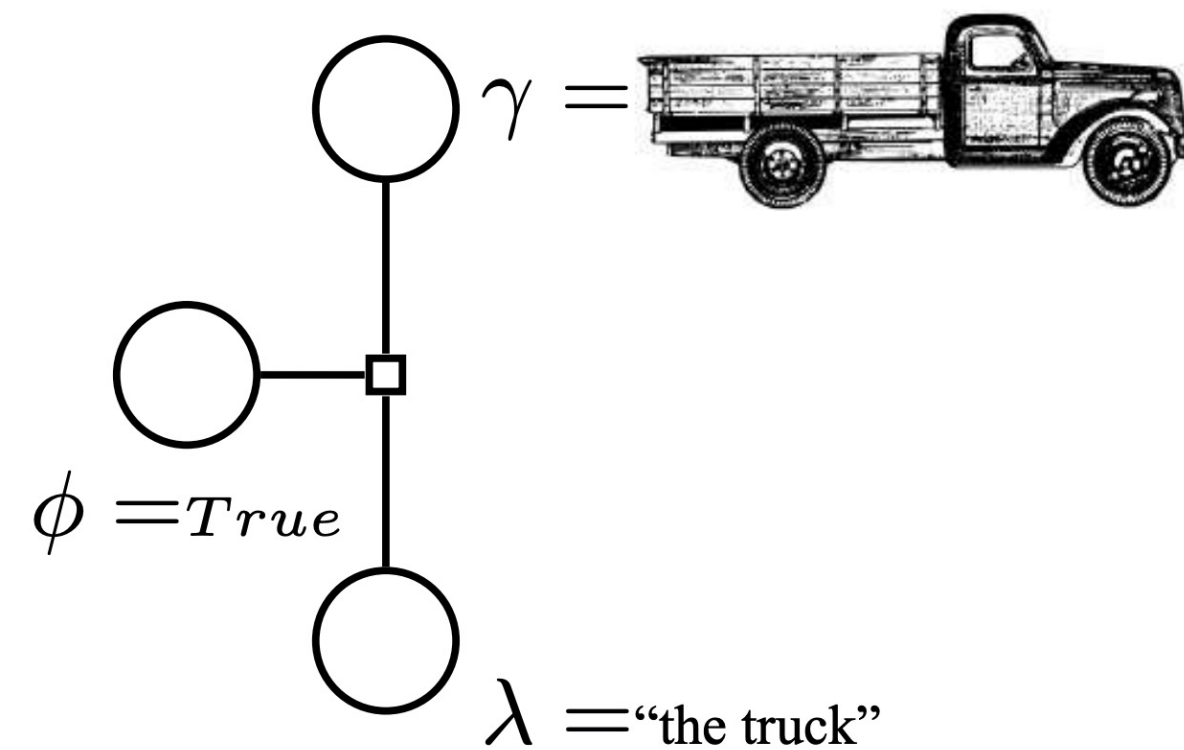# Remainder of the Semester (a lot and a little)

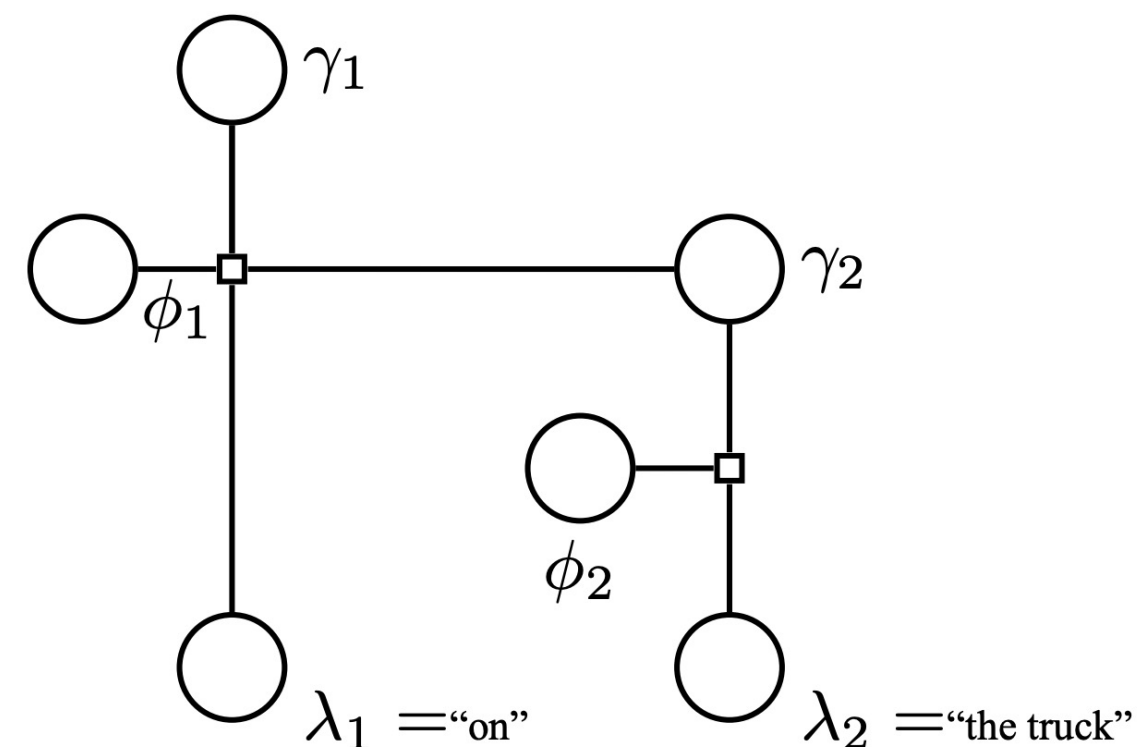| | |
|---|---|
| Mar 30: RL | Apr 1: Multimodal RL |
| Apr 6: Project Hours (R5) | Apr 8: Project Hours (R5) |
| Apr 13: Fusion and Co-Learning | Apr 15: — No Class — Carnival — |
| Apr 20: New Research Directions | Apr 22: TBD?!? |
| Apr 27: Project Hours (Final) | Apr 29: Project Hours (Final) |
| May 4: Guest (Bias in V+L): Mark Yatskar @ UPenn | May 6: Guest (Robotics): Chris Paxton @ NVIDIA |
| May 11: Project Presentations | May 13: Reports Due |

**Carnegie Mellon University** Language Technologies Institute

# Instruction Following



Language Technologies Institute

http://people.csail.mit.edu/stefie10/publications/tellex13.pdf

# Symbol Grounding

$OBJ(f = \texttt{the truck})$



$\gamma = $ [truck image]

$\phi = True$

$\lambda = \text{"the truck"}$

$PLACE_2(r = \texttt{on}$
$\qquad l1 = OBJ_1(f = \texttt{the truck}))$

$\gamma_1$

$\phi_1$

$\gamma_2$

$\phi_2$

$\lambda_1 = \text{"on"}$ $\qquad \lambda_2 = \text{"the truck"}$

$EVENT_1(r = \texttt{Put},$
$\qquad\qquad l = OBJ_2(f = \texttt{the pallet}),$
$\qquad\qquad l2 = PLACE_3(r = \texttt{on},$
$\qquad\qquad\qquad\qquad l = OBJ_4(f = \texttt{the truck})))$

(a) SDC tree

$\gamma_1$

$\phi_1$

$\gamma_2$ $\qquad \gamma_3$

$\phi_2$ $\qquad \phi_3$ $\qquad \gamma_4$

$\phi_4$

$\lambda_1^r$ $\quad \lambda_2^f$ $\qquad\qquad \lambda_3^r$ $\quad \lambda_4^f$

"Put" $\quad$ "the pallet" $\qquad$ "on" $\quad$ "the truck"

(b) Induced Model

$$p(\Phi|\Gamma,\text{SDCs}, m) = p(\phi_1|\gamma_1, \gamma_2, \gamma_3, \lambda_1^r = \text{Put}, m) \times$$
$$p(\phi_2|\gamma_2, \lambda_2^f = \text{the pallet}, m) \times p(\phi_3|\gamma_3, \gamma_4, \lambda_3^r = \text{on}, m) \times$$
$$p(\phi_4|\gamma_4, \lambda_4^f = \text{the truck}, m)$$

(c) Factorization

**Carnegie Mellon University** Language Technologies Institute $\qquad$ http://people.csail.mit.edu/teller/pubs/TellexEtAIAImagazine2011.pdf

# Sim2Real Language -> Control

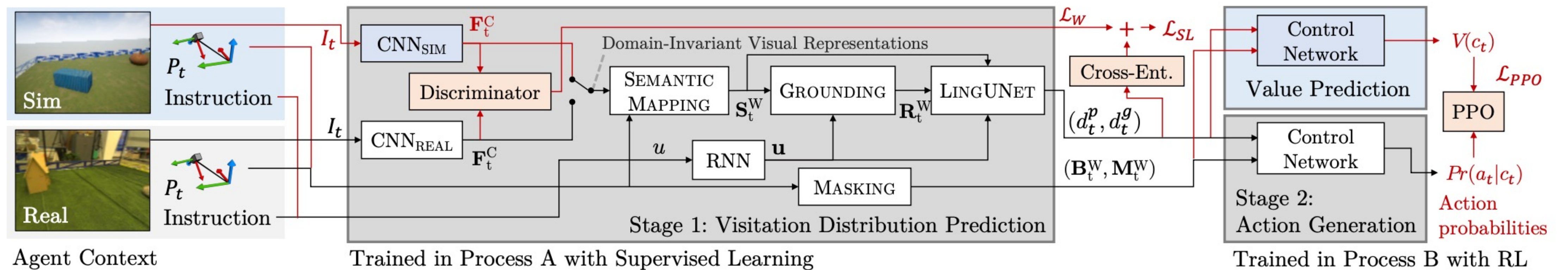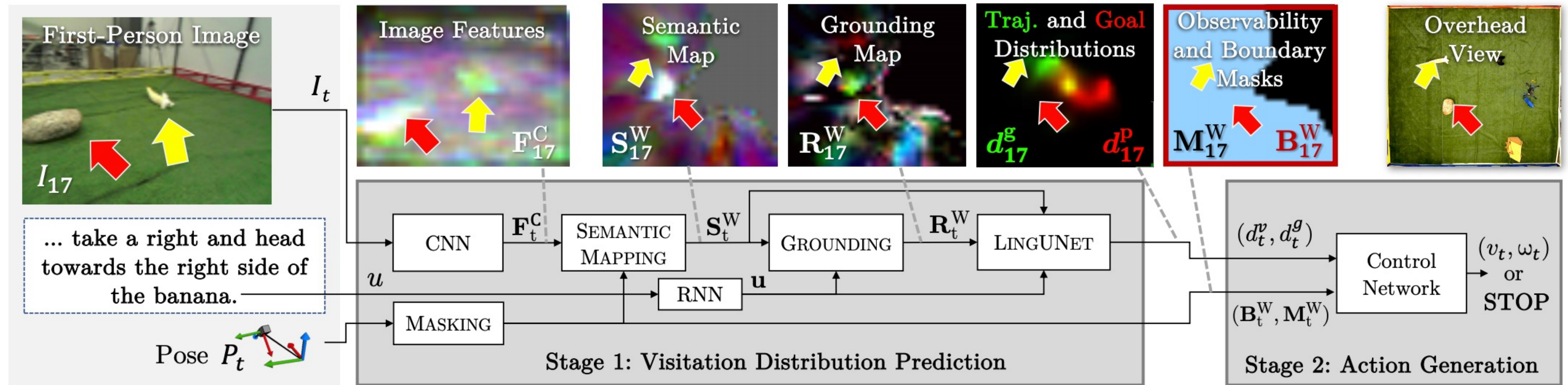Once near the rear of the gorilla, turn right and head towards the rock stopping once near it



Input Image

3rd Person View
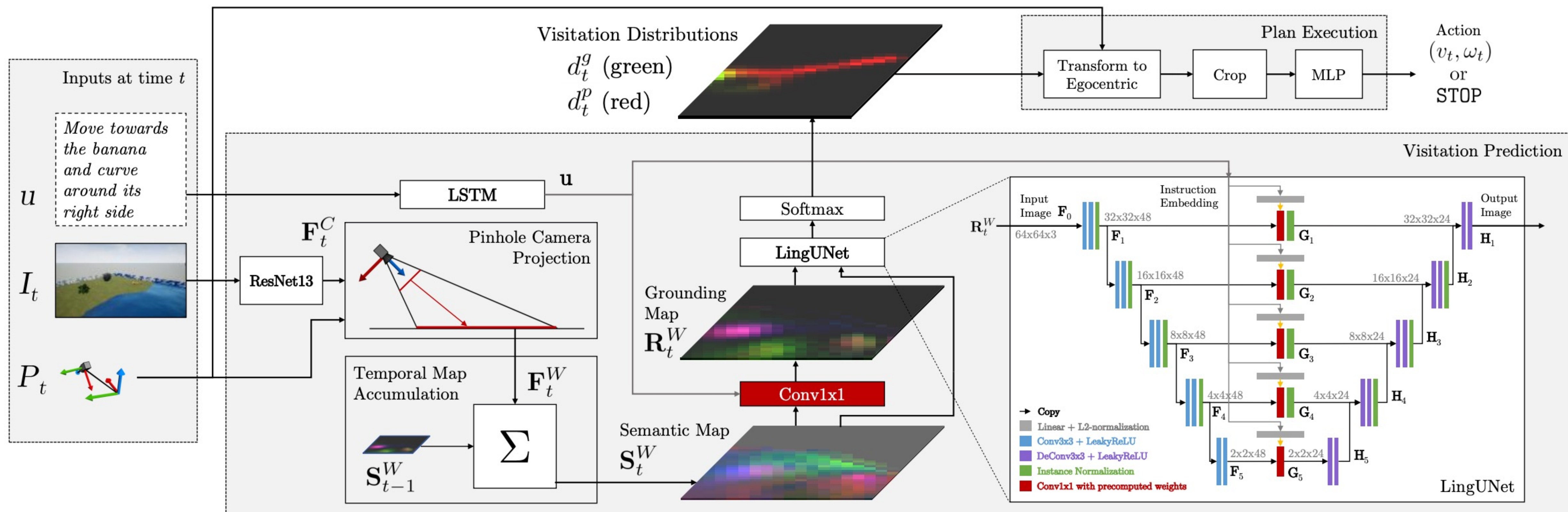
Continuous velocity commands
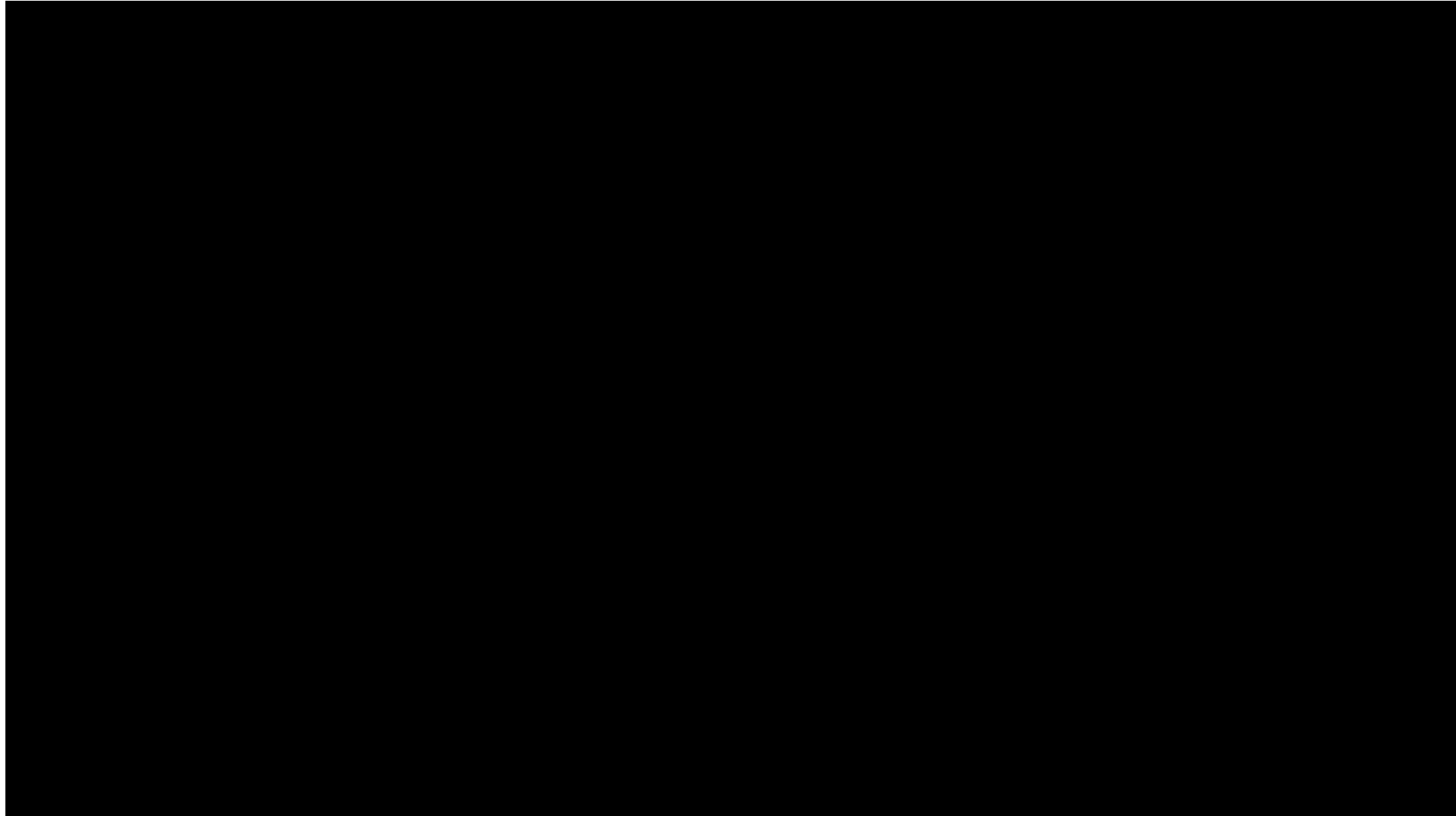
Forward Velocity

Yaw Rate

# Building a Map



First-Person Image — $I_t$ — $I_{17}$

... take a right and head towards the right side of the banana.

$u$

Pose $P_t$

Image Features — $\mathbf{F}^{C}_{17}$

Semantic Map — $\mathbf{S}^{W}_{17}$

Grounding Map — $\mathbf{R}^{W}_{17}$

Traj. and Goal Distributions — $d^{g}_{17}$ — $d^{p}_{17}$

Observability and Boundary Masks — $\mathbf{M}^{W}_{17}$ — $\mathbf{B}^{W}_{17}$

Overhead View

CNN — $\mathbf{F}^{C}_{t}$ — Semantic Mapping — $\mathbf{S}^{W}_{t}$ — Grounding — $\mathbf{R}^{W}_{t}$ — LinguNet

RNN — $\mathbf{u}$

Masking

$(d^{v}_t, d^{g}_t)$

$(\mathbf{B}^{W}_t, \mathbf{M}^{W}_t)$

Control Network — $(v_t, \omega_t)$ or STOP

**Stage 1: Visitation Distribution Prediction**

**Stage 2: Action Generation**

Sim — $P_t$ Instruction

Real — $P_t$ Instruction

$I_t$ — CNN$_{\text{SIM}}$ — $\mathbf{F}^{C}_{t}$

Discriminator

$I_t$ — CNN$_{\text{REAL}}$ — $\mathbf{F}^{C}_{t}$

Domain-Invariant Visual Representations

Semantic Mapping — $\mathbf{S}^{W}_{t}$ — Grounding — $\mathbf{R}^{W}_{t}$ — LinguNet

$u$ — RNN — $\mathbf{u}$

Masking

$\mathcal{L}_W$ — $+$ → $\mathcal{L}_{SL}$

Cross-Ent.

$(d^{p}_t, d^{g}_t)$

$(\mathbf{B}^{W}_t, \mathbf{M}^{W}_t)$

Control Network — Value Prediction — $V(c_t)$

$\mathcal{L}_{PPO}$ — PPO

Control Network — Stage 2: Action Generation — $Pr(a_t | c_t)$ Action probabilities

**Stage 1: Visitation Distribution Prediction**

Agent Context

Trained in Process A with Supervised Learning

Trained in Process B with RL

**Carnegie Mellon University** Language Technologies Institute

# LingUNet

https://arxiv.org/abs/1811.04179

# Sim-to-Real VLN

Sim-to-Real Transfer for
Vision-and-Language Navigation

1

# Knowledge Acquisition

# SoundScapes

# Why?

## Language that affects the world

*Remove the cream from the middle of the Oreo…*



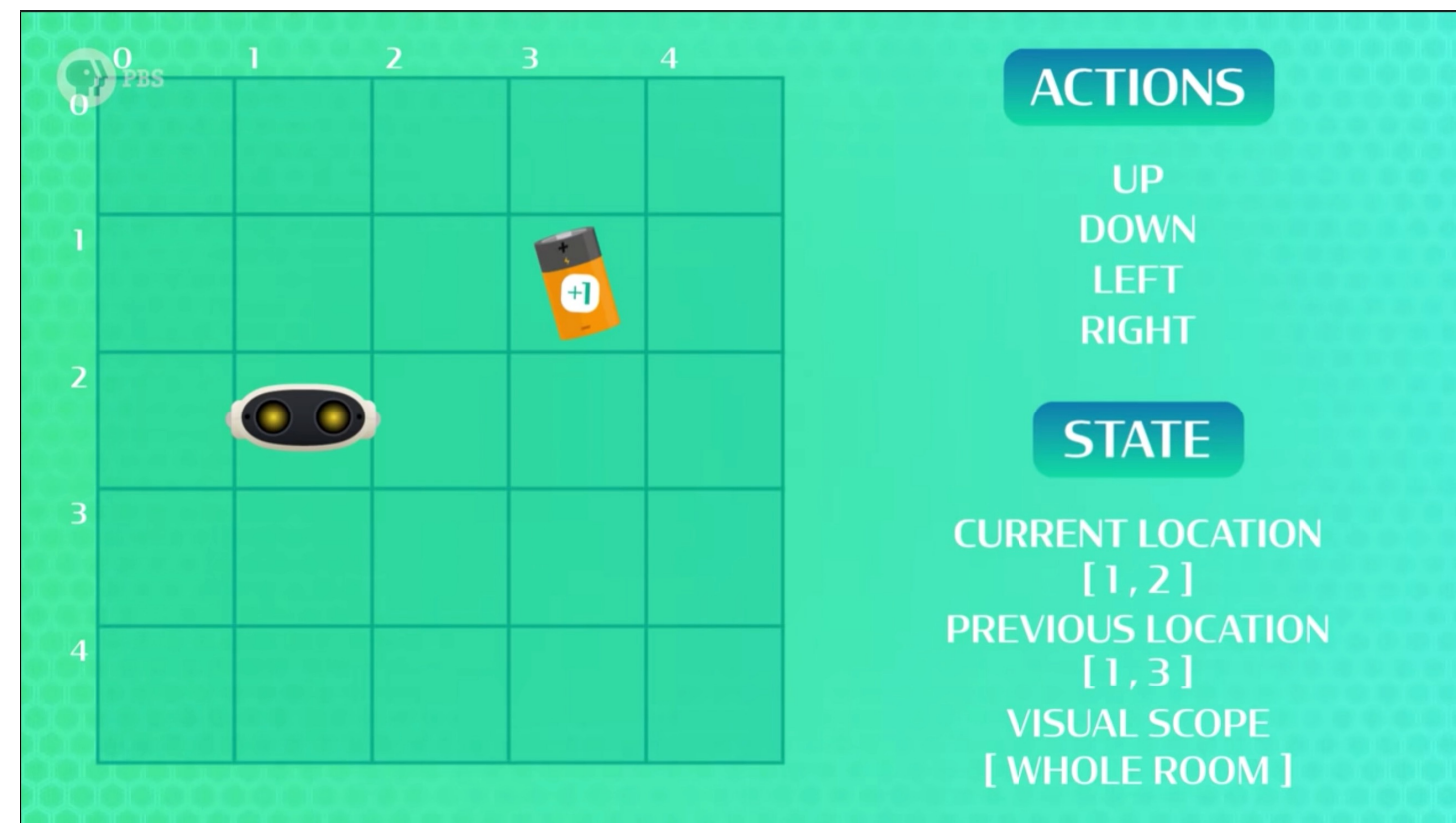HERB (Siddhartha Srinivasa)

## Access to Broader Semantics

What's it like to drive a bus?



How many hours of watching to achieve same level of performance as 30m of practice?

https://www.youtube.com/watch?v=SqGdH8lWvbA

# What does interaction mean?

Grid World?



Reinforcement Learning: Crash Course AI#9
https://www.youtube.com/watch?v=nIglv4IfJ6s

Graph Navigation?



Anderson 2018

Manipulation?



Paxton 2019

1. How does the agent move?
2. How many arms or legs does it have?
3. How many fingers (if any) do the grippers have?
4. How many joints do the limbs have?
5. What about physics? Real motor noise?

…

Carnegie Mellon University Language Technologies Institute

# Every Dimension Interacts



1. How rich or abstract is the language?
2. How complex is the visual field?
3. Is the vision 2D, 3D, Lidar, … ?
4. What kind of supervision do you have?
 …

# Choose your own adventure

# Sequential and Online Modeling

Action Recognition



**3D Conv**

"Action Summary"

$$p(\text{Action}|v_0, ..., v_t)$$

Embodied

Action

Action

??

$$p(v_t|v_0, ..., \text{Action})$$

Requirement: Have a goal

# What is a "goal"?

"Put the green dog on the table"

$$p(v_t | v_0, ..., \text{Action})$$

$$p(v_t | v_0, ..., v_{t-1}, a_0, ..., a_t)$$

$v_t =$

# Planning
## Pre- and Post-Conditions

Task 4: Must locate object,
to move to object

Task 3: Must move to object,
to hold object

Task 2: Must hold object,
to place object

Task 1: Recognize Success

Instances of "green dog sculpture on table"



So what are we actually optimizing? What's our actual goal?

**Carnegie Mellon University** Language Technologies Institute

# Let's Start Simple

# Instruction Following

## Explicit Action Supervision

Walk out of the bedroom through the open door into the hallway

Turn the corner and walk into the dining area.

Pass the dining table and walk into the living room area towards the television.

Stop near the chair and open sliding doors to outside

# V+L -> A

Turn left

LEFT

and go straight

FWD

FAST Agent

FAST Agent

FAST Agent

Does this actually need vision?

Does this understand plans?

No, this is ~Semantic Parsing

# V+L -> A

Walk out of the bedroom through the open door into the hallway

LEFT

FWD



FAST Agent

FAST Agent

FAST Agent

Does this actually need vision?        Yes

Does this understand plans?        Maybe, probably not

# First Major Question: Alignment



Exit the bedroom and go towards the table. Go to the stairs on the left of the couch. Wait on the third step.

Ma et al, "Self-Monitoring Navigation Agent via Auxiliary Progress Estimation" ICLR 2019

# Alignment

Exit the bedroom and go towards the table. Go to the stairs on the left of the couch. Wait on the third step.





### Textual grounding

Exit the bedroom and go towards the table. Go to the stairs on the left of the couch. Wait on the third step.

embedded ↓ instructions

$x_1$ $x_2$ $x_3$ • • • $x_L$

Positional encoding

Soft-attention

Grounded instruction

$\hat{x}_t$

$\alpha_t$

$h_{t-1}$

Projection $W_x$        Projection $W_v$

Prev. action $a_{t-1}$

$\hat{x}_t$     $h_t$     LSTM     $\hat{v}_t$

$c_t$

Action selection

$a_t$

distance to goal $p_t^{pm}$     $h_t^{pm}$

Progress monitoring

### Visual grounding

feature     extraction

$v_{t,1}$  $v_{t,2}$  • • •  $v_{t,K}$

$g(v_{t,1})$ $g(v_{t,2})$ • • • $g(v_{t,K})$

Soft-attention

$\beta_t$

Grounded img features

$\hat{v}_t$

**Baseline (panoramic action space)**

validation seen                    validation unseen

Agent step

instruction: $x_1, x_2, ..., x_L$

**Self-Monitoring**

validation seen                    validation unseen

**Carnegie Mellon University** Language Technologies Institute

# Lots of Data

## Lots and lots of aligned data?

### Wait, remember the bus driver question?



Our starting point is in a living room, we're facing towards a long beige sofa, and in front of the sofa there are three glass coffee tables, turn around and exit through the doorway that's in front of you, walk pass the bed that's on your right and then turn left, we're now facing towards another living room, and on the left there's an open door, walk towards that open door enter the bathroom that's in front of you, turn towards the right into the shower area. and that's your destination.

| | Number of: | | | Includes: | | |
|---|---|---|---|---|---|---|
| | Lang | Instruct | Words | Paths | Text | Ground | Demos |
| CVDN | 1 | 2K[†] | 167K | 7K | ✓ | | |
| R2R | 1 | 22K | 625K | 7K | ✓ | | |
| Touchdown | 1 | 9K | 1.0M | 9K | ✓ | ✓[‡] | |
| REVERIE | 1 | 22K | 388K | 7K | ✓ | ✓[‡] | |
| RxR | 3 | 126K | 9.8M | 16.5K | ✓ | ✓ | ✓ |

[†]The number of dialogues. [‡]Grounding limited to one object per instruction.

Ku et al. Room-Across-Room: Multilingual Vision-and-Language Navigation with Dense Spatiotemporal Grounding — EMNLP 2020



**Carnegie Mellon University** Language Technologies Institute

# What if you make a mistake?

**The Frontier**

Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019
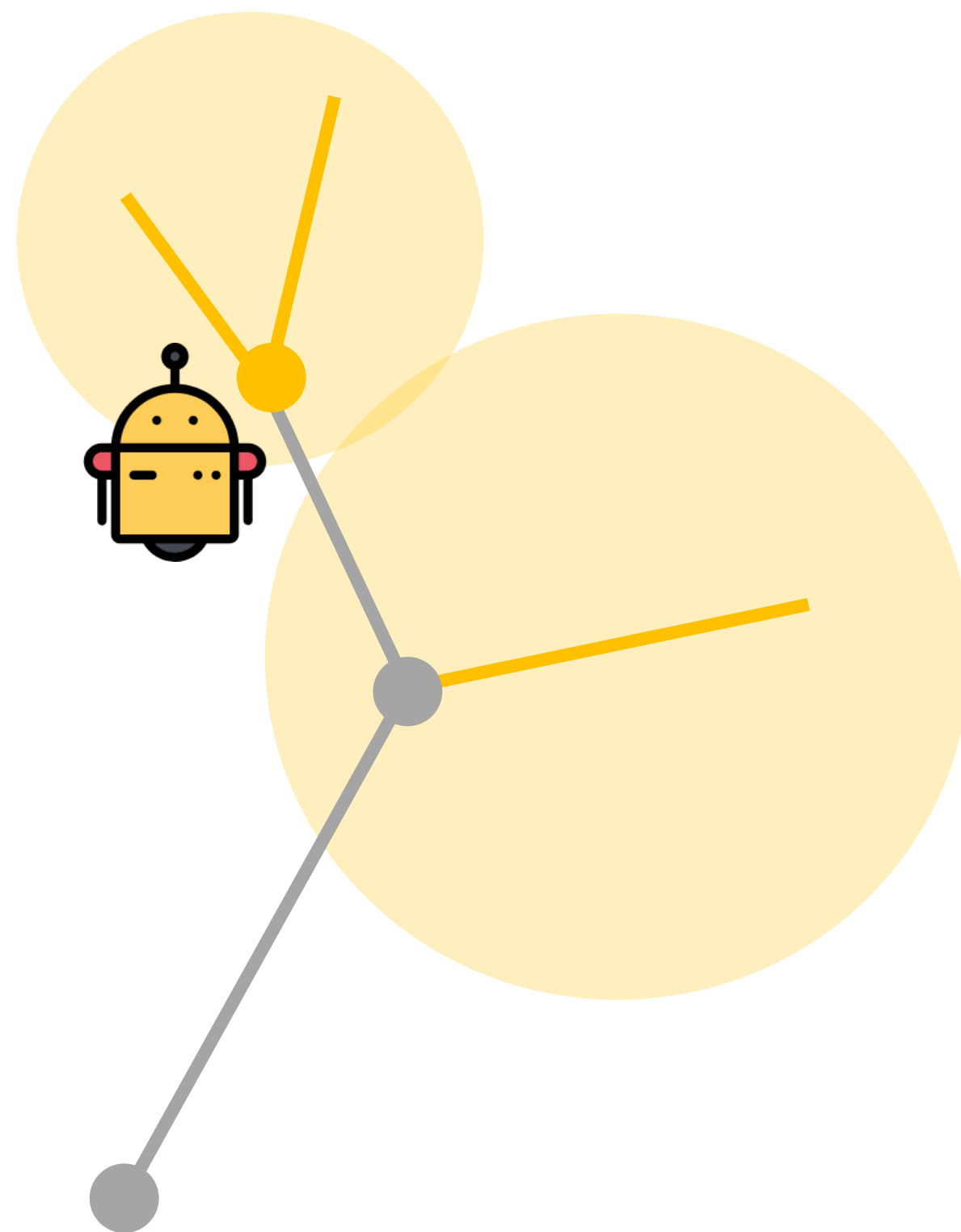
# What if you make a mistake?

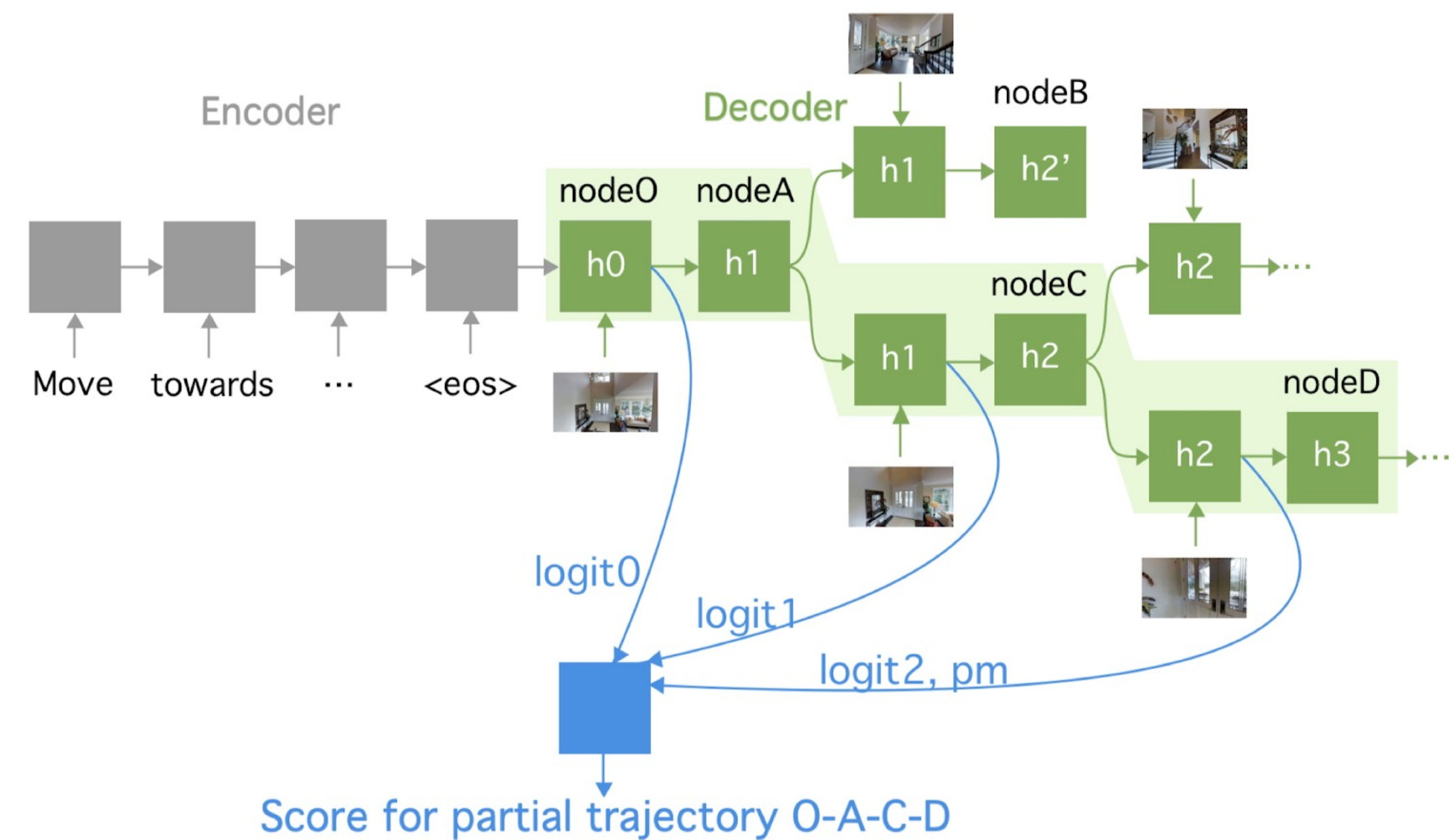**The New Frontier**

Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019

**Carnegie Mellon University** Language Technologies Institute

# What if you make a mistake?

**The Frontier**

Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019

# What if you make a mistake?

## Eventually ...



Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019

# What if you make a mistake?

??

**1. Did I reach the target?**

Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019

# What if you make a mistake?

1. **Did I reach the target?**
2. **Am I lost?**

????

Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019

Carnegie Mellon University Language Technologies Institute

# What if you make a mistake?

1. **Did I reach the target?**
2. **Am I lost?**
3. **Should I backtrack?**

Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019

# What if you make a mistake?

1. **Did I reach the target?**
2. **Am I lost?**
3. **Should I backtrack?**
4. **Where to backtrack to?**

Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019

# What if you make a mistake?

A lot of the visual observations and actions have no correspondence to the language



Encoder

Decoder    nodeB

nodeO  nodeA

Move  towards  ···  <eos>

logit0
logit1
logit2, pm

Score for partial trajectory O-A-C-D

Ke 2019, Tactical Rewind: Self-Correction via Backtracking in Vision-and-Language Navigation - CVPR 2019

**Carnegie Mellon University** Language Technologies Institute

# Underspecification

Leave the bedroom

LEFT

FWD



Does this actually need vision?     Yes

Does this understand plans?     Maybe?

# Why does this question matter?

Because in general, we can't supervise everything

Hey Siri, remind me to do my laundry

if(detergent)                    else

remind at home              remind to buy detergent when at store

Hey Siri-bot, do my laundry

Go to hamper…

# ALFRED

Action Learning From Realistic Environments and Directives

# Seven High-level Tasks

Paths are generated by planner



Pick & Place



Double Place



Stack



Examine



Heat



Cool



Rinse

# Data collection

Tuple

(Stack, Fork, Cup, CounterTop, Kitchen3)

Planner

(x,y,z) | is_fork(x) ^ is_cup(y) ^ on(x, y) ^ is_counter(z) ^ on(y, z)

Sample

Execute



Annotate

Language        Language        Language

Language        Language        Language        Language        Language        Language

# Example Language



Goal: "Put a clean bowl of water on the kitchen island"

Instructions:
"**Turn right and begin walking across the room, then hang a left and walk over to the far side of the kitchen island.** Pick up the dirty bowl that is closest to the bottle of wine on the kitchen island. Turn left and take a step forward, then turn left and walk up to the sink. Put the dirty bowl in the sink and turn on the water, after a couple seconds turn the water off and remove the now clean bowl filled with water. Turn around and take a step forward so you are facing the kitchen island. Put the clean bowl of water on the island on the left corner."

# Action Space

*Wash the cup*



**Put In**

**Toggle**

- Masks for object interaction
- Discrete actions (no torques)

"Place a heated apple slice on the large table"

create_slice(apple)          heat(apple_slice)    place(apple_slice, table)

collect(knife)  locate(apple)    slice(apple)

# End-to-End Models

*Turn around and move to the stove, then turn left to face the counter to the left of the stove. Pick up the sharp knife with the yellow handle from the counter…*



LSTM

PickupObject

# Action Spaces

**Choose a view**

**Outline an Object**

**Grasp an Object**

`PickupObject`

# Pick-up

## What's hidden in that?

If I gave you one of these and labeled it, could you abstract to the others?

Does "pick up" mean the same thing for all of these?



Does "pick up" correspond to a specific action sequence?



Mousavian et al. 6-DOF GraspNet: Variational Grasp Generation for Object Manipulation — ICCV 2019

# Simplify with Blocks and Coordinates

*Put the orange block to the right of the green block*



Why?

Is this a useful training datum?

("Put the orange block to right of the green block",
0.35)

We no longer have a discrete grounding
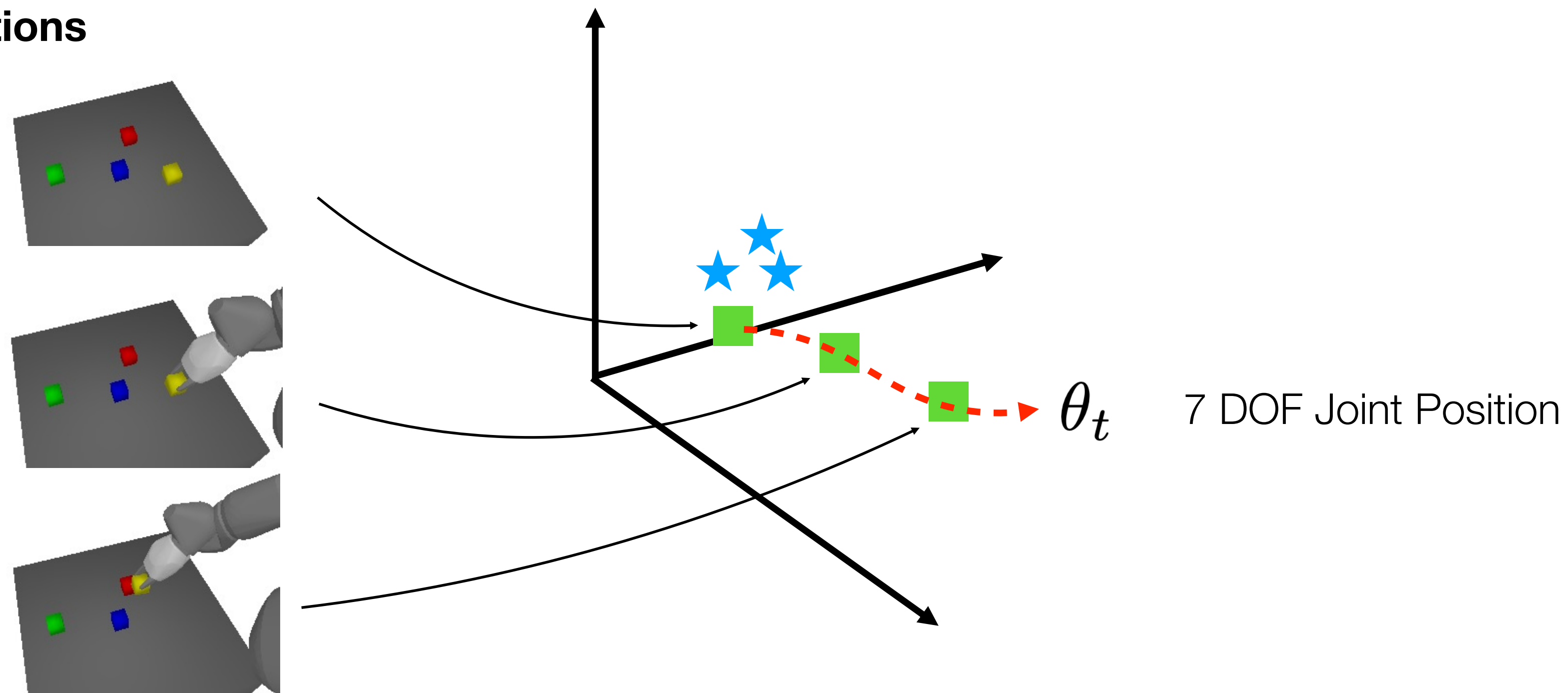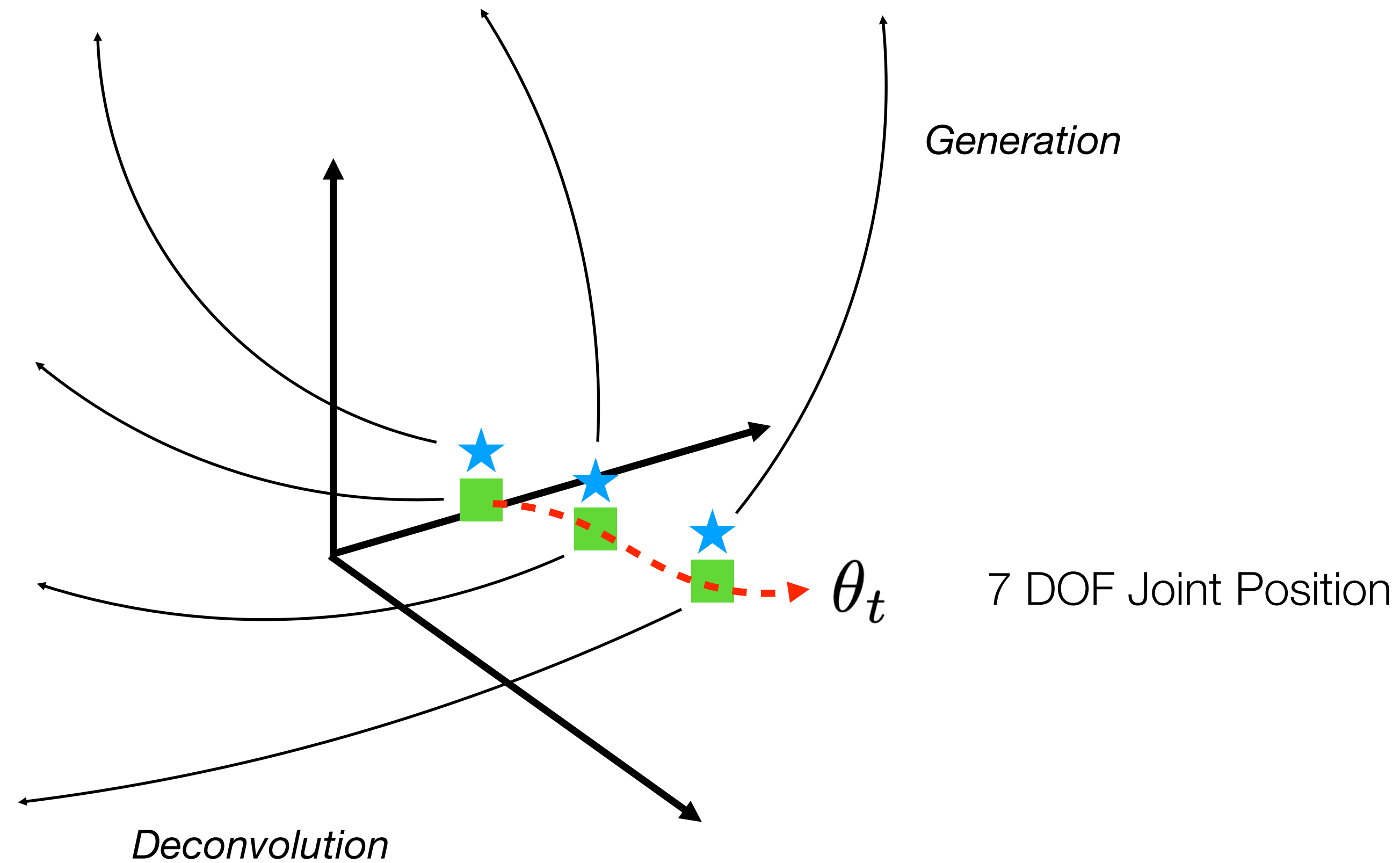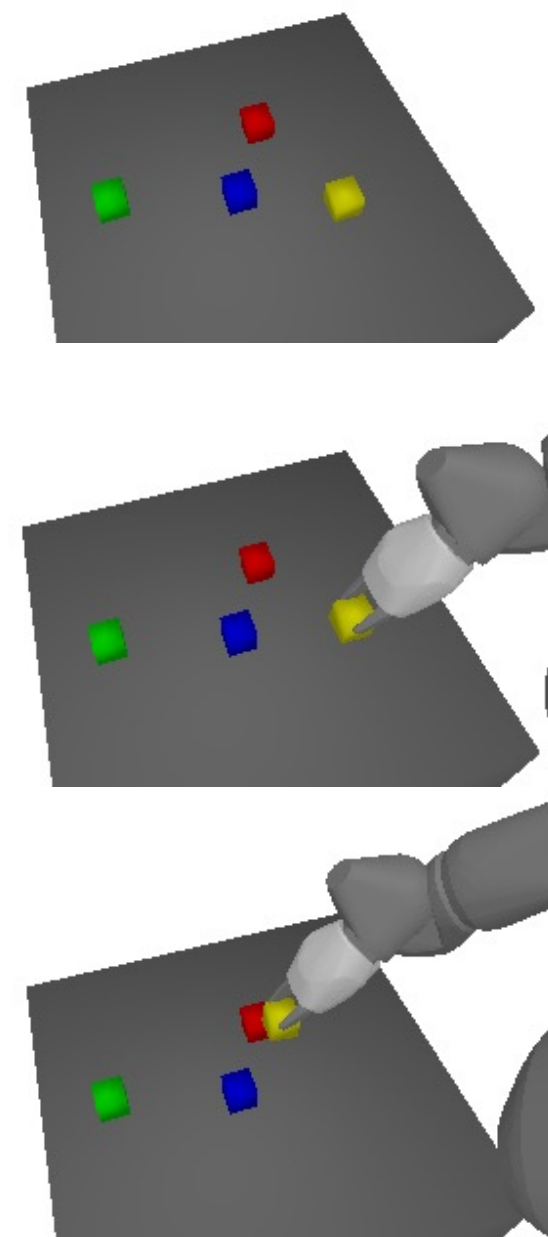
# Simple Blocks

# A Shared Semantic Space

**Language**

*"take the yellow object from the table and place it on top of the red object"*

`move_to(yellow)  grasp(yellow)  … release(yellow)`

**Observations**

Paxton et al. Prospection: Interpretable Plans From Language By Predicting the Future ICRA 2019

# A Shared Semantic Space

**Language**

*"take the yellow object from the table and place it on top of the red object"*

`move_to(yellow)  grasp(yellow)  … release(yellow)`
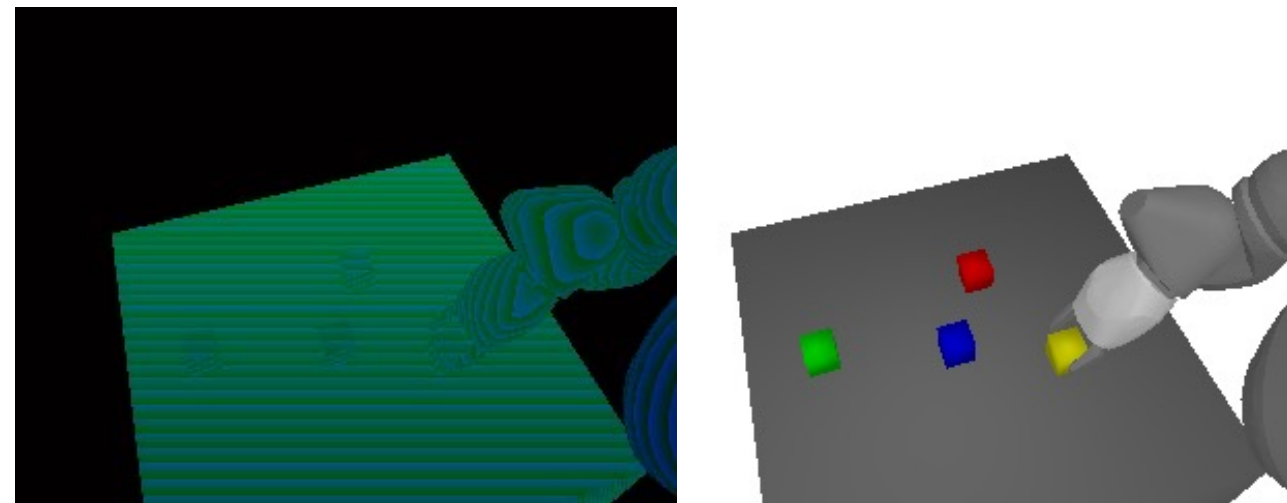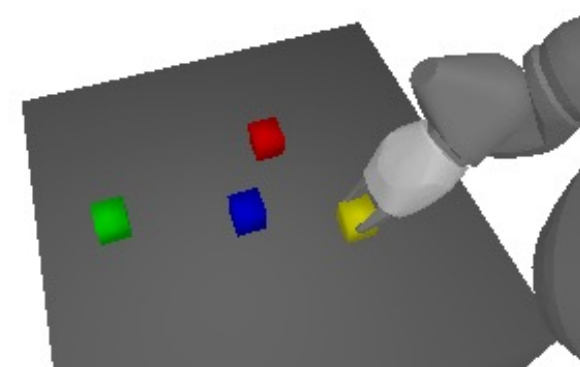
**Observations**



$\theta_t$    7 DOF Joint Position

# A Shared Semantic Space

**Language**

*"take the yellow object from the table and place it on top of the red object"*

```
move_to(yellow)  grasp(yellow)  … release(yellow)
```

*Generation*

**Observations**

$\theta_t$    7 DOF Joint Position

*Deconvolution*

Paxton et al. Prospection: Interpretable Plans From Language By Predicting the Future ICRA 2019

# Predicting the Future

**Goal:**

*take the yellow object from the table and place it on top of the red object*
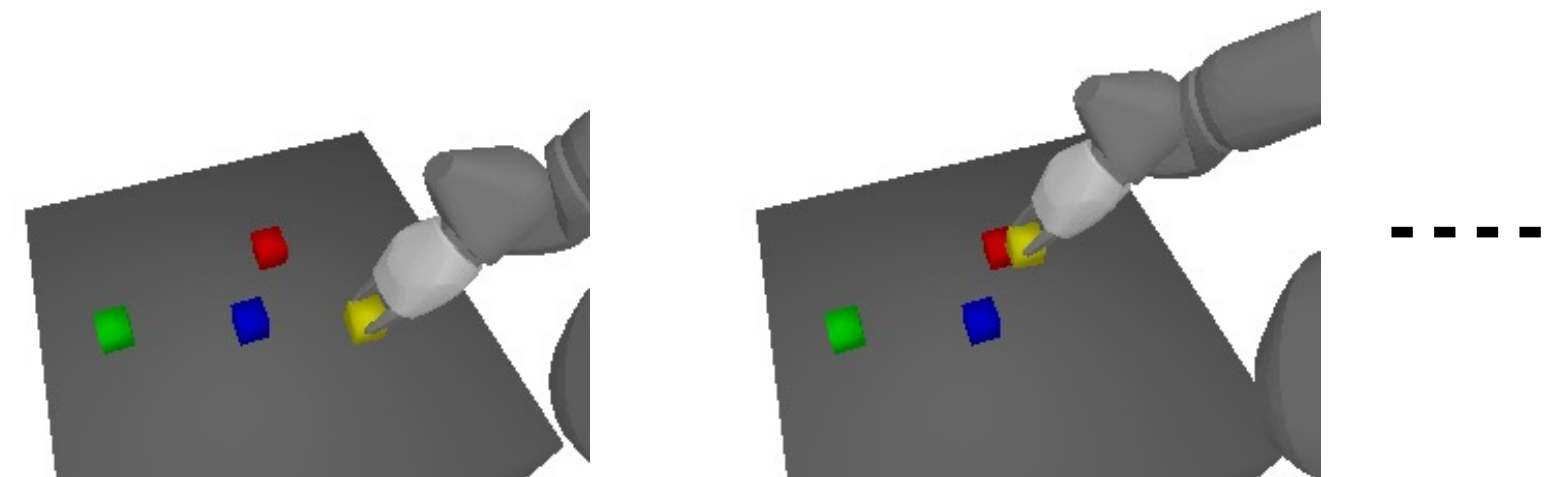
**Current World**



$h_t \longrightarrow$ `grasp(yellow)`



**Interpretable Possible Futures**

`lift(yellow)`   `move(yellow, red)`

# Objectives

Latent Space   $Z_t$

| Reconstruction | Pose | SubGoal | Block pos |
|---|---|---|---|
| $||\hat{W}_t - W_t||_2^2$ | $C_{actor}(\hat{\theta}_t, \theta_t)$ | $C_G(\hat{G}_t, G_t)$ | $C_{obj}(z_t)$ |



predicted
current

```
move(yellow,
    red)
```

x #steps in horizon

**Carnegie Mellon University** Language Technologies Institute

Paxton et al. Prospection: Interpretable Plans From Language By Predicting the Future ICRA 2019

# Long Tails



**Templates:**

   put the yellow one on the green block

**Humans:**

   move the yellow cube to the right until it is on top of the green cube with the front
   half  of  the  yellow  cube  touching the far half of the top of the green cube

# Simple UNet Sim2Real



Sagar Gubbi, Anirban Biswas, Raviteja Upadrashta, Vikram Srinivasan,
Partha Talukdar, Bharadwaj Amrutur

INDIAN INSTITUTE OF SCIENCE

Carnegie Mellon University Language Technologies Institute

https://arxiv.org/abs/2012.13693

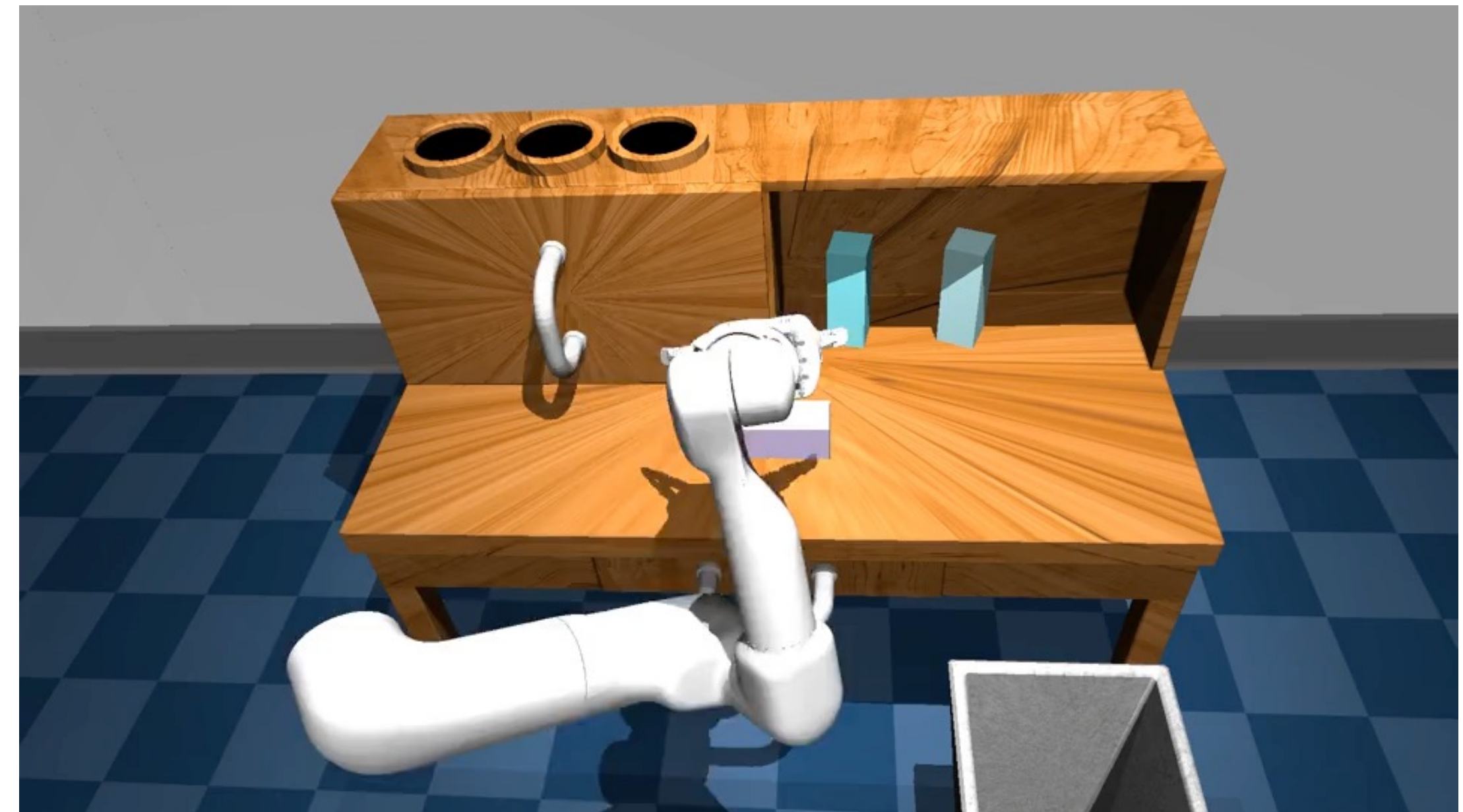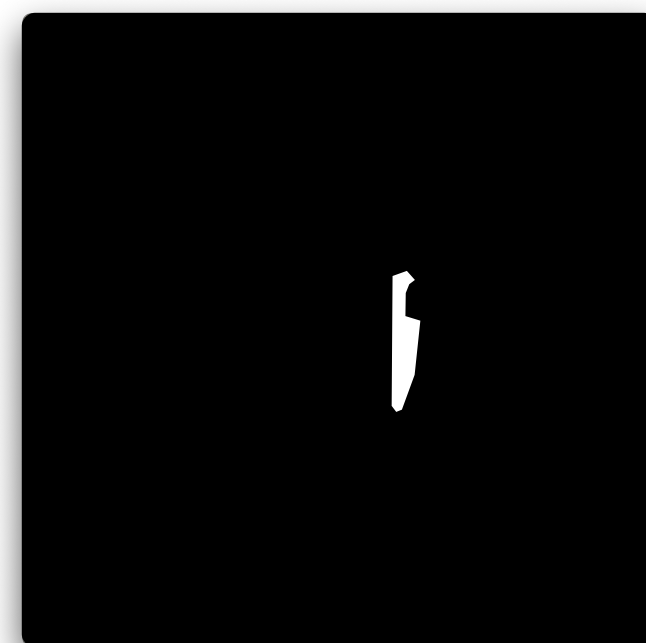# Where does semantics come from?

Someone labeled it?

$$p(a|v_0, ..., v_t)$$

Self-Play and Physical Affordances?
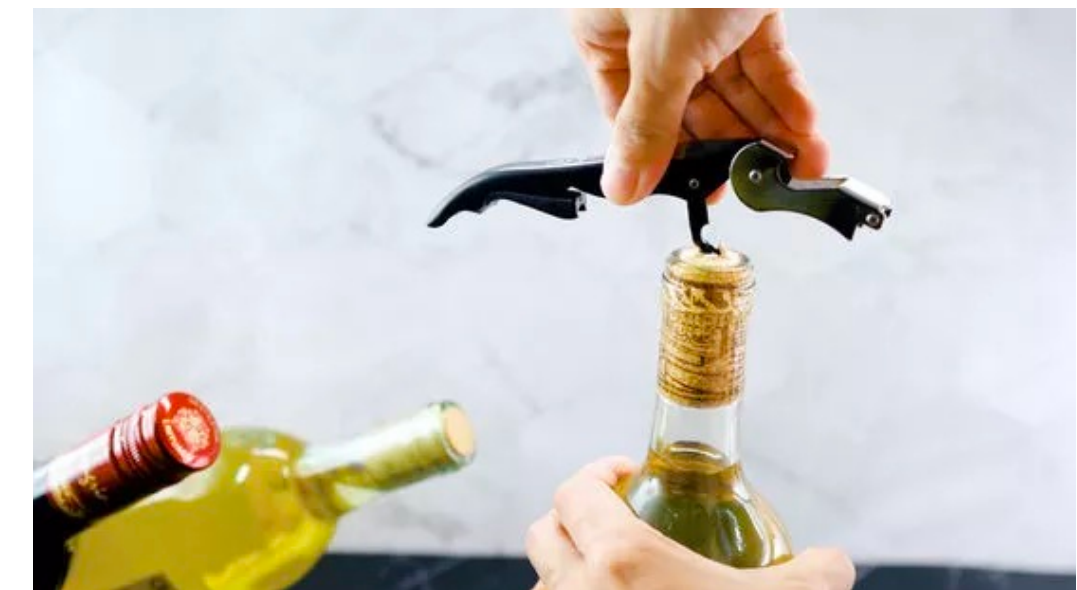
Simulator Definitions?



PickupObject



Lynch et al. — Learning Latent Plans from Play — CoRL 2019

# Embodiment

- Choose your own adventure — Lots of noise

- What does it mean to succeed?

- Where do concepts come from?

- What's the role of exploration?

- Language is woefully underspecified

- + Everything that makes vision and robotics hard

**Carnegie Mellon University** Language Technologies Institute

# Practical Comments — Don't try and solve it all

- How much error is due to underspecification / *TASK* planning failures?
  - Prediction?
  - Tracking?
- How much error is due to *CONTROL* planning failures?
  - Kinematics?
  - Grounding?
- How much is due to novel scenarios?
  - Unseen environments/worlds
  - New Language?
  - Novel task composition?