

# Linguistic scaffolds for policy learning

---

Jacob Andreas

Berkeley → Microsoft Semantic Machines → MIT

# Linguistic scaffolds for policy learning

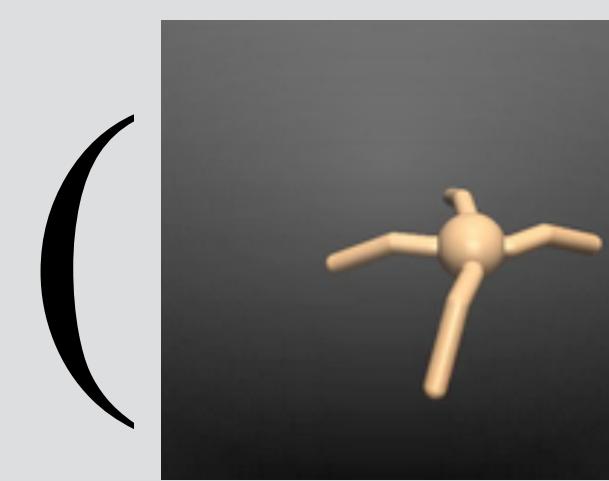
## (what can language do for RL?)

---

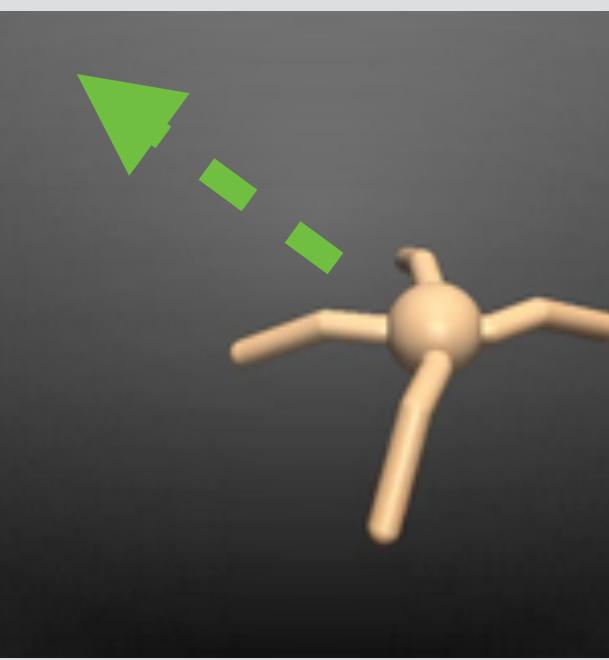
Jacob Andreas

Berkeley → Microsoft Semantic Machines → MIT

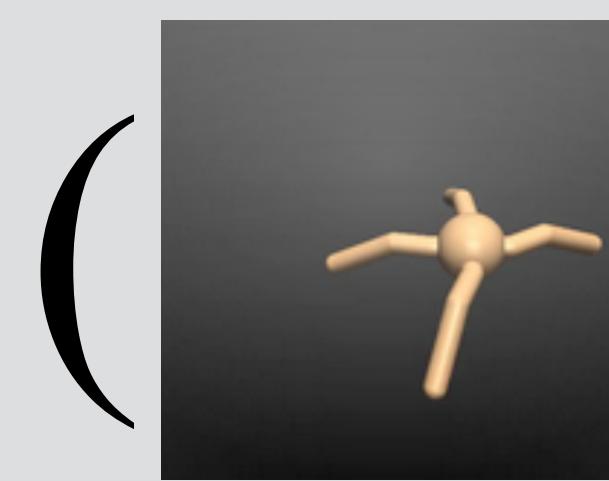
# An NLPer's view of RL



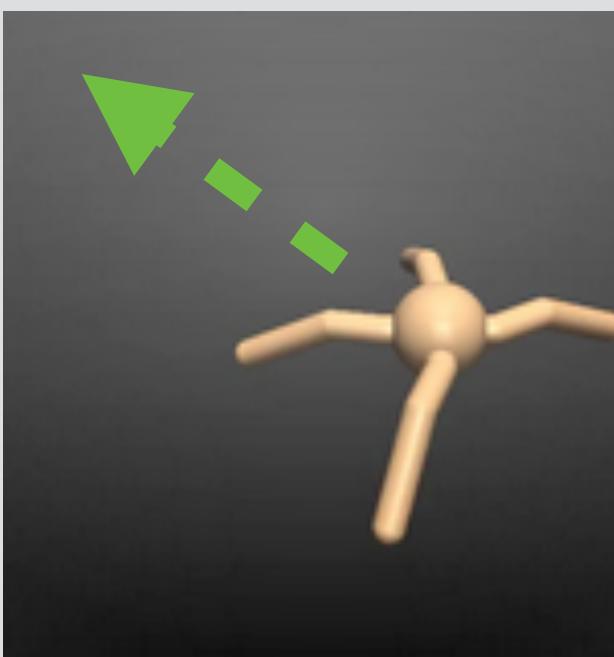
$(\text{ }$ ,  $R)$   $\rightarrow$



# An NLPer's view of RL

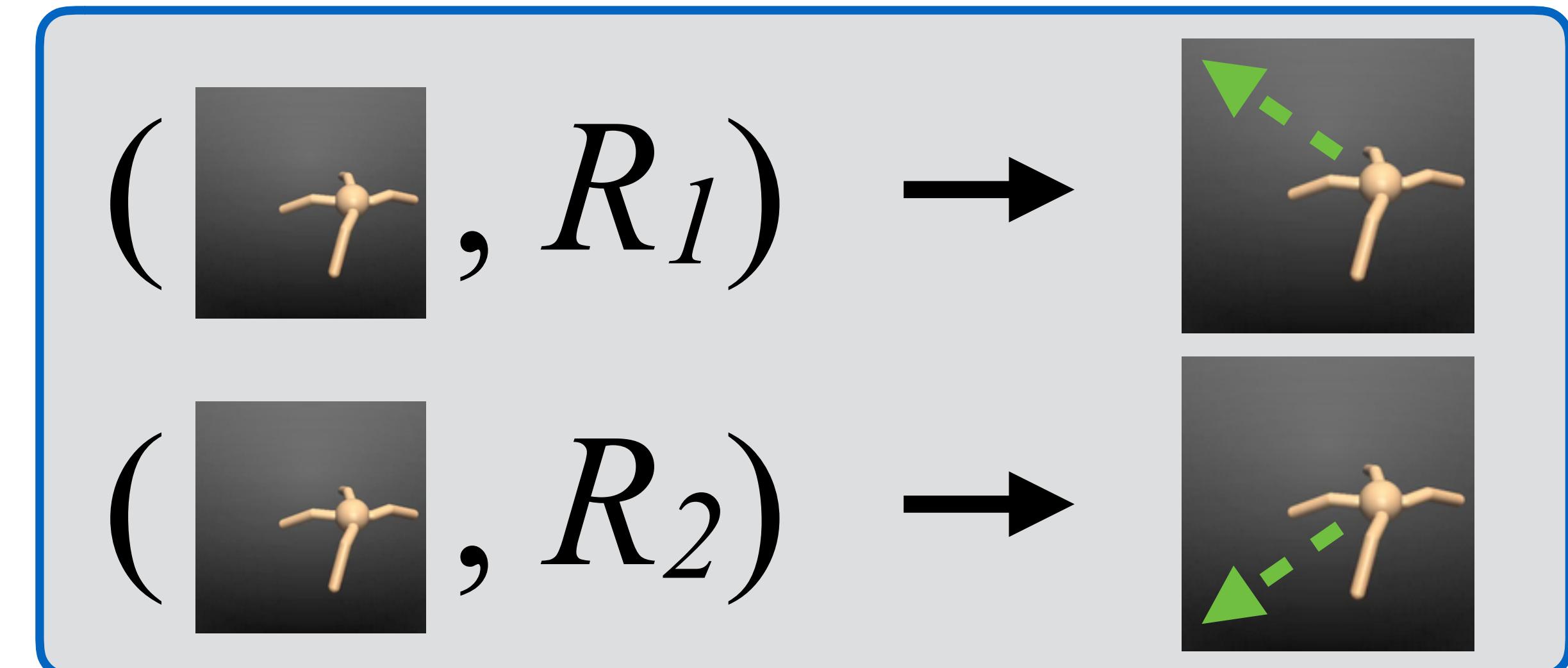
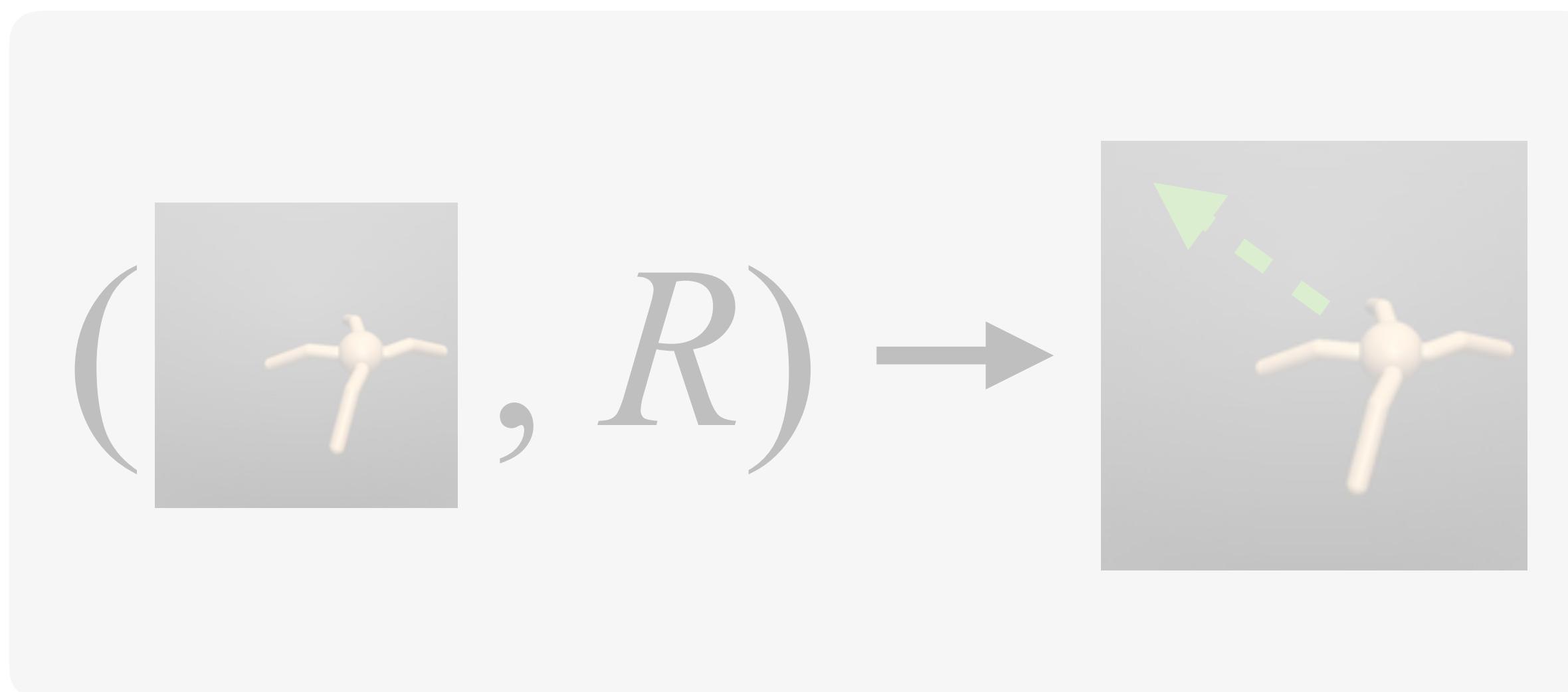


$(\text{ }, R)$   $\rightarrow$



memorize 1 reward fn

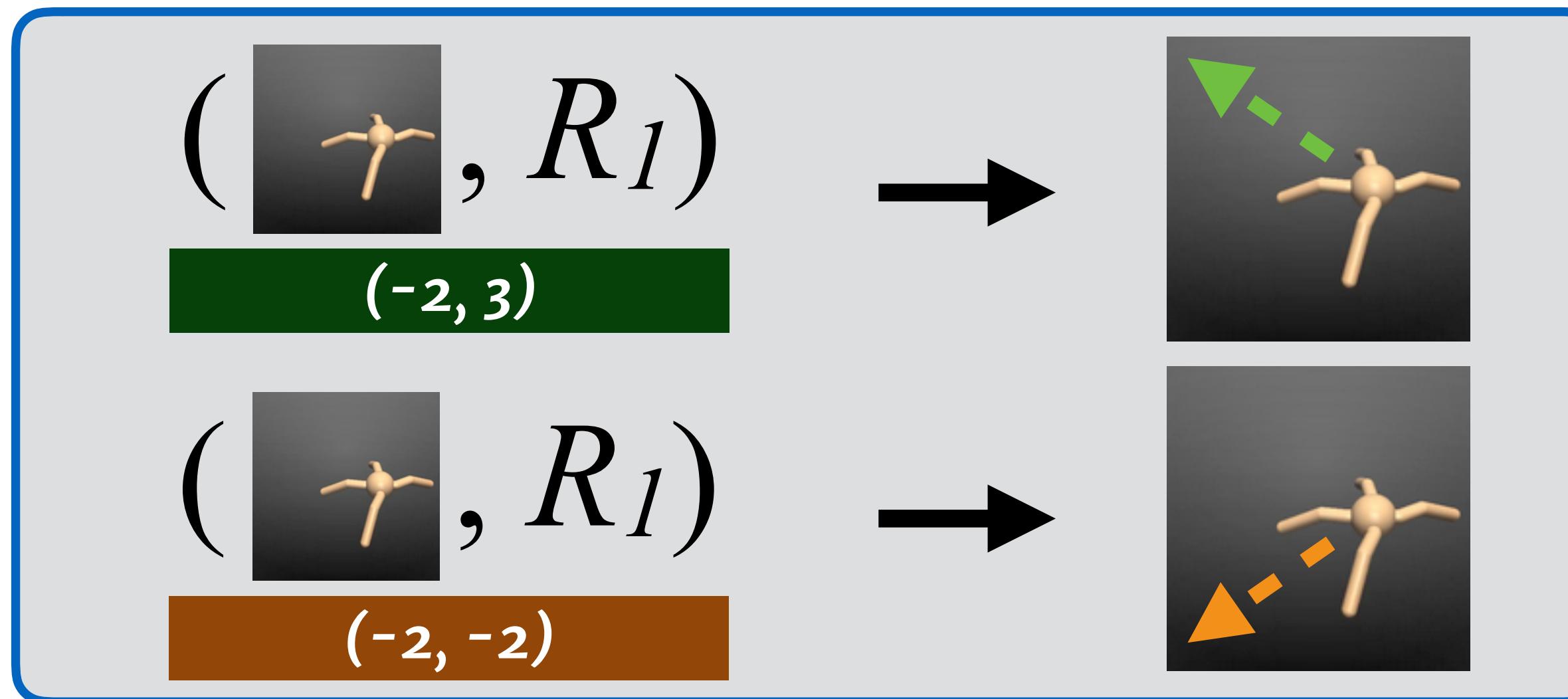
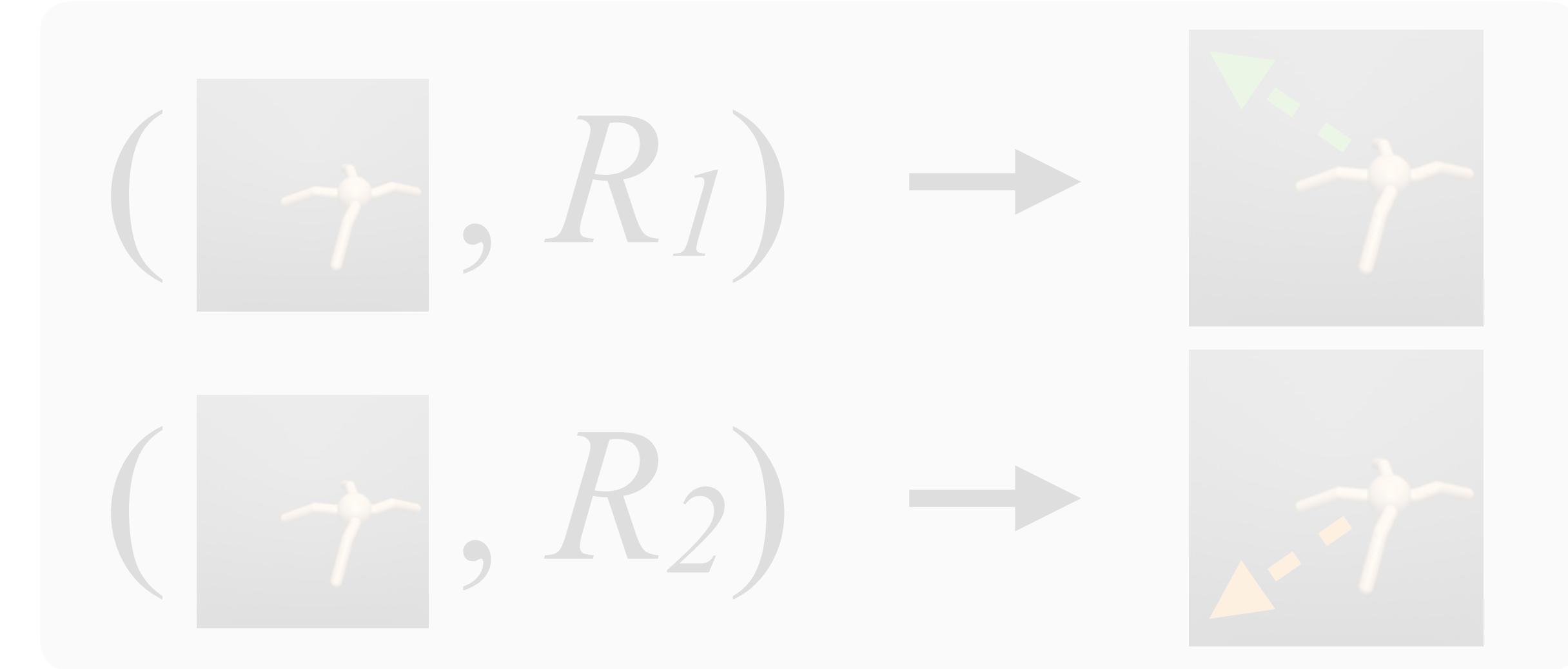
# An NLPer's view of RL



memorize k reward fns

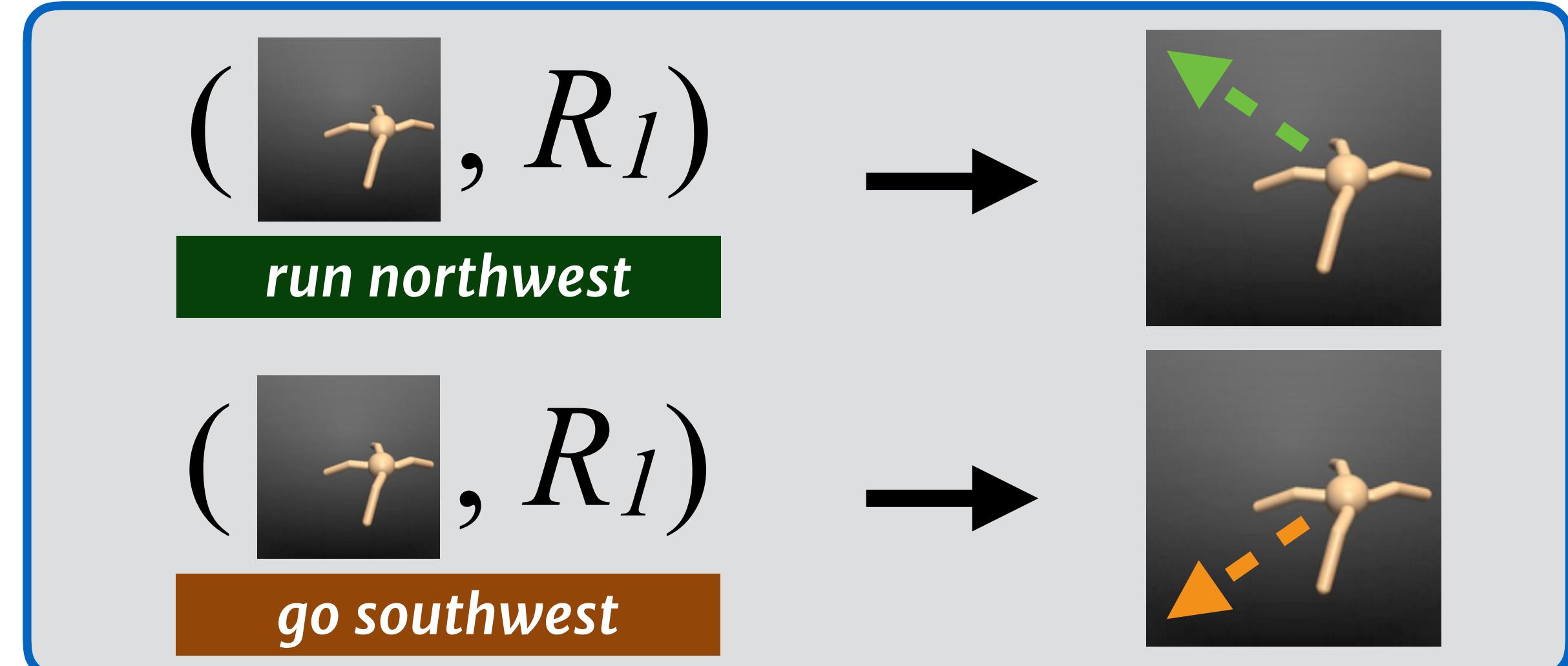
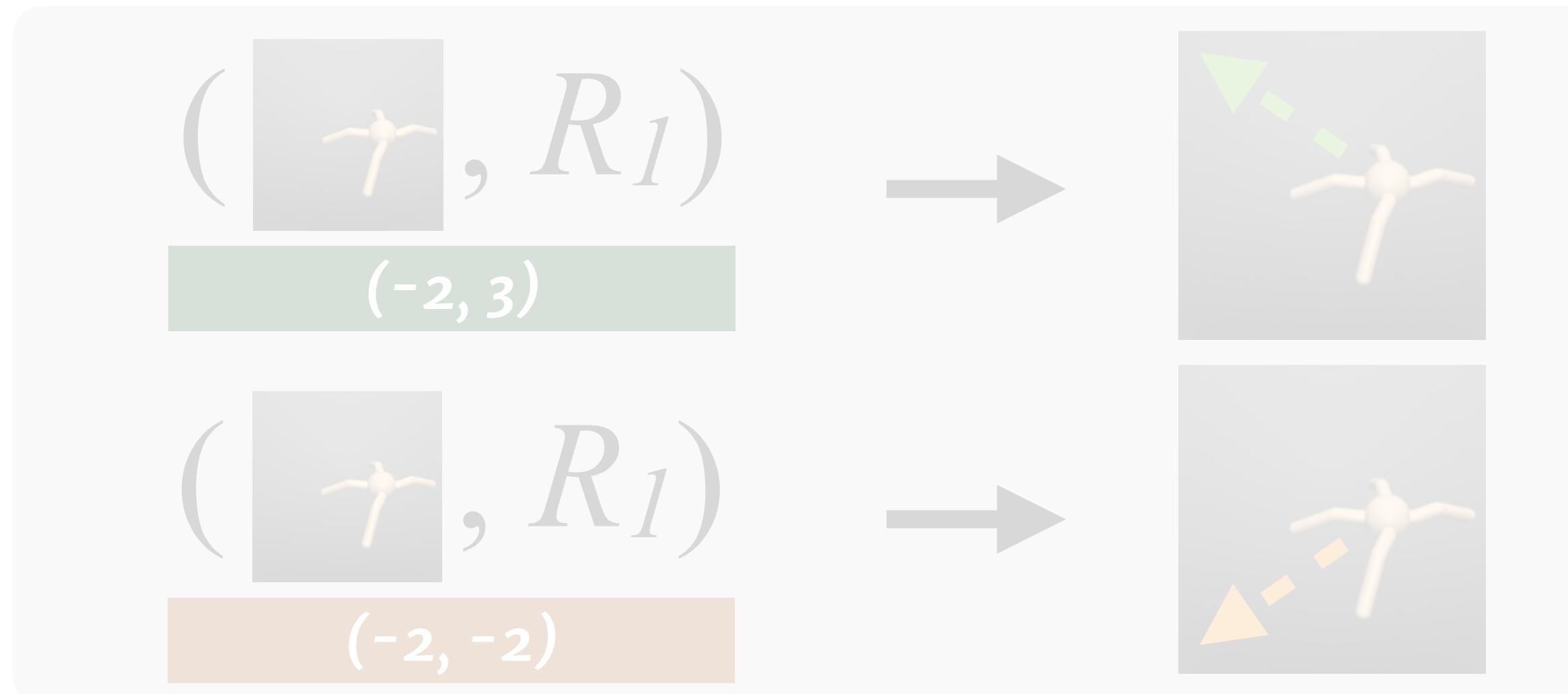
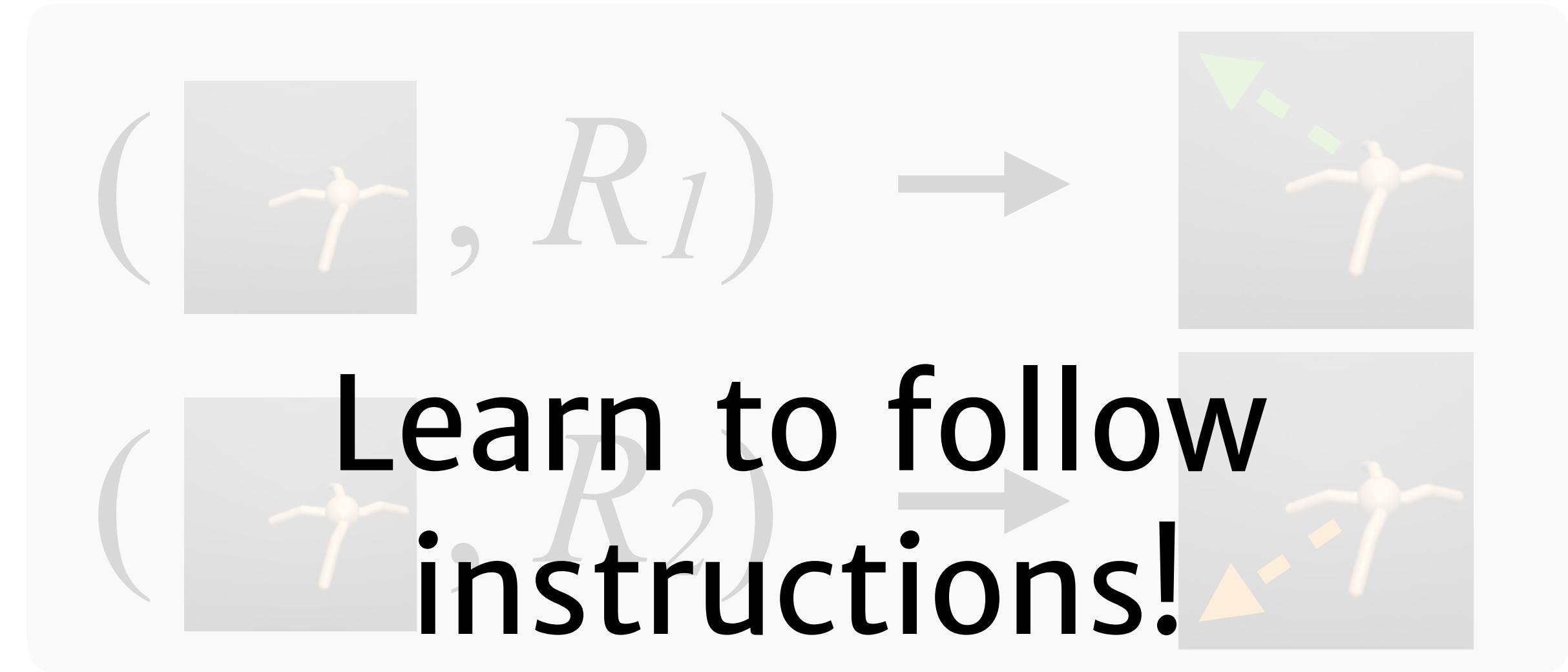
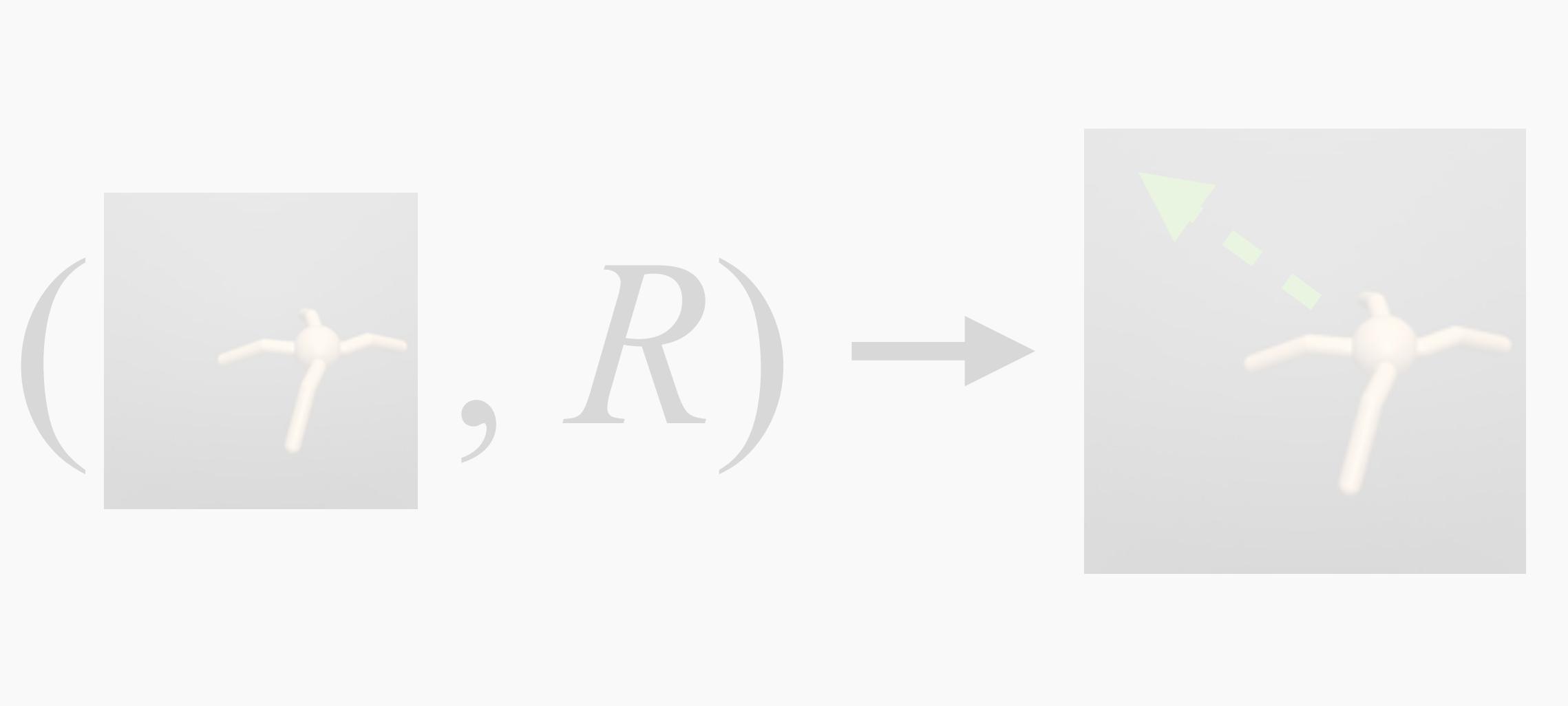
[e.g. Taylor & Stone 09]

# An NLPer's view of RL

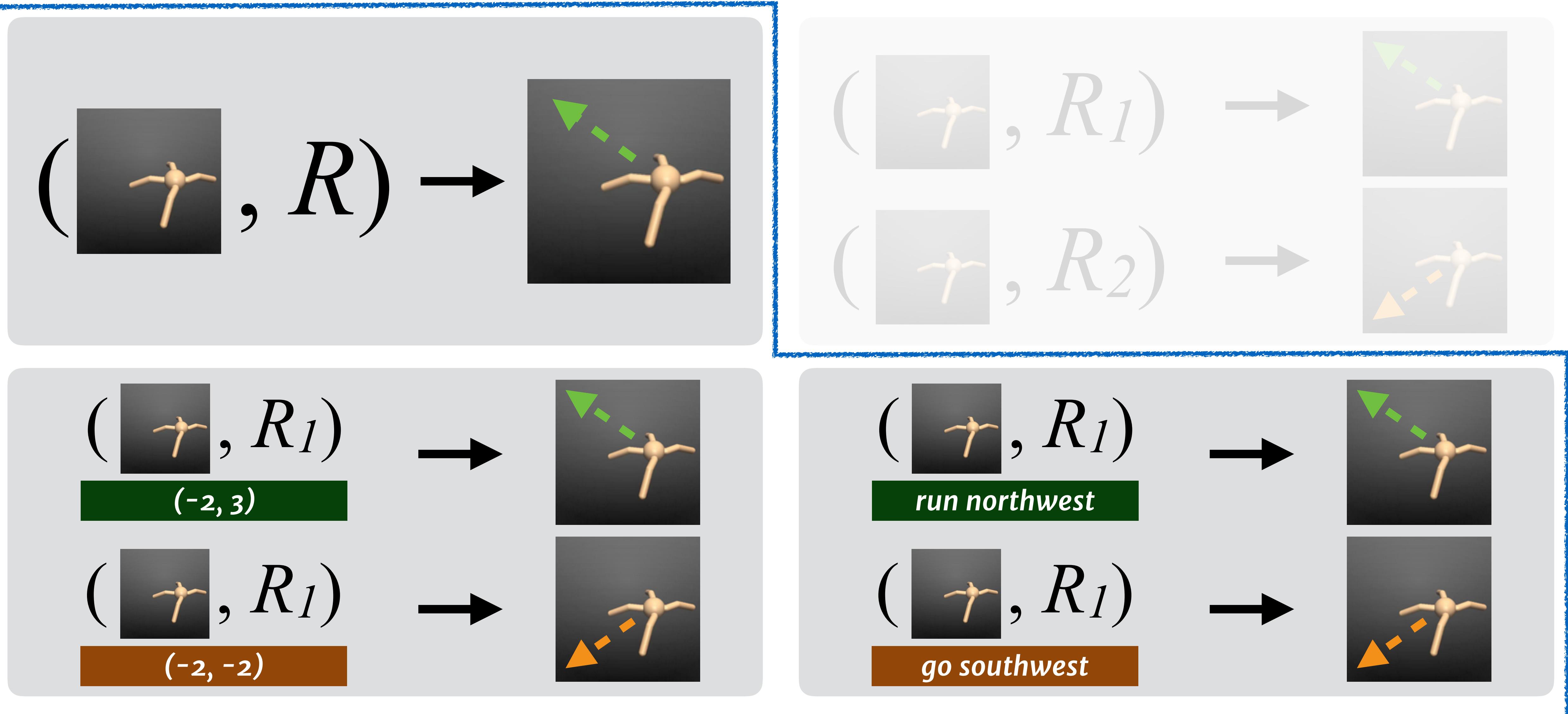


[e.g. Schaul et al. 15]

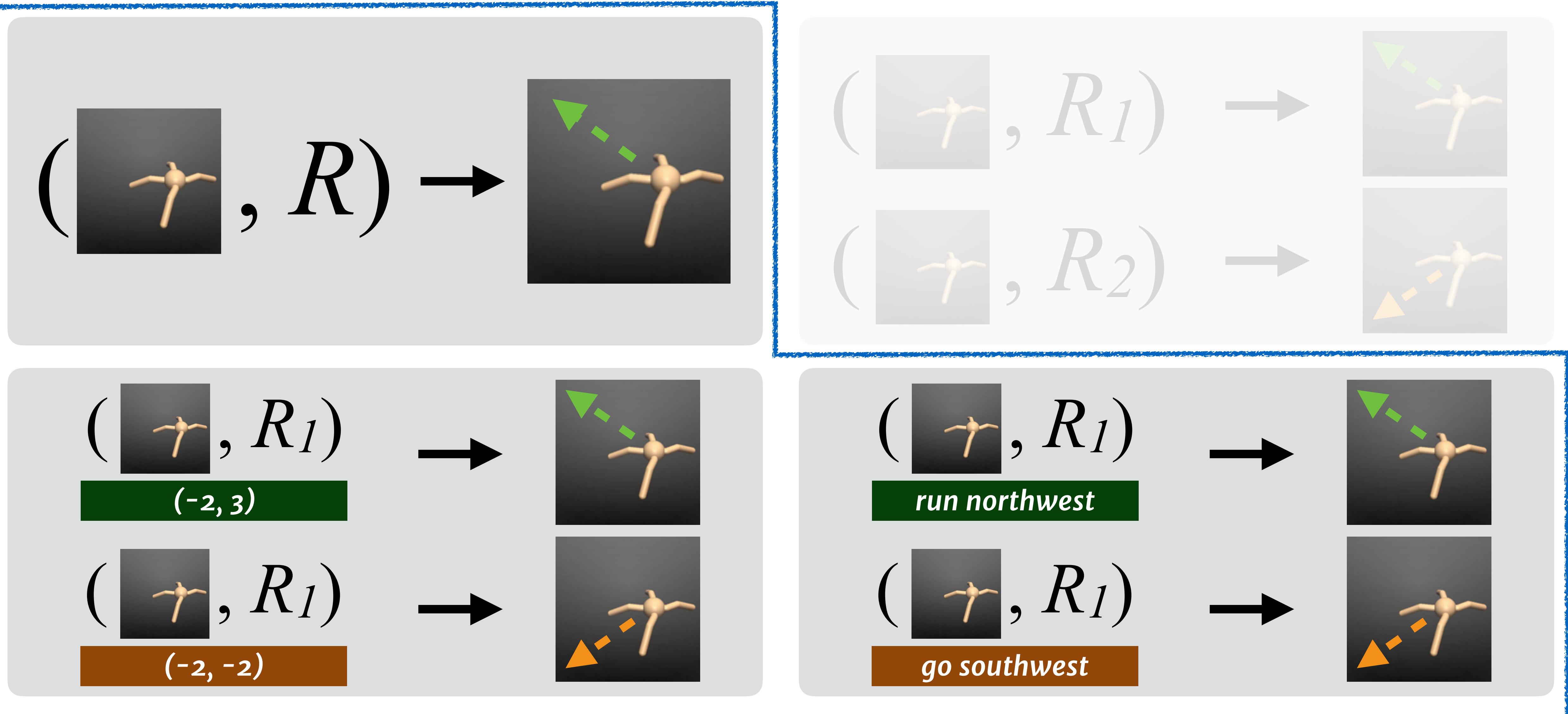
# An NLPer's view of RL



# Instructions as observations

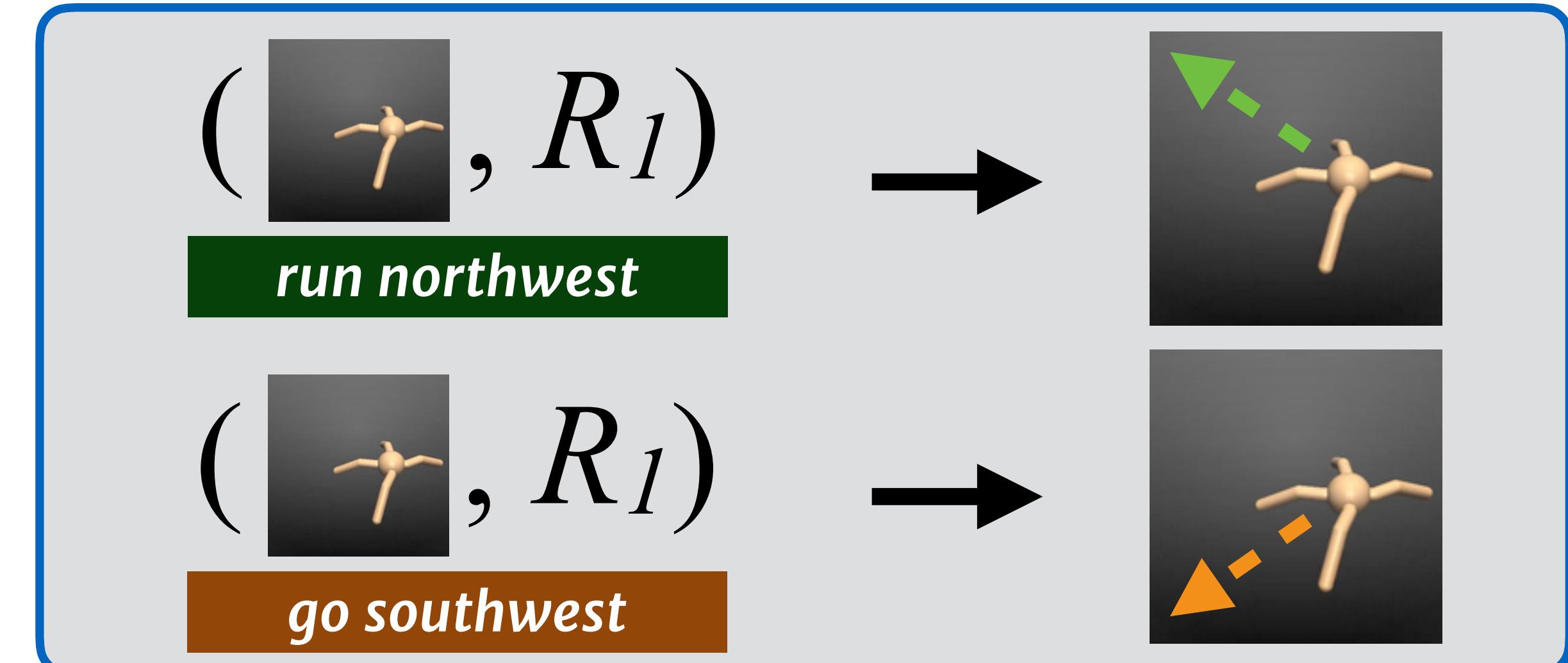


# Instructions as observations



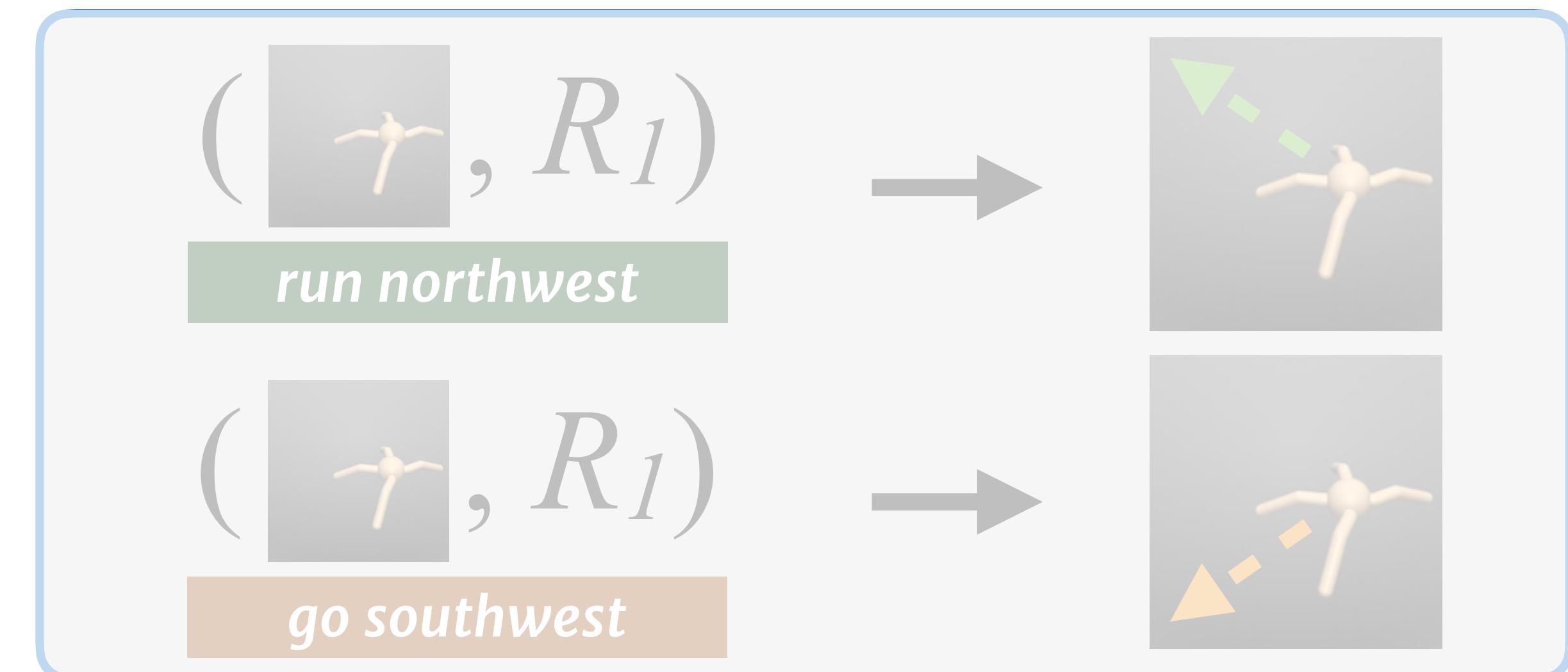
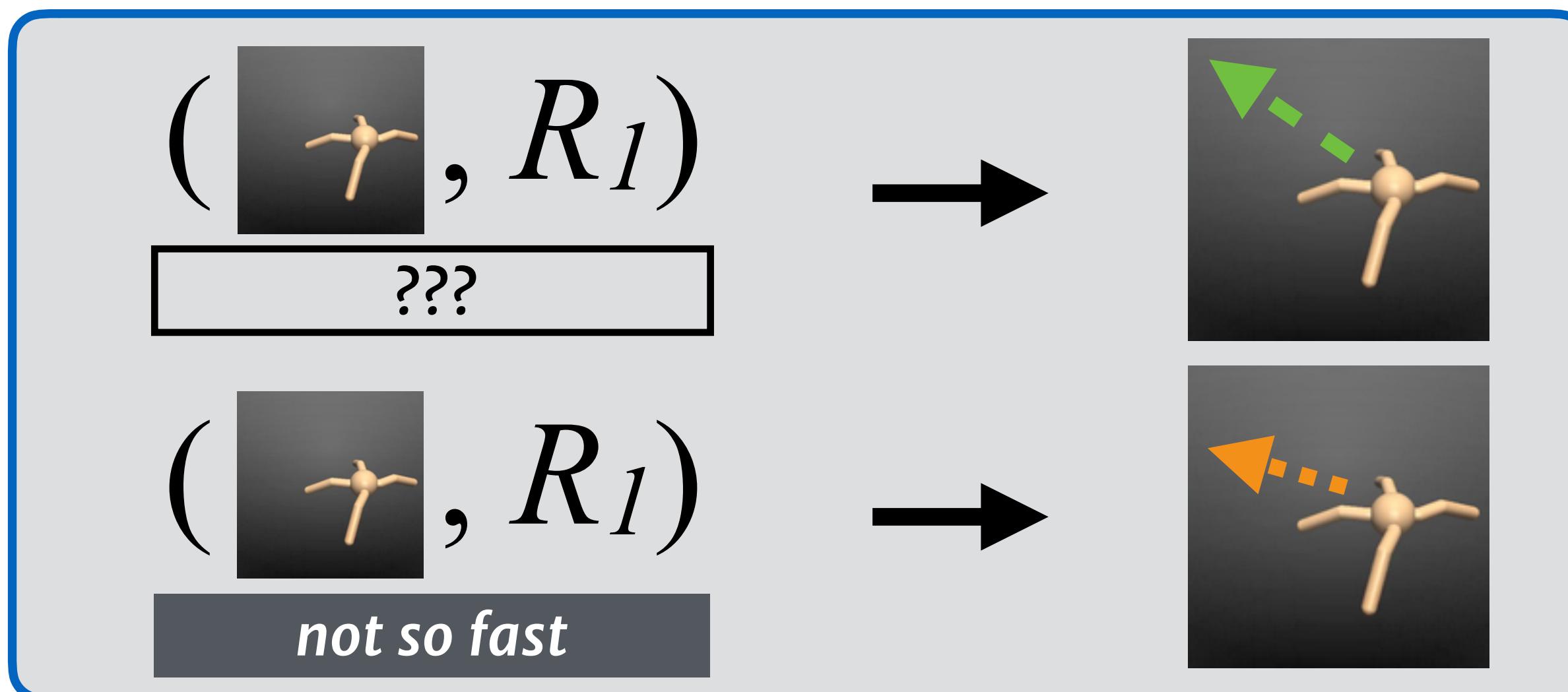
# Beyond observations

(1) Instructions are moves in a game, not observations of an environment.



# Beyond goals

(2) There's more to language learning than instruction following!



Language use as gameplay

# Generation & understanding

---

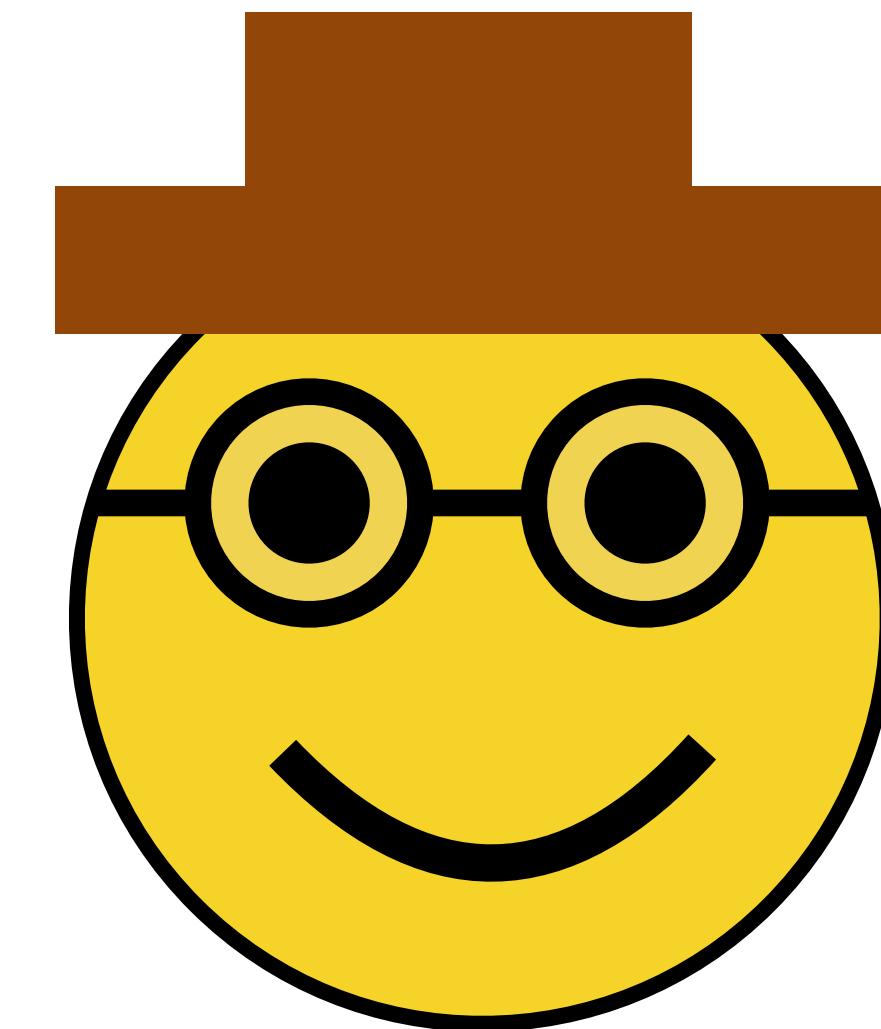


*Turn right and walk through the kitchen. Go right into the living room and stop by the rug.*

[Anderson et al. 18]

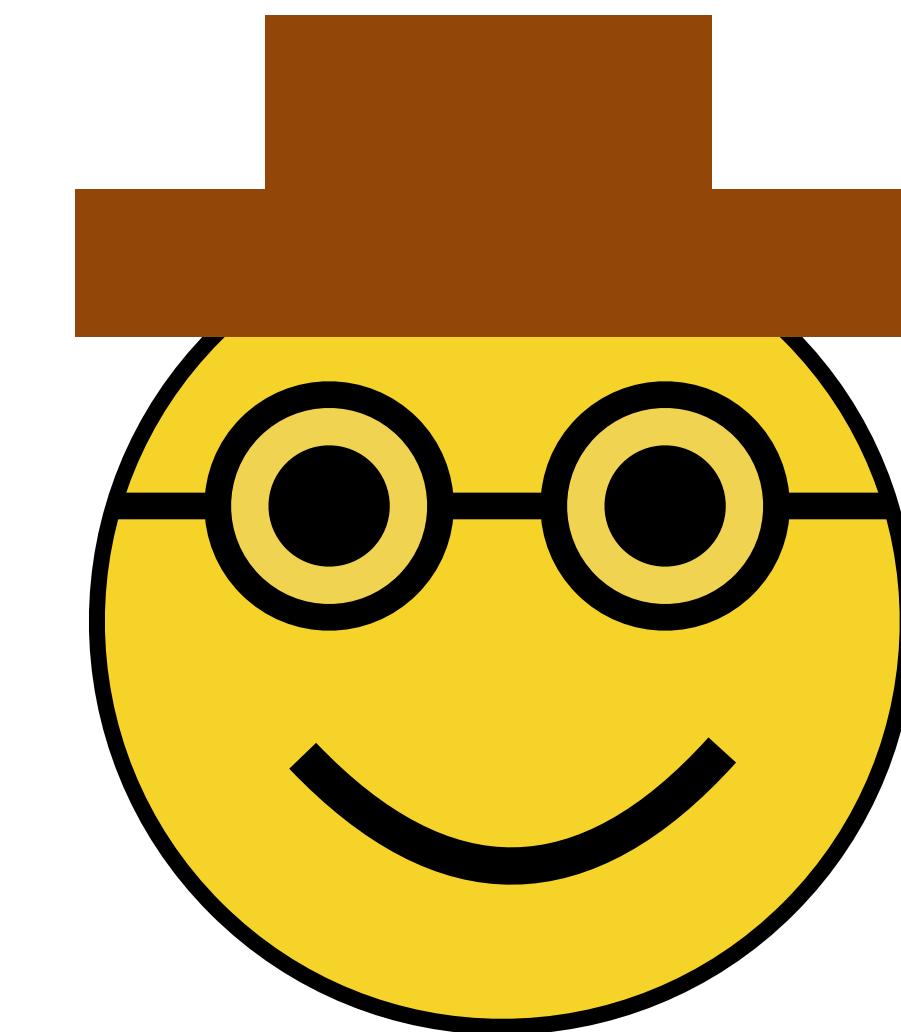
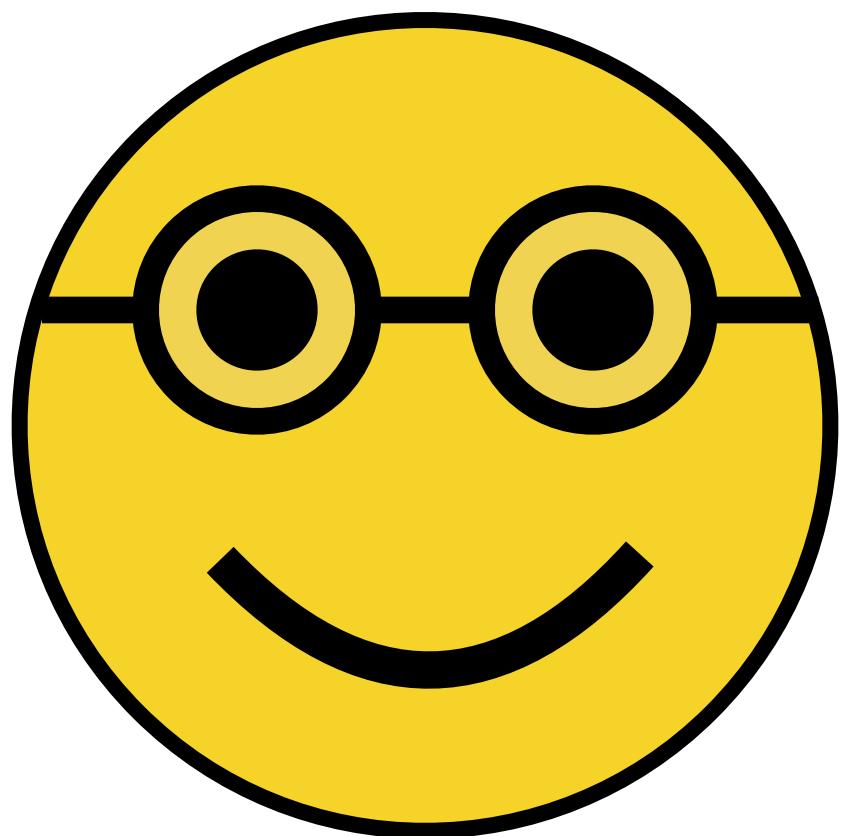
# A reference game

---



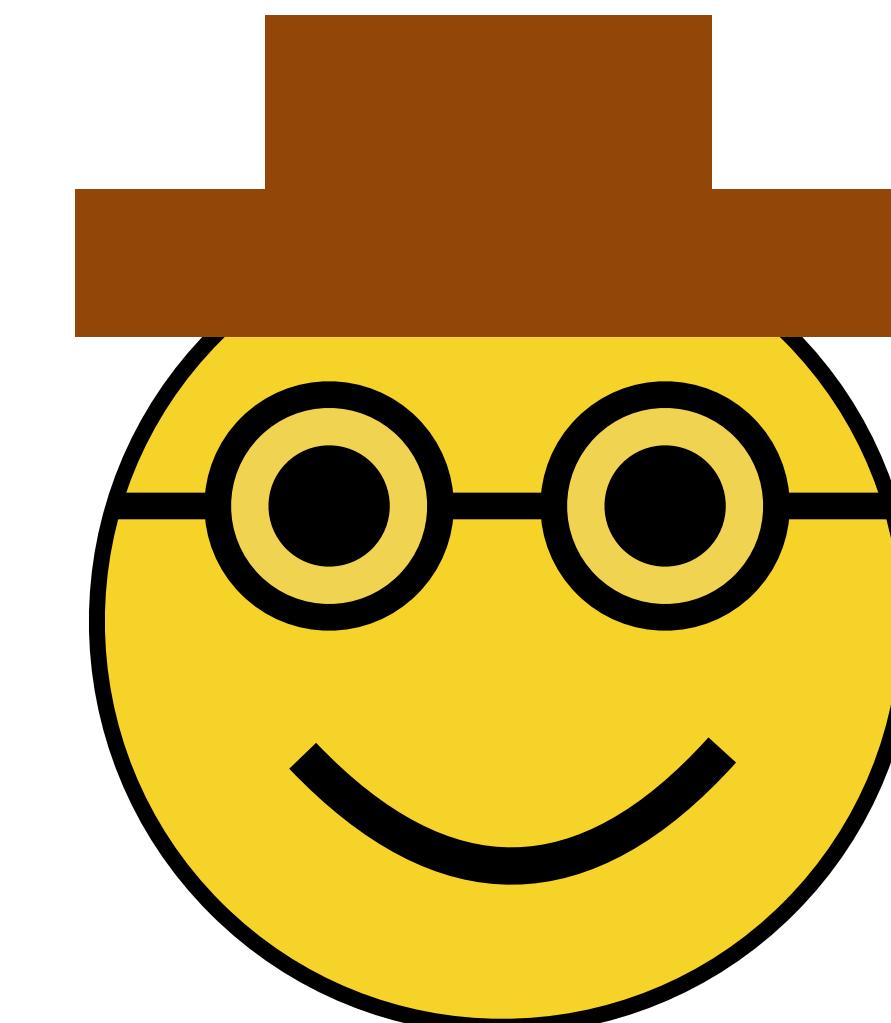
*“glasses”*

---



*“glasses”*

---



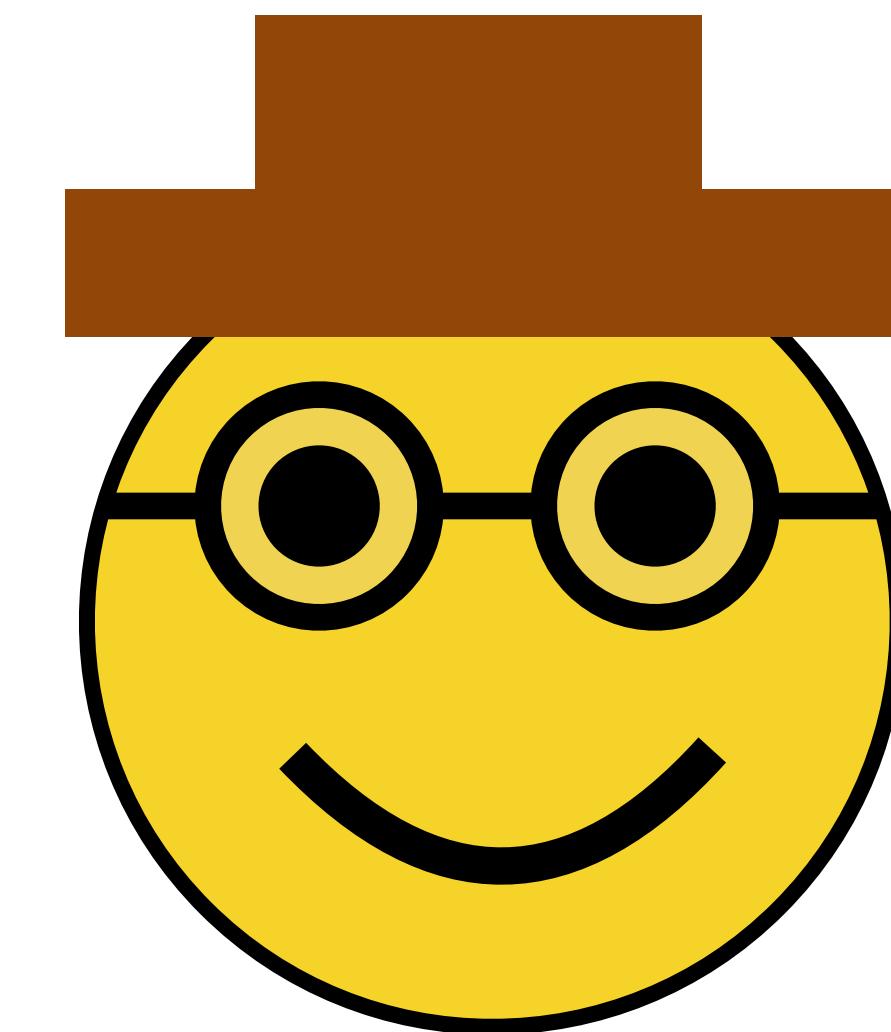
*“glasses”*

---



*“glasses”*

---



# The rational speech acts model

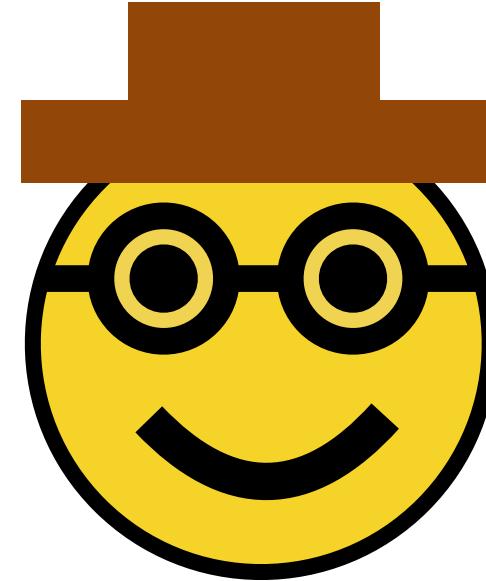
---

$L_o(. \mid glasses)$



1/2

$L_o(. \mid hat)$



0

1/2

1

# The rational speech acts model

---

$L_o(\cdot \mid \text{glasses})$	$1/2$	$1/2$
$L_o(\cdot \mid \text{hat})$	$0$	$1$
$S_1(\text{glasses} \mid \cdot) \propto L_o(\cdot \mid \text{glasses})$	$1$	$1/3$
$S_1(\text{hat} \mid \cdot)$	$0$	$2/3$

# The rational speech acts model

---

$$L_1(\cdot | \text{glasses}) \propto S_1(\text{glasses} | \cdot)$$

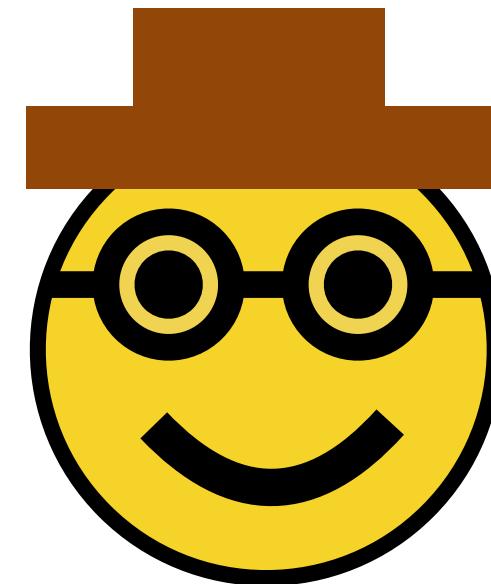
3/4

$$L_1(\cdot | \text{hat})$$

0

1/4

1



$$S_1(\text{glasses} | \cdot) \propto L_0(\cdot | \text{glasses})$$

1

1/3

$$S_1(\text{hat} | \cdot)$$

0

2/3

# Pragmatics

---

*Q: Do you know what time it is?*

# Pragmatics

---

*Q: Do you know what time it is?*

*A: Yes*

# Pragmatics

---

*Q: Do you know what time it is?*

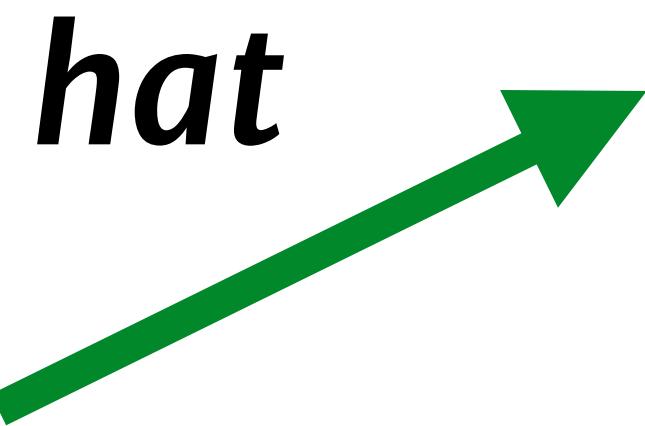
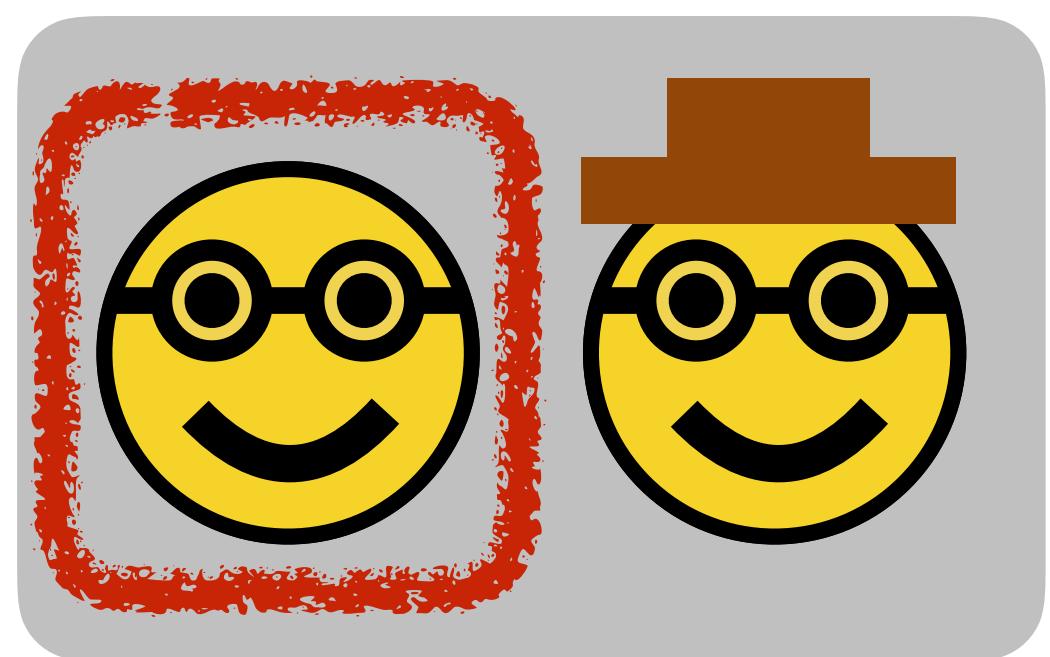
*A: Yes*

*I find his cooking very interesting.*

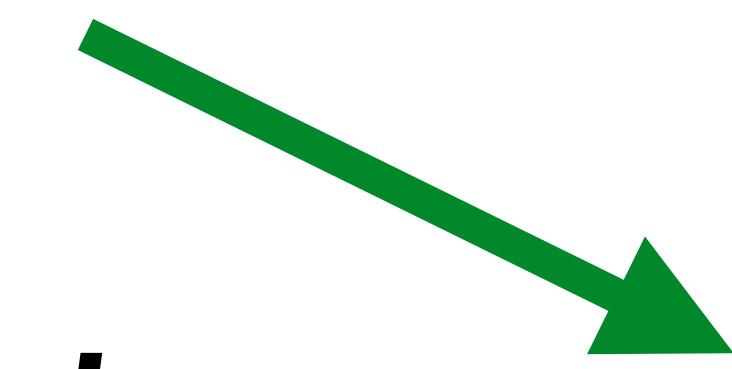
# RSA game tree

---

speaker

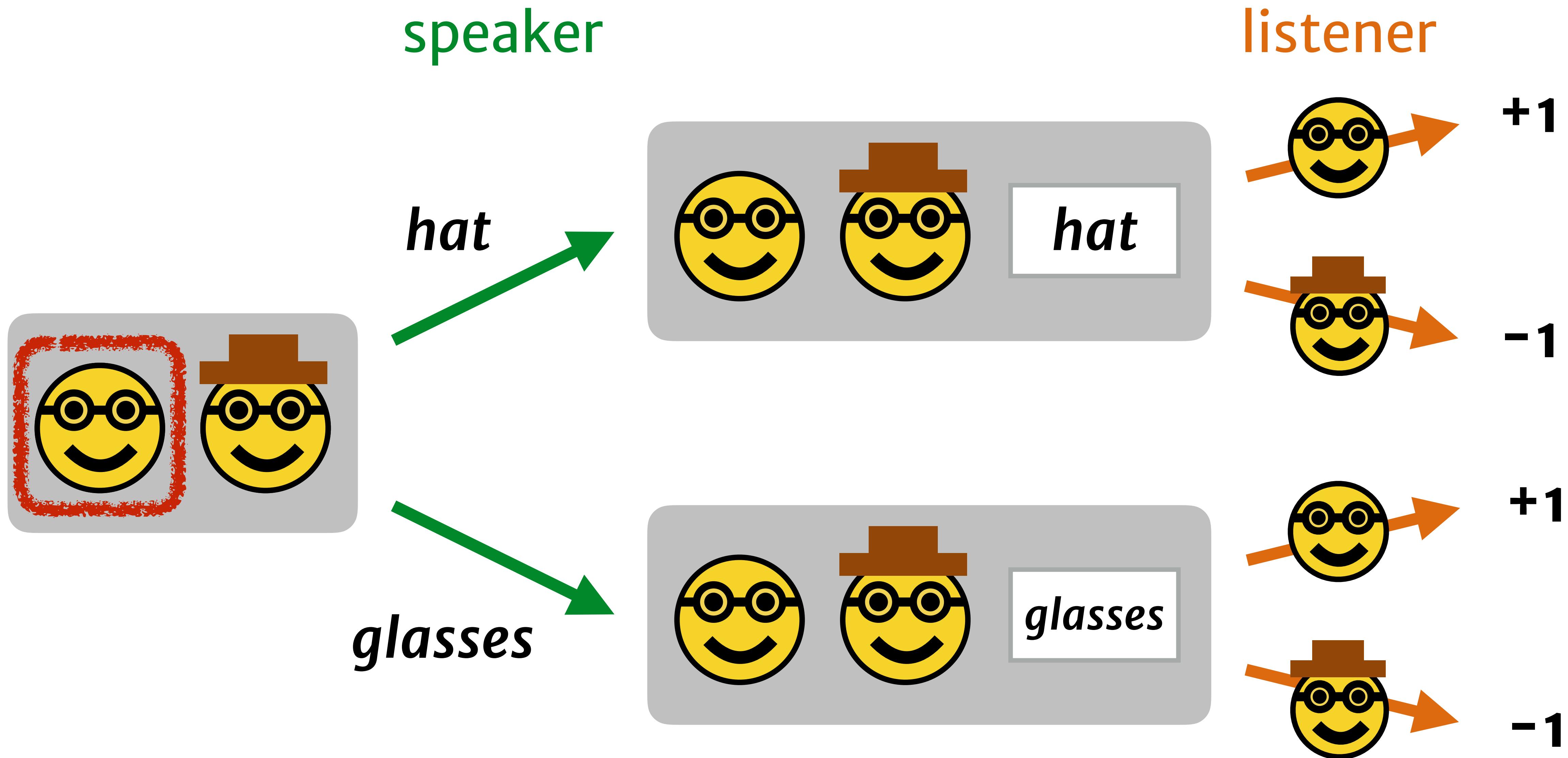


*hat*

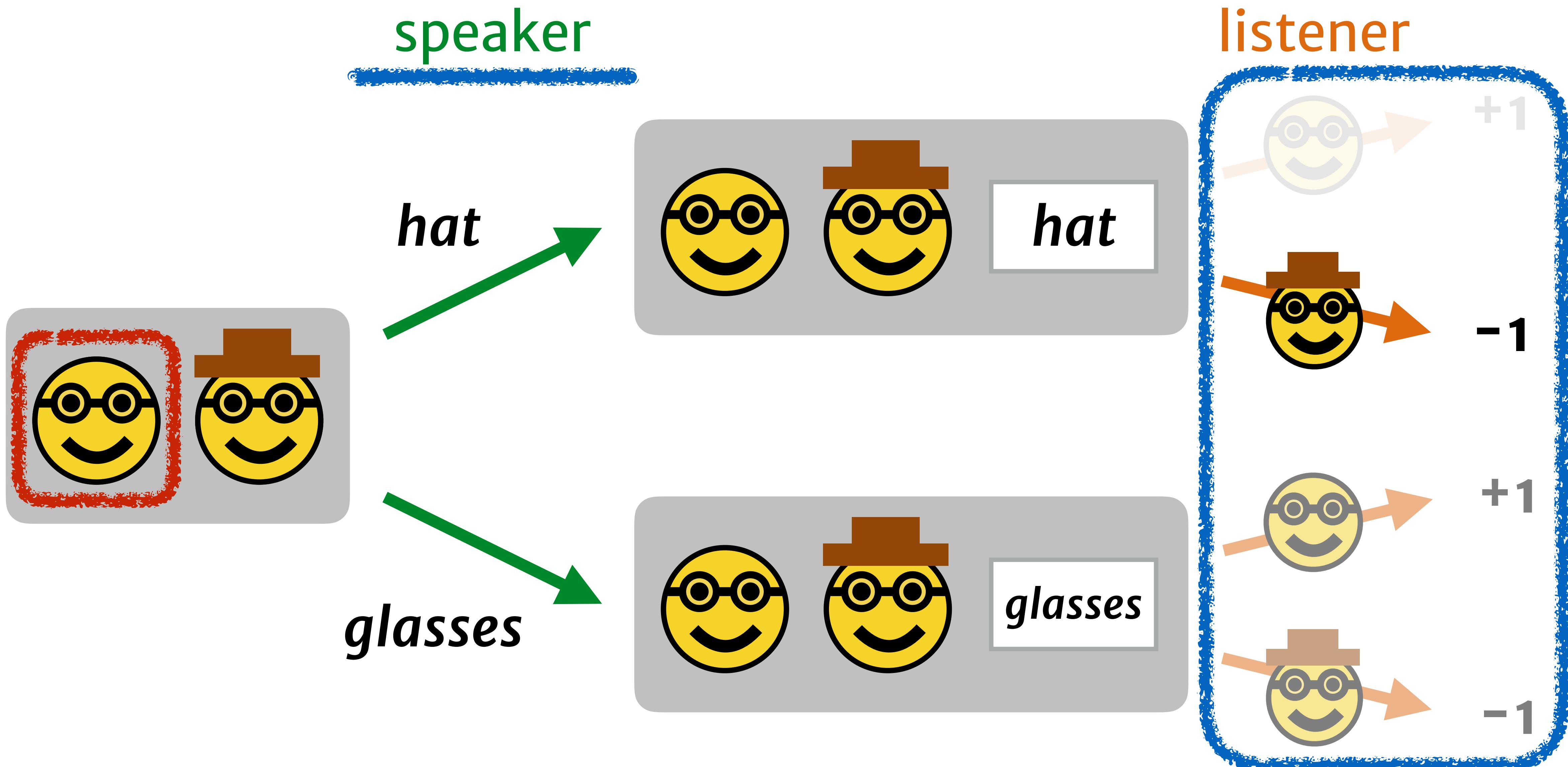


*glasses*

# RSA game tree: as speaker



# RSA game tree: as speaker

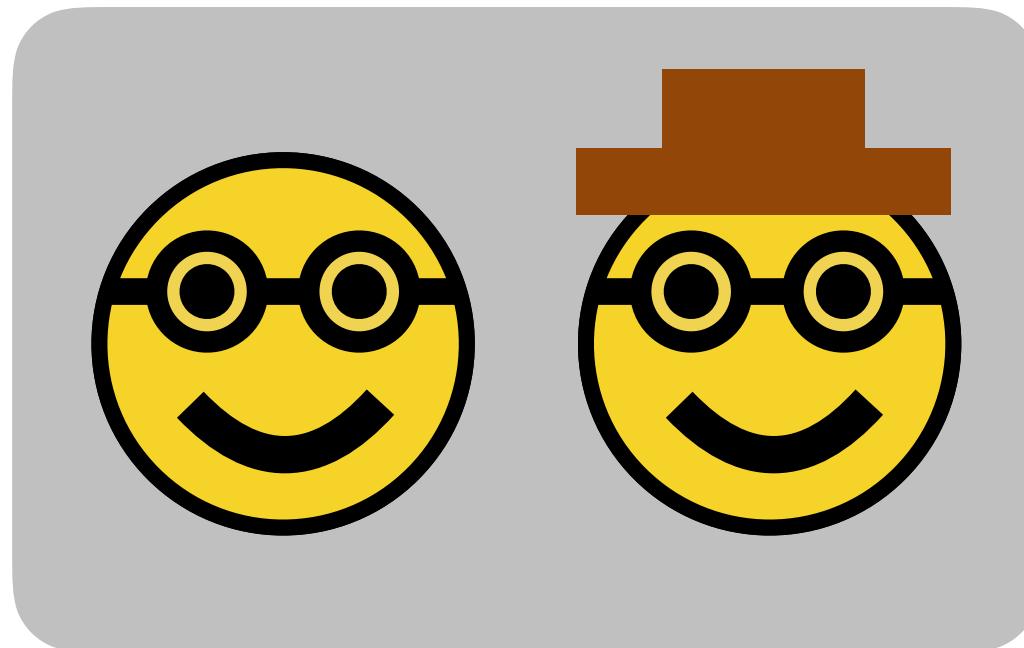


# RSA game tree: as listener

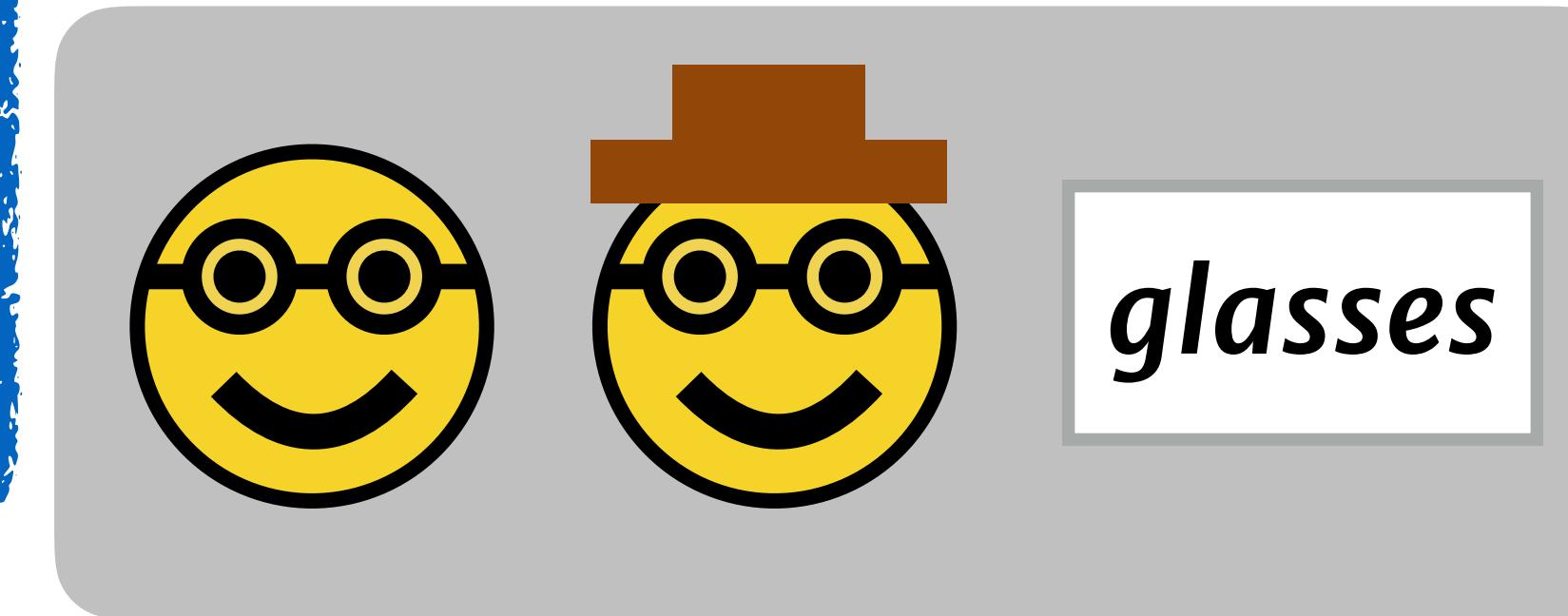
speaker

listener

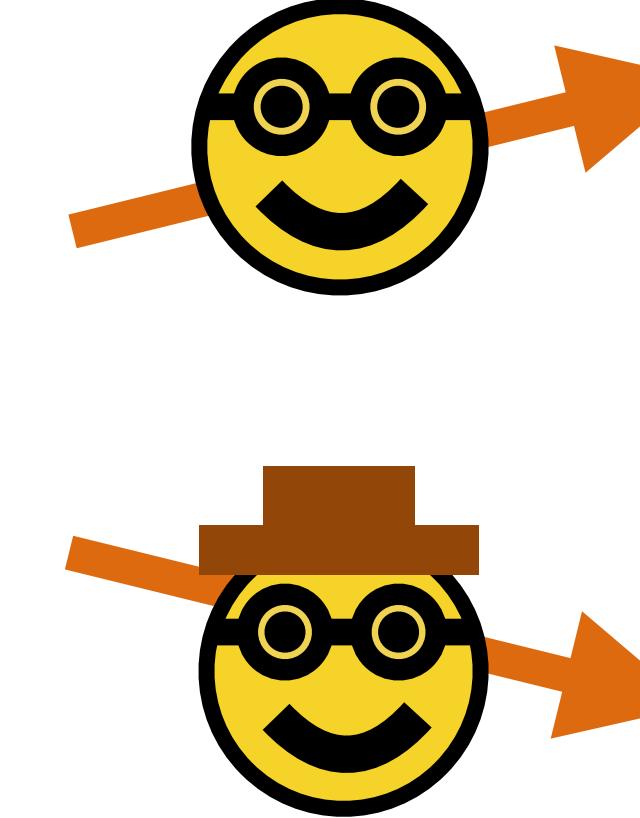
?



*glasses*



*glasses*



?

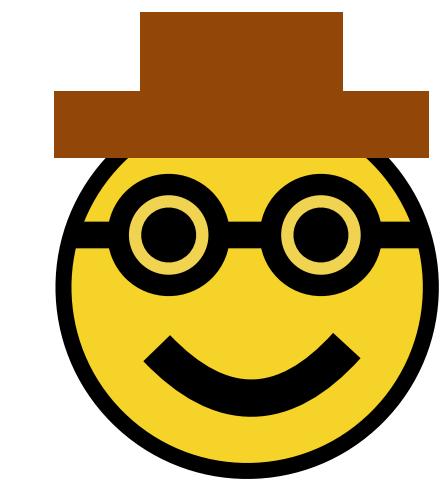
?

# A recipe for pragmatic language understanding

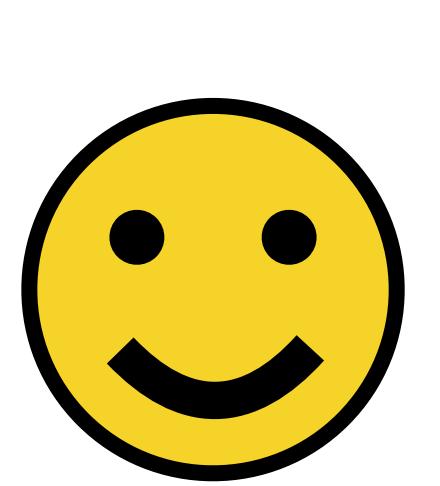
## 1. Train a base **speaker** model



smiley



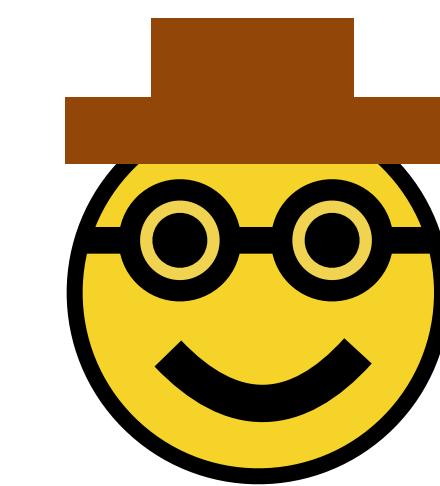
glasses  
man



plain

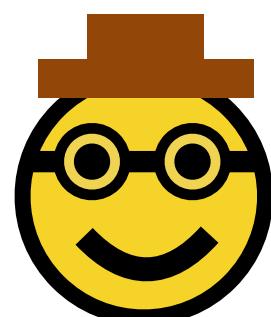
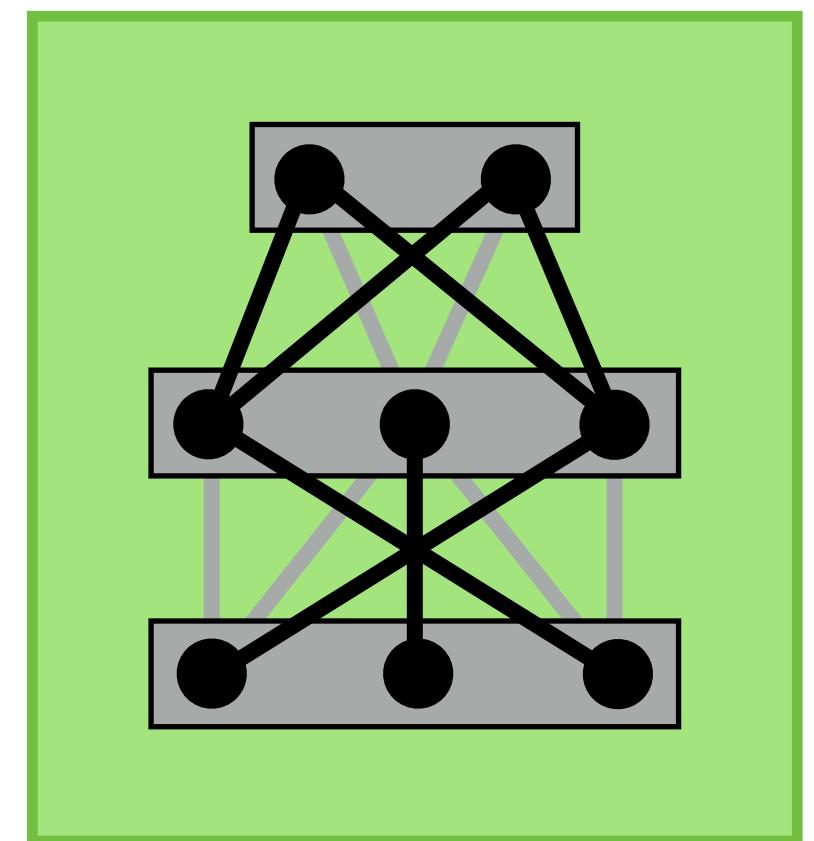
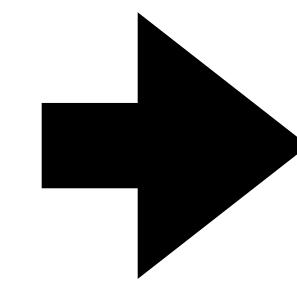


glasses



hat &  
glasses

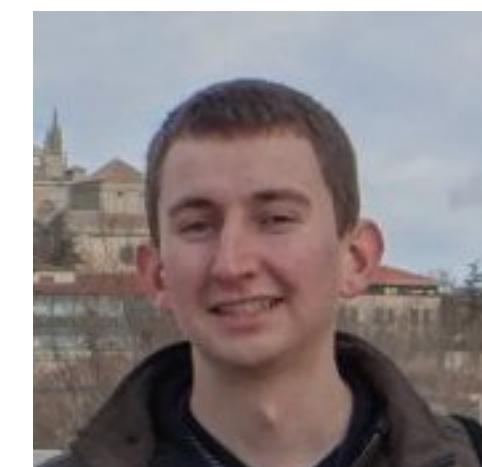
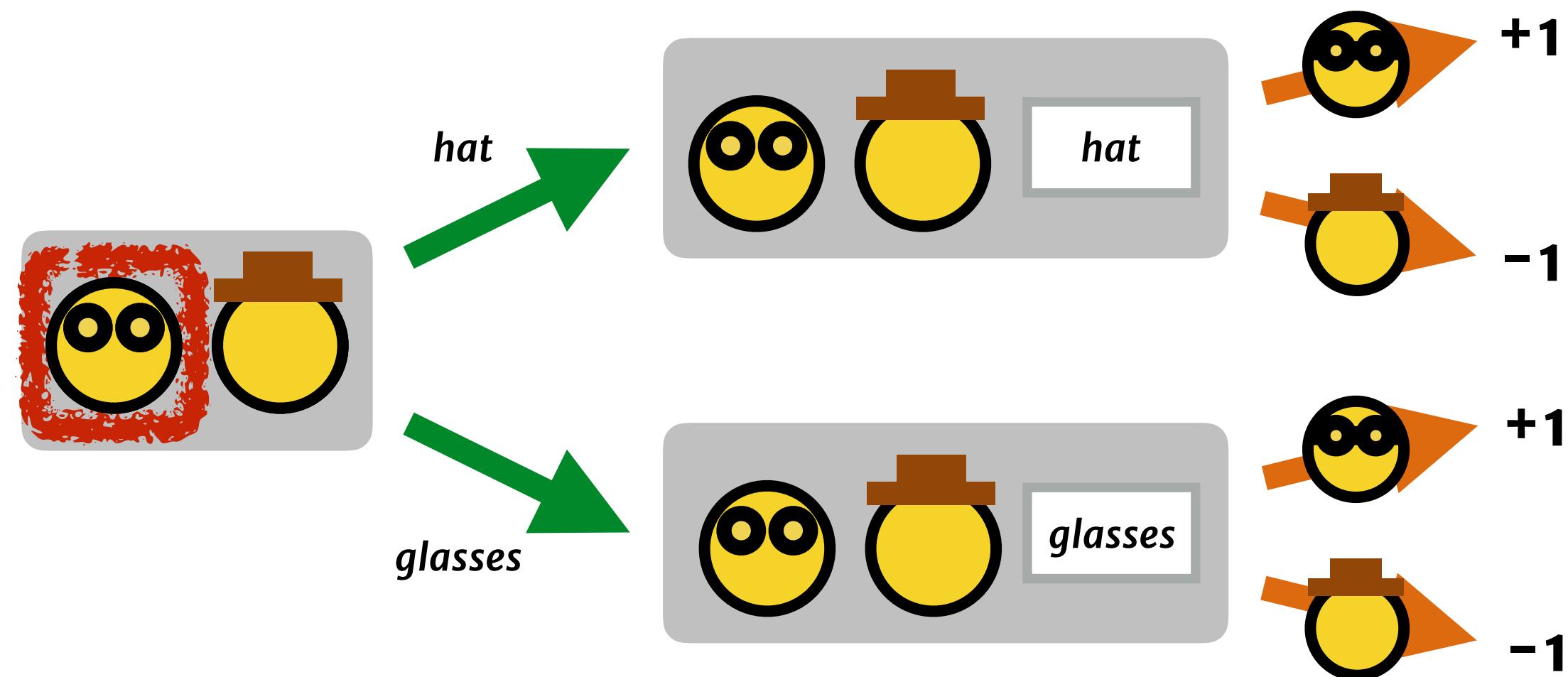
guy with  
hat      hat &  
glasses      glasses  
                man



# A recipe for pragmatic language understanding

1. Train a base **speaker** model

2. Solve this POMDP:



Daniel  
Fried



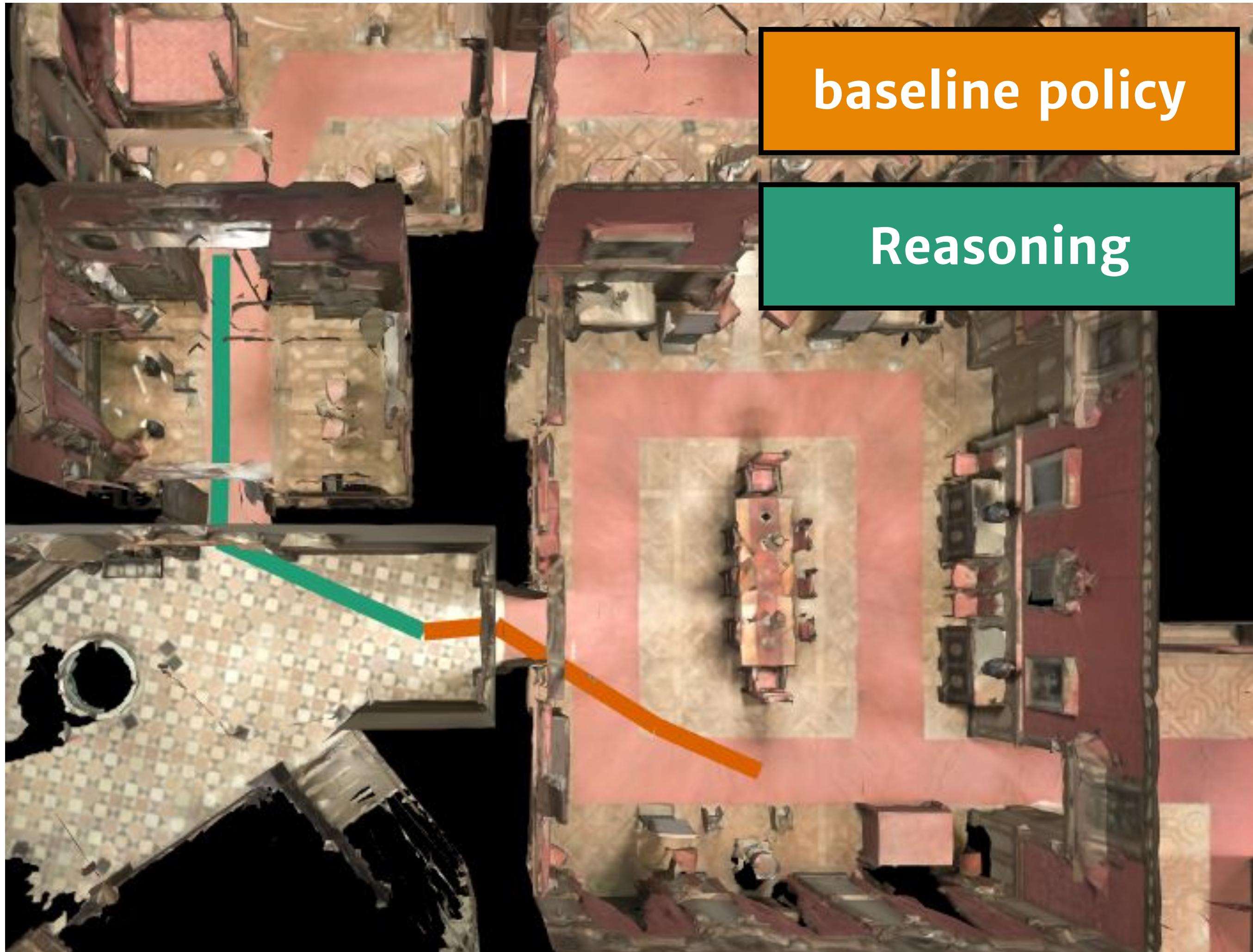
Ronghang  
Hu



Volkan  
Cirik

Speaker—follower models for vision-and-language navigation. NeurIPS 18.

# Application: instruction following



*human: Go through the door on the right and continue straight. Stop in the next room in front of the bed.*

# Application: instruction generation

**seq2seq:** *Walk past the dining room table and chairs and wait there.*

**reasoning:** *Walk past the dining room table and chairs and take a right into the living room. Stop once you are on the rug.*

**human:** *Turn right and walk through the kitchen. Go right into the living room and stop by the rug.*



# Lesson

---

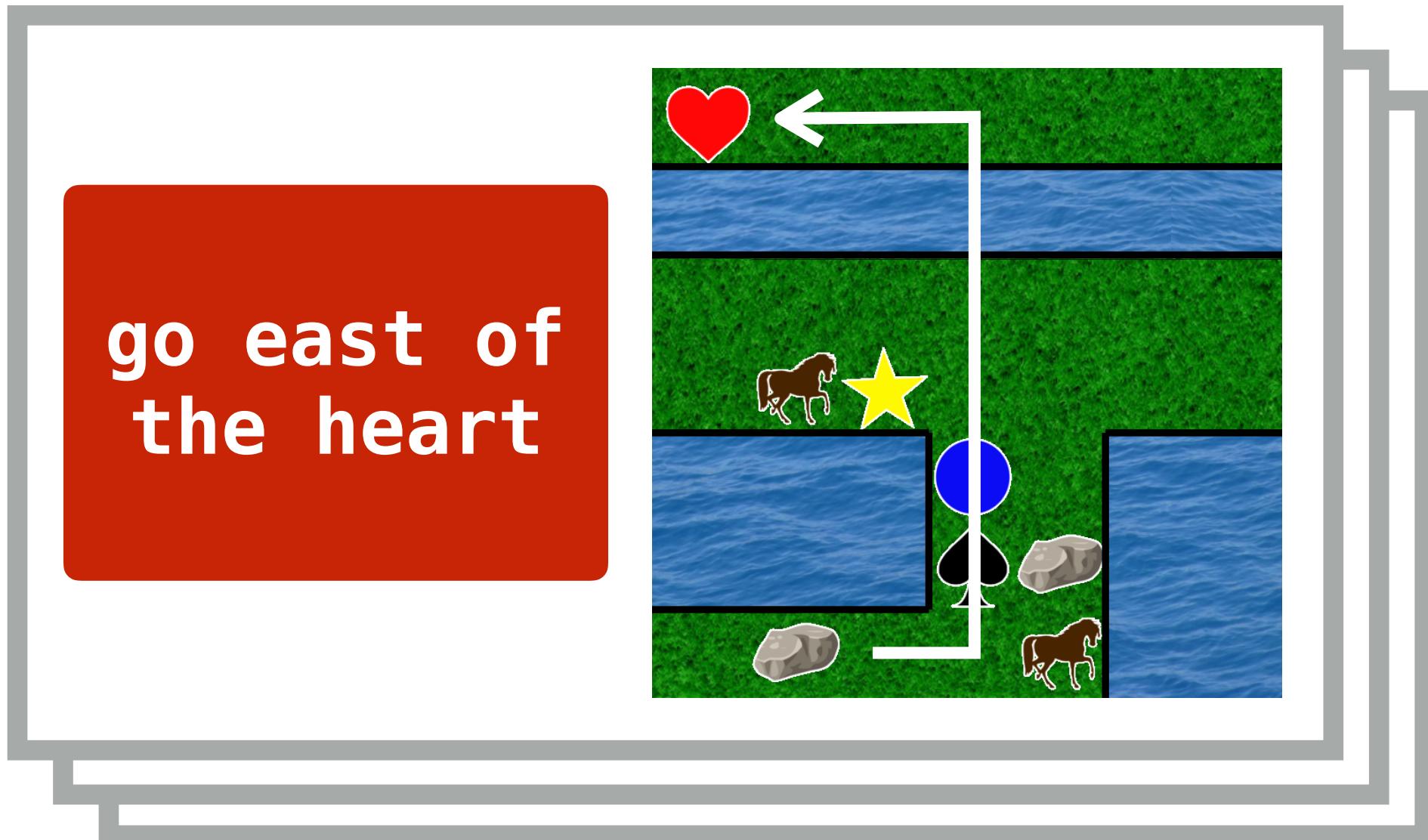
**Utterances are chosen to facilitate  
correct interpretation in context.**

(This makes the learning problem easier!)

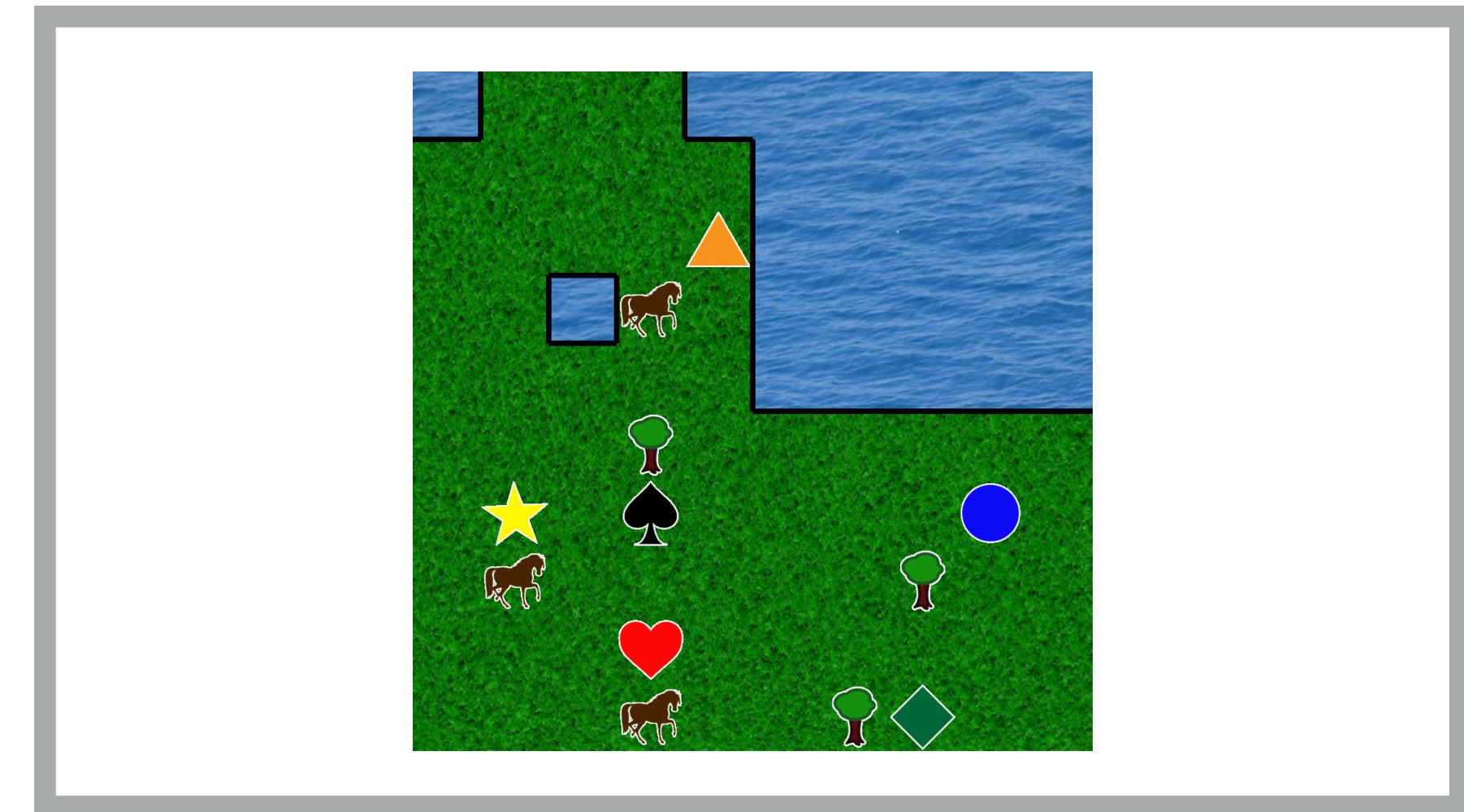
Language as a scaffold  
for learning

# What else is an instruction follower good for?

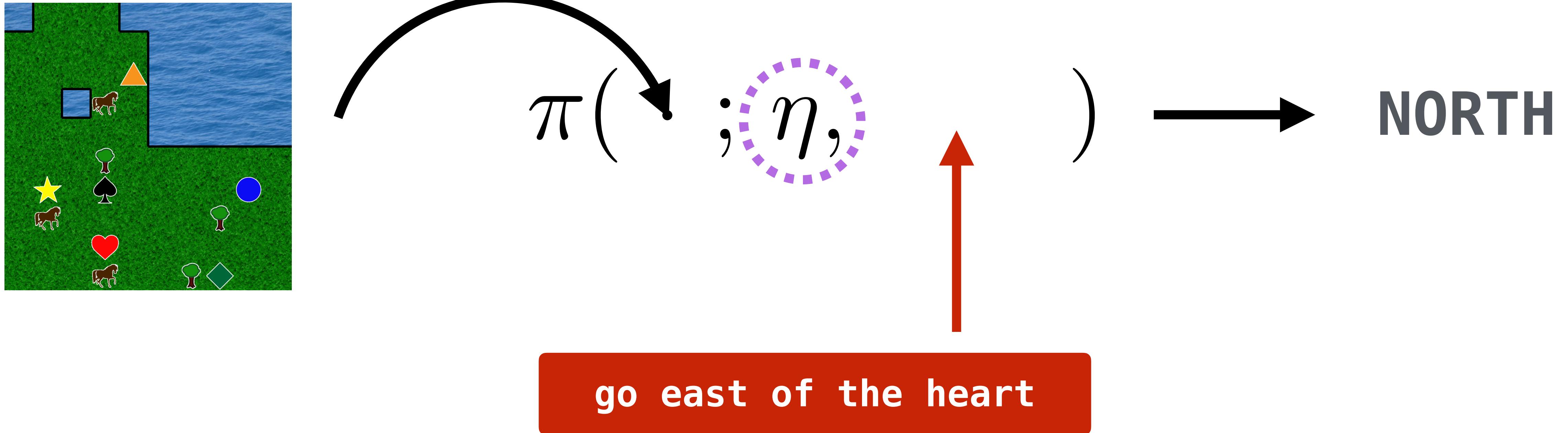
## Language learning



## Reinforcement learning

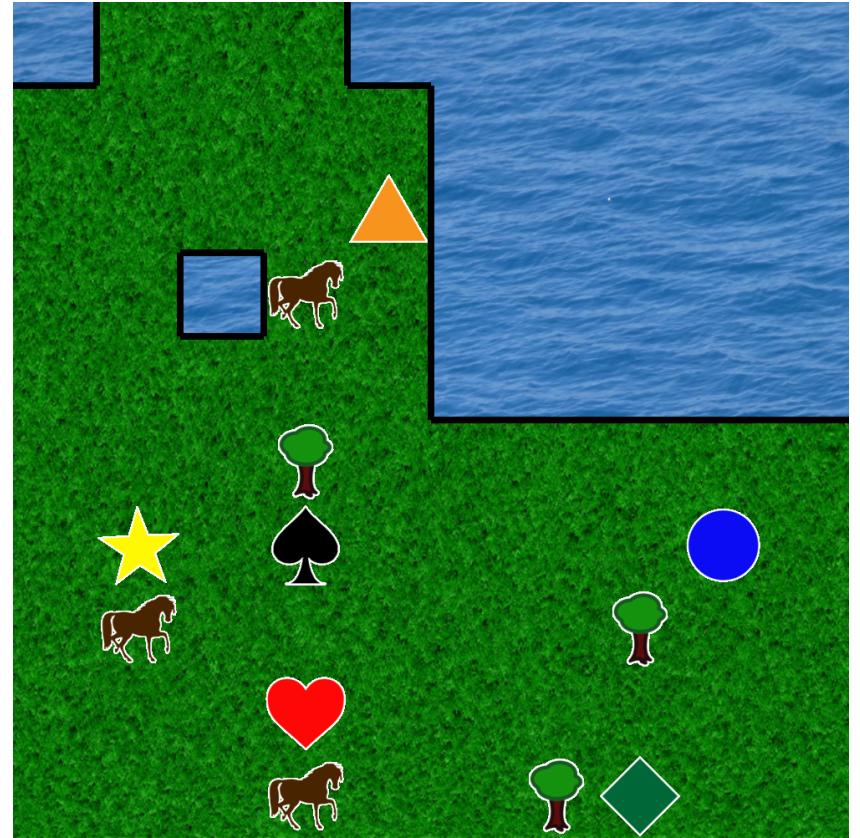


# Pretraining via language learning



# (Standard) reinforcement learning

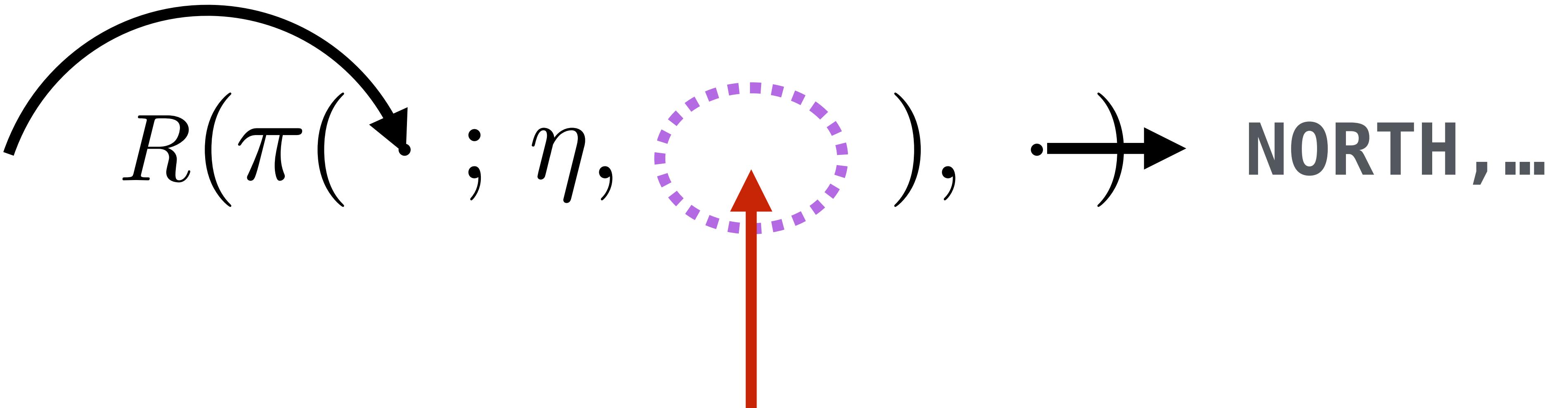
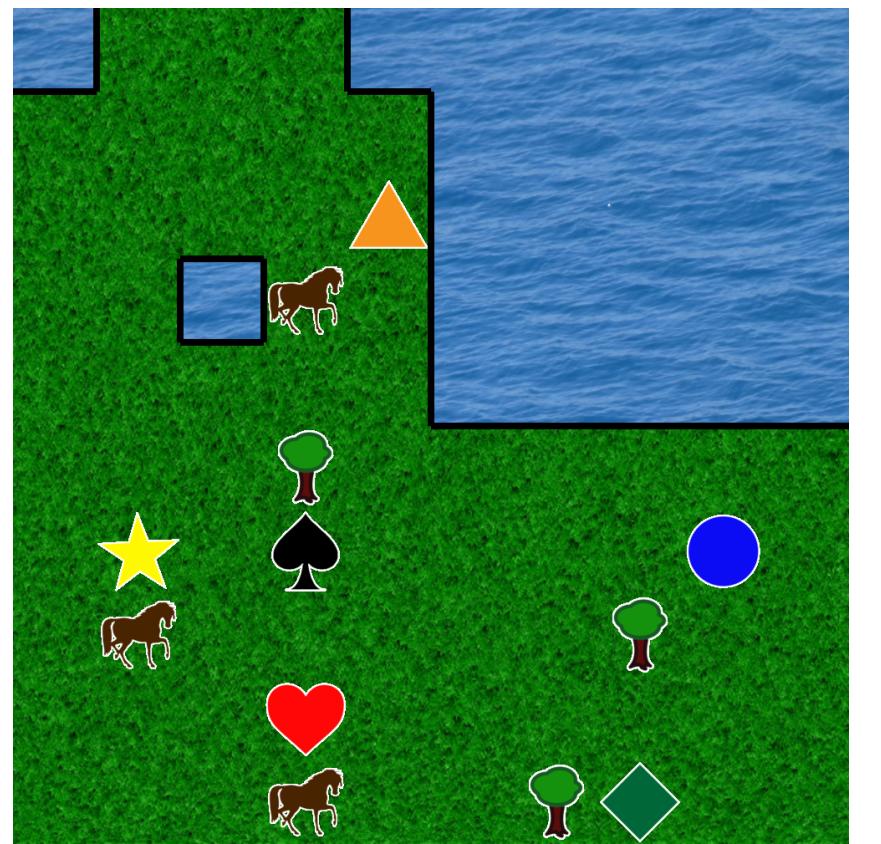
---



$$R(\pi( \xrightarrow{\hspace{1cm}} ; \eta, \textcircled{?} ), \xrightarrow{\hspace{1cm}} \text{???})$$

# Concept learning

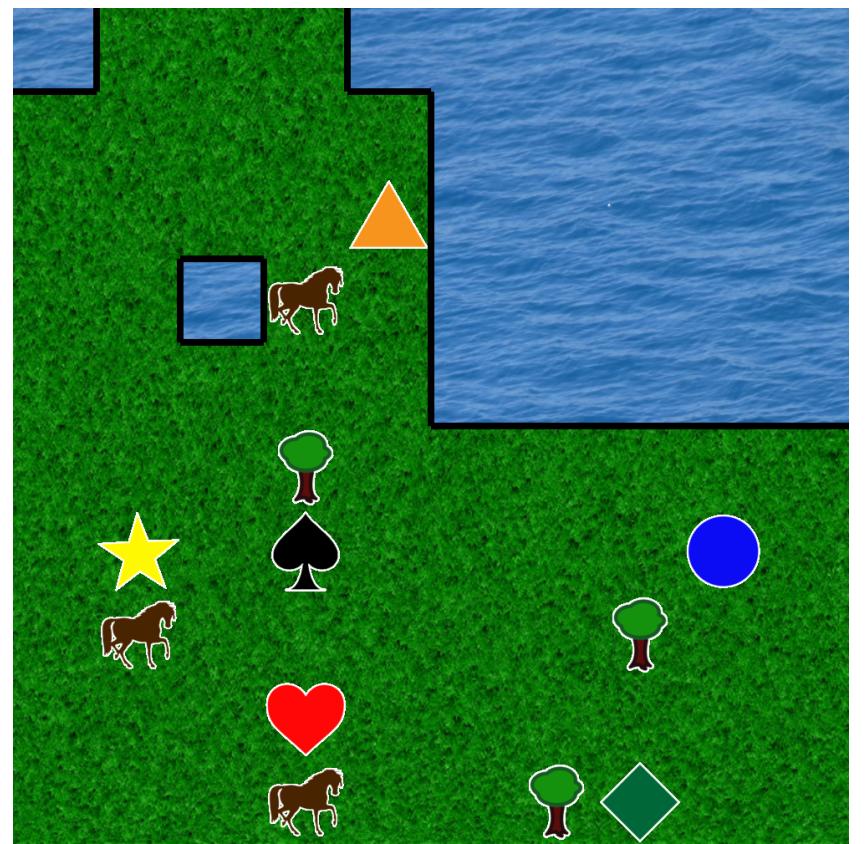
---



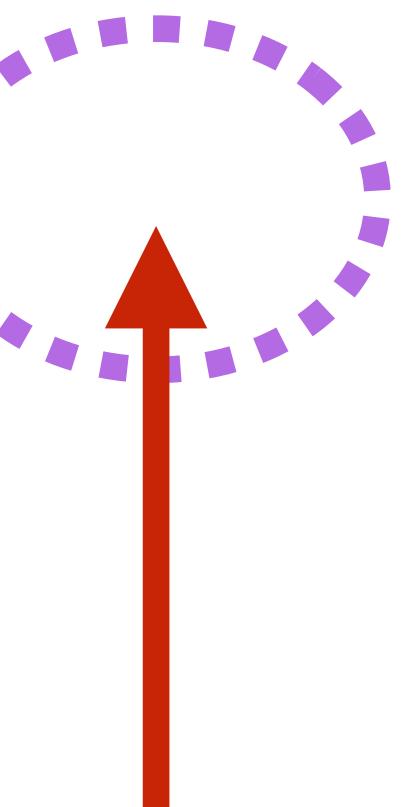
find the horse

# Concept learning

---



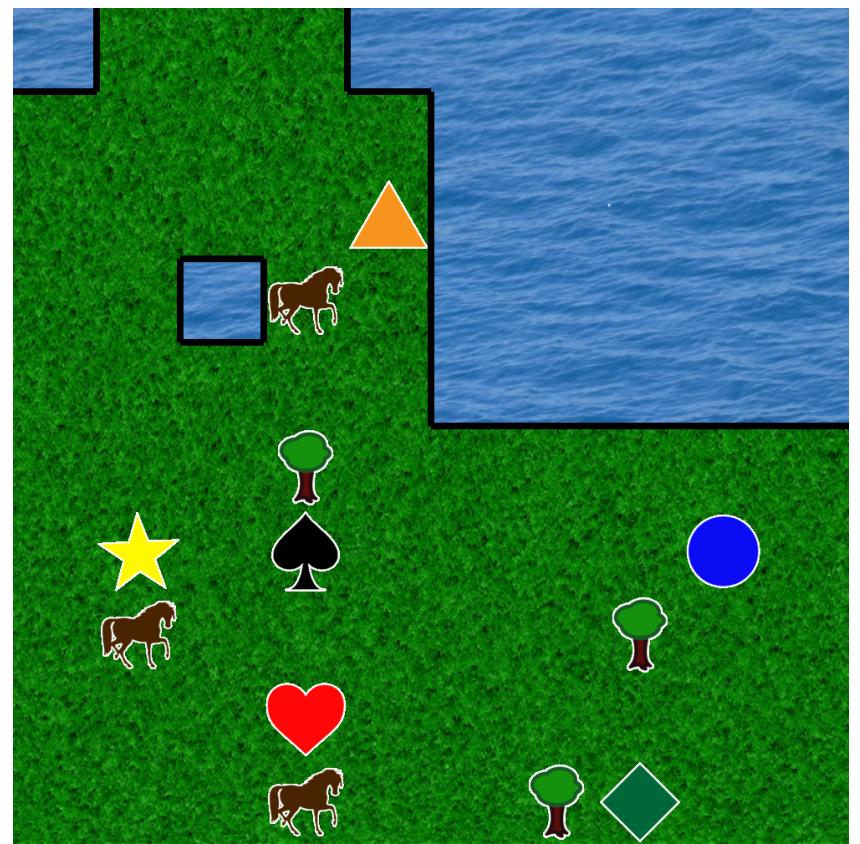
$R(\pi( \xrightarrow{\hspace{1cm}} ; \eta, \xrightarrow{\hspace{1cm}} ), \xrightarrow{\hspace{1cm}} \text{NORTH}, \dots)$



find the horse

-0.52

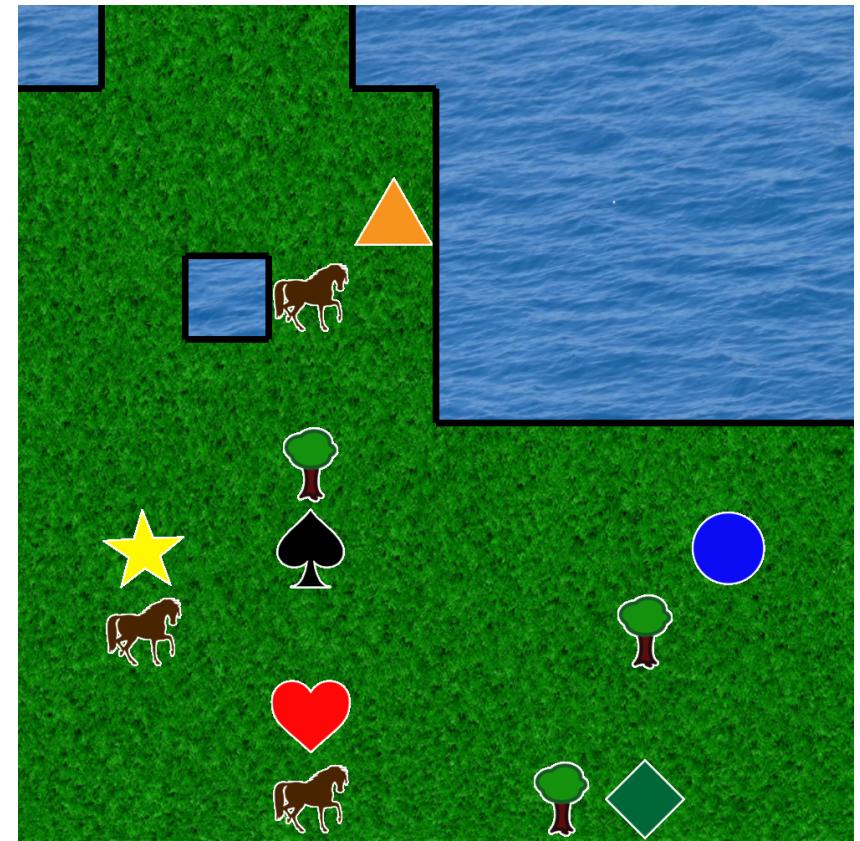
# Concept learning



$R(\pi( \xrightarrow{\hspace{1cm}} ; \eta, \circlearrowleft ), \cdot) \rightarrow \text{SOUTH, ...}$

find the horse	-0.52
left of heart	0.33

# Concept learning

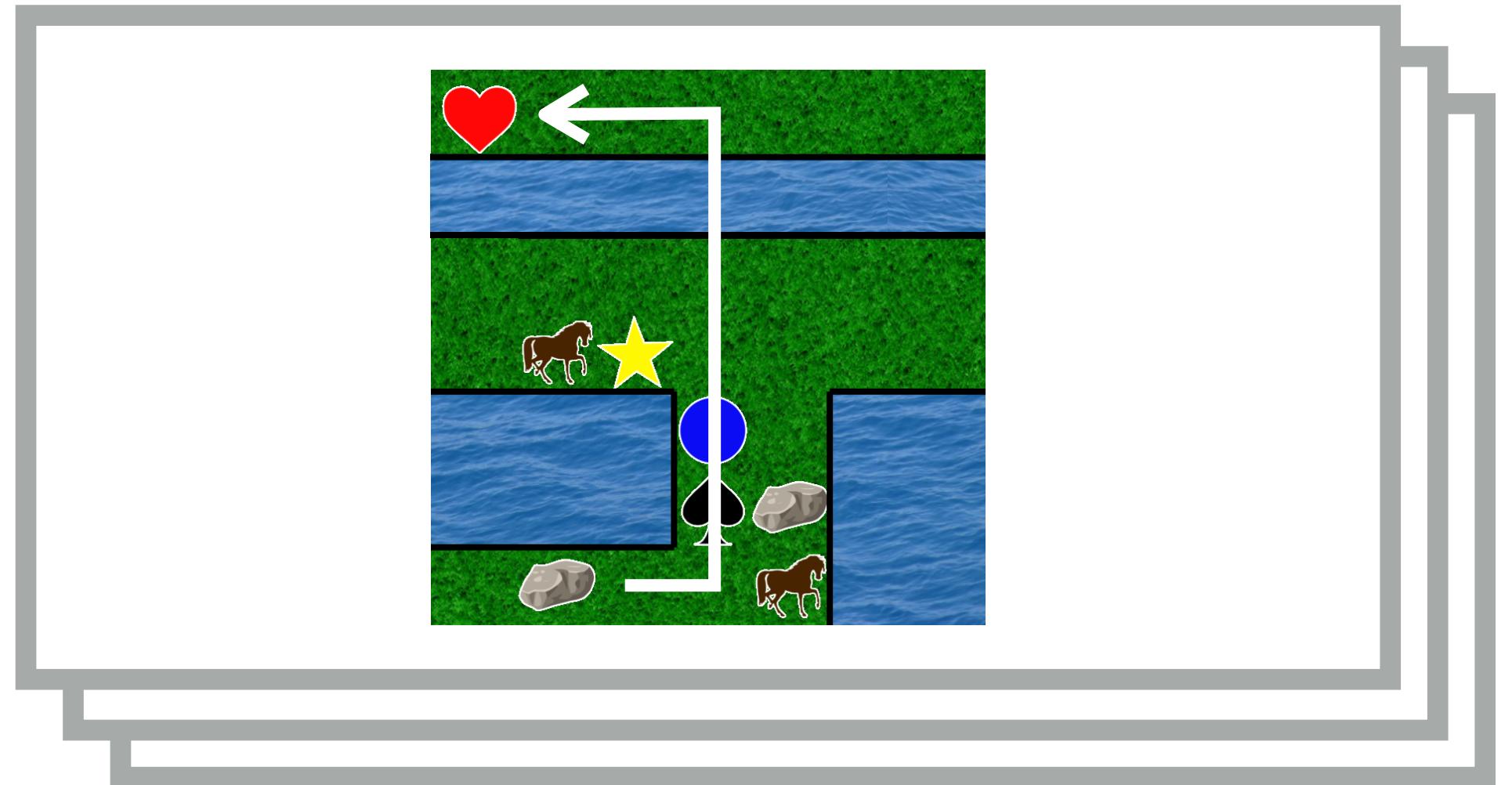


$R(\pi( \xrightarrow{\hspace{1cm}} ; \eta, \xrightarrow{\hspace{1cm}} ), \xrightarrow{\hspace{1cm}} \text{SOUTH}, \dots)$

find the horse	-0.52
left of the heart	0.33
heart east side	0.95

# As multitask learning

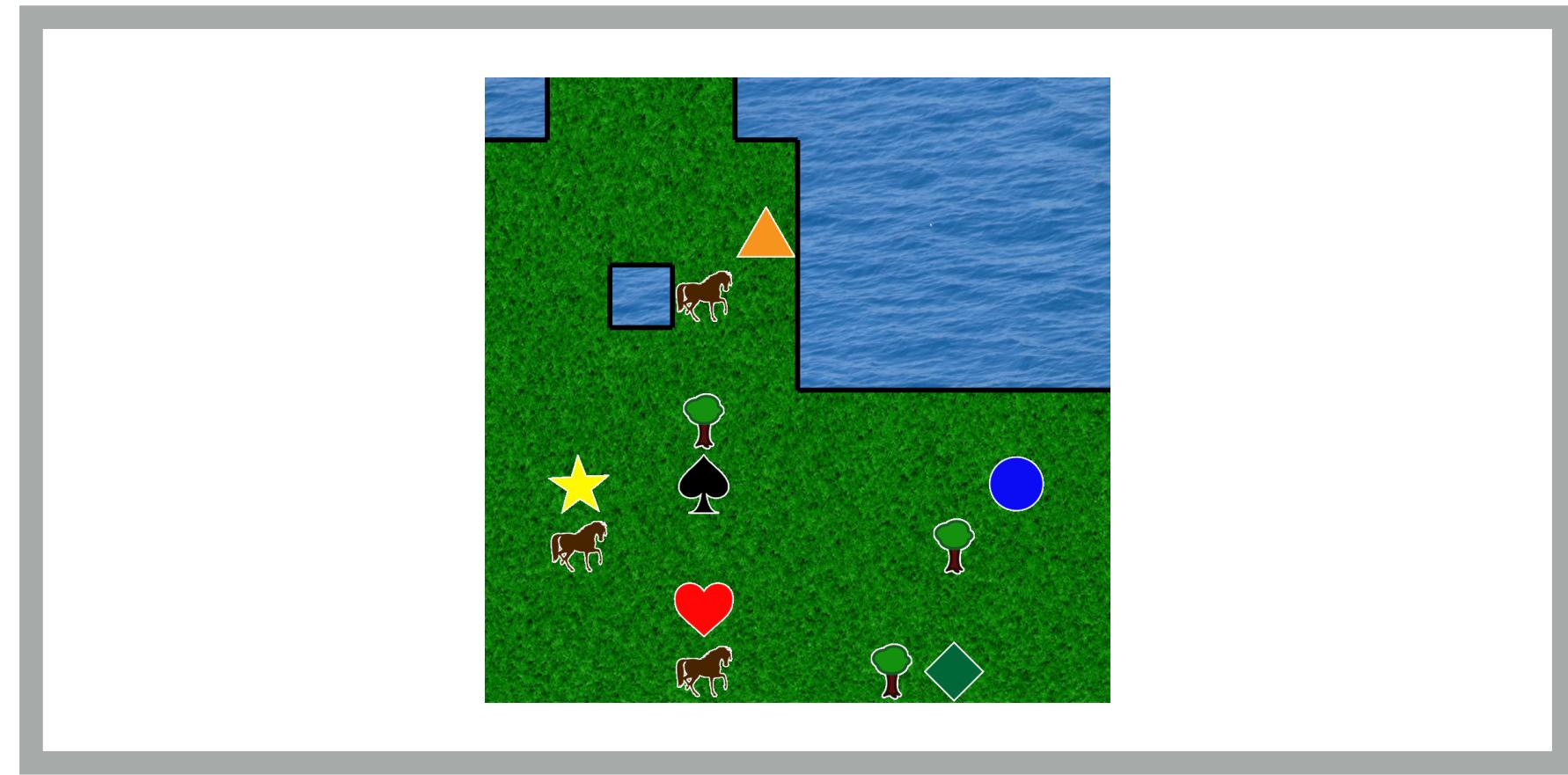
## Language learning



$$\arg \min_{\eta} R(\pi(\nearrow | \text{[image]} ; \eta, \uparrow))$$

go east of the heart

## Reinforcement learning



$$\arg \min R(\pi(\searrow | \text{[image]} ; \eta, \uparrow))$$

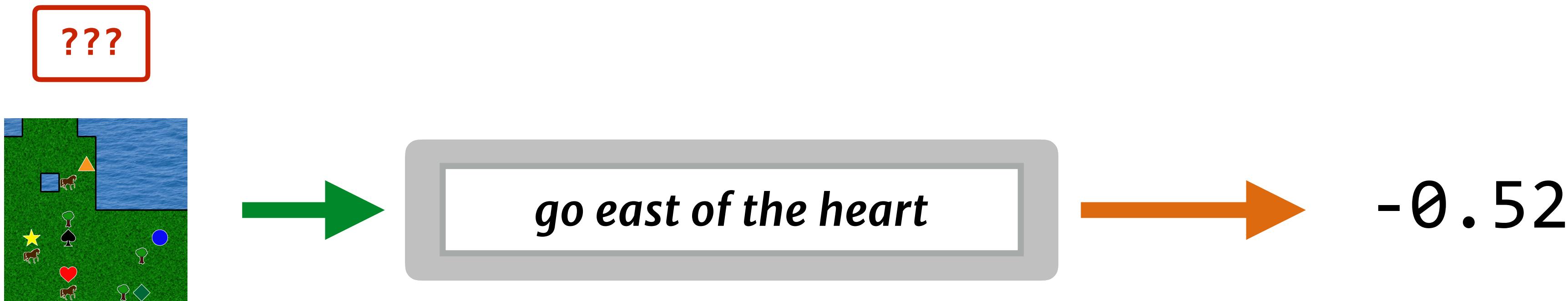
find the triangle

# As a language game...

---

$\arg \min$

$\pi \circ R$

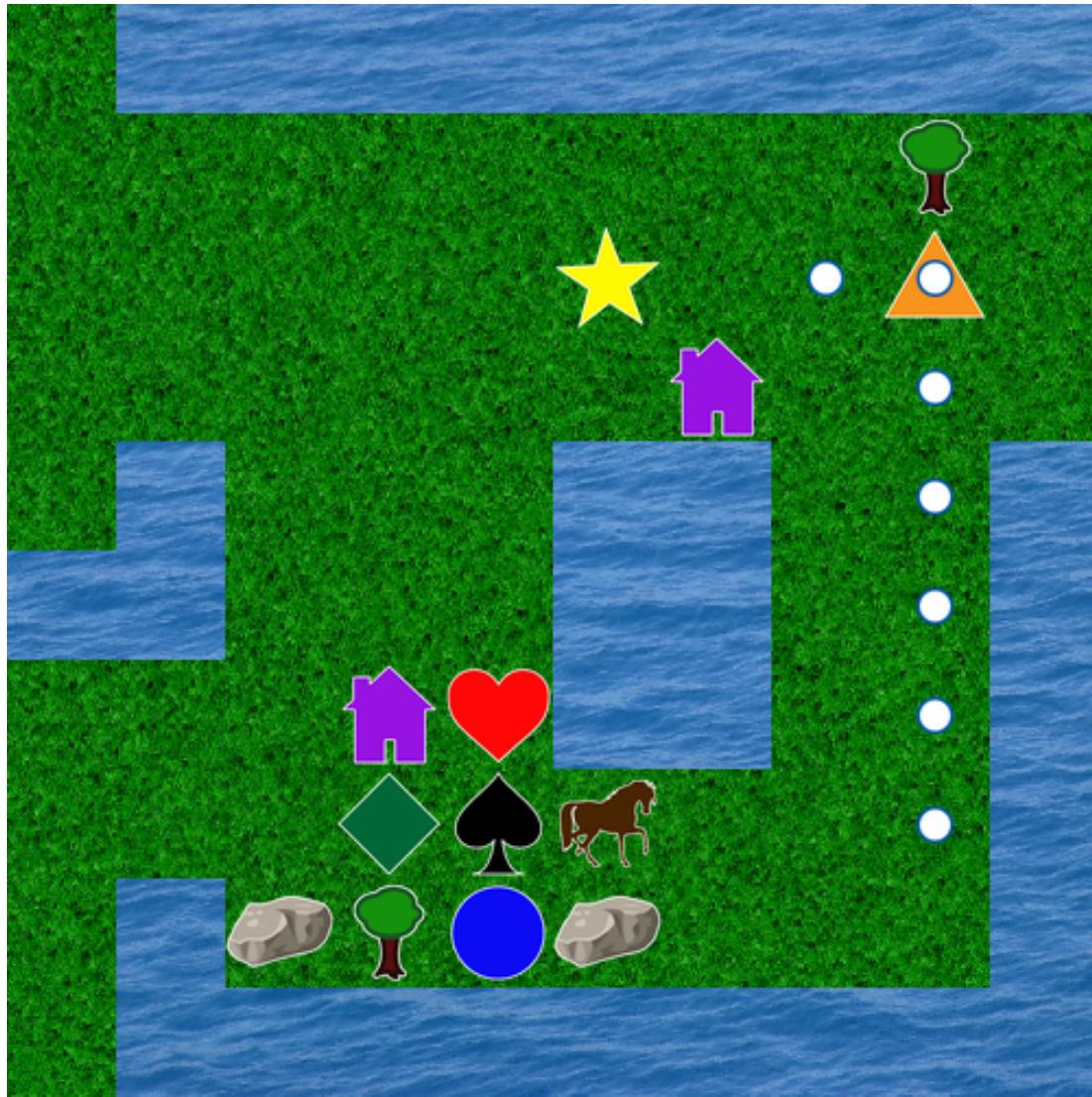


**speaker model**

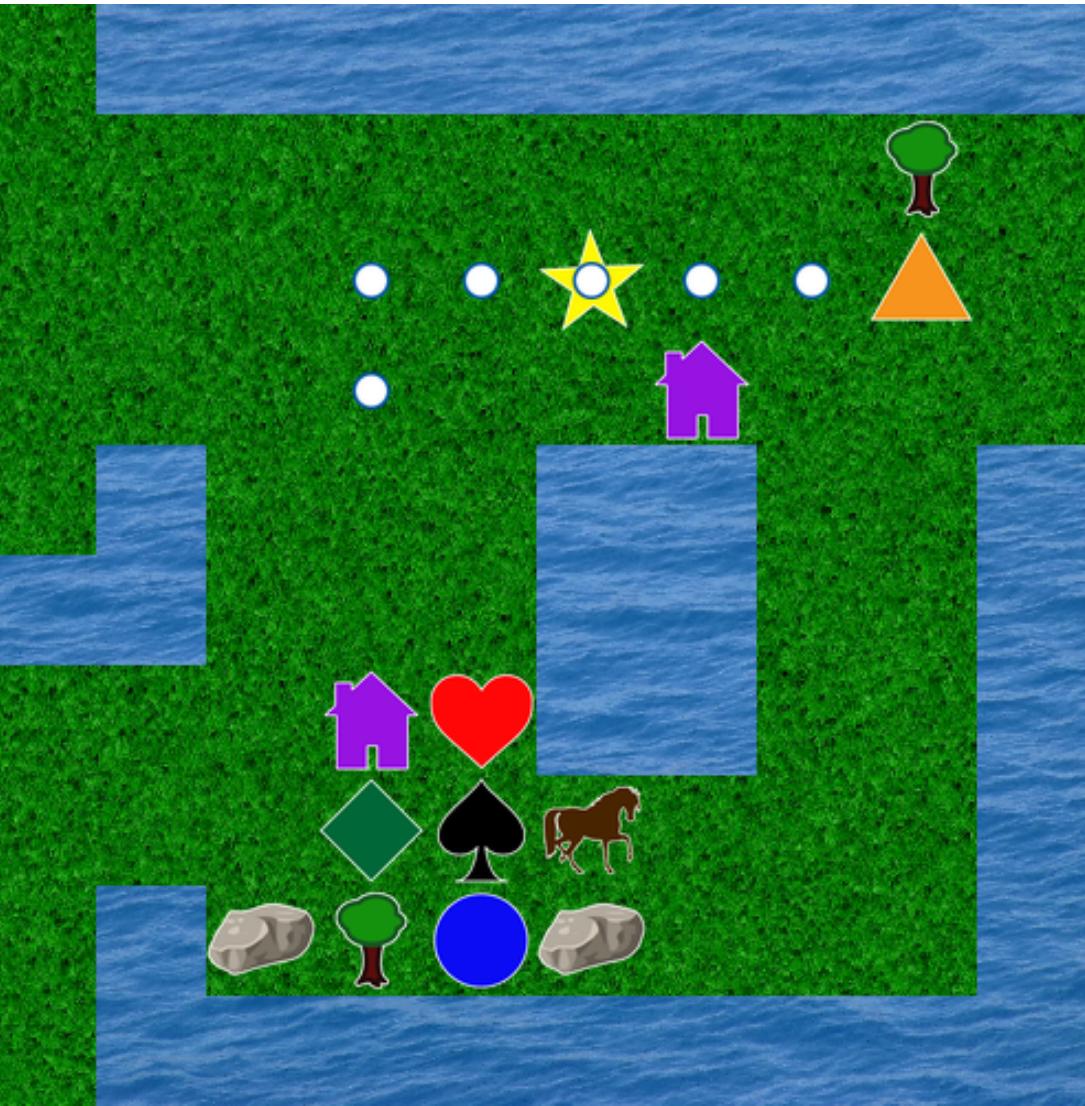
**listener loss**



# Results



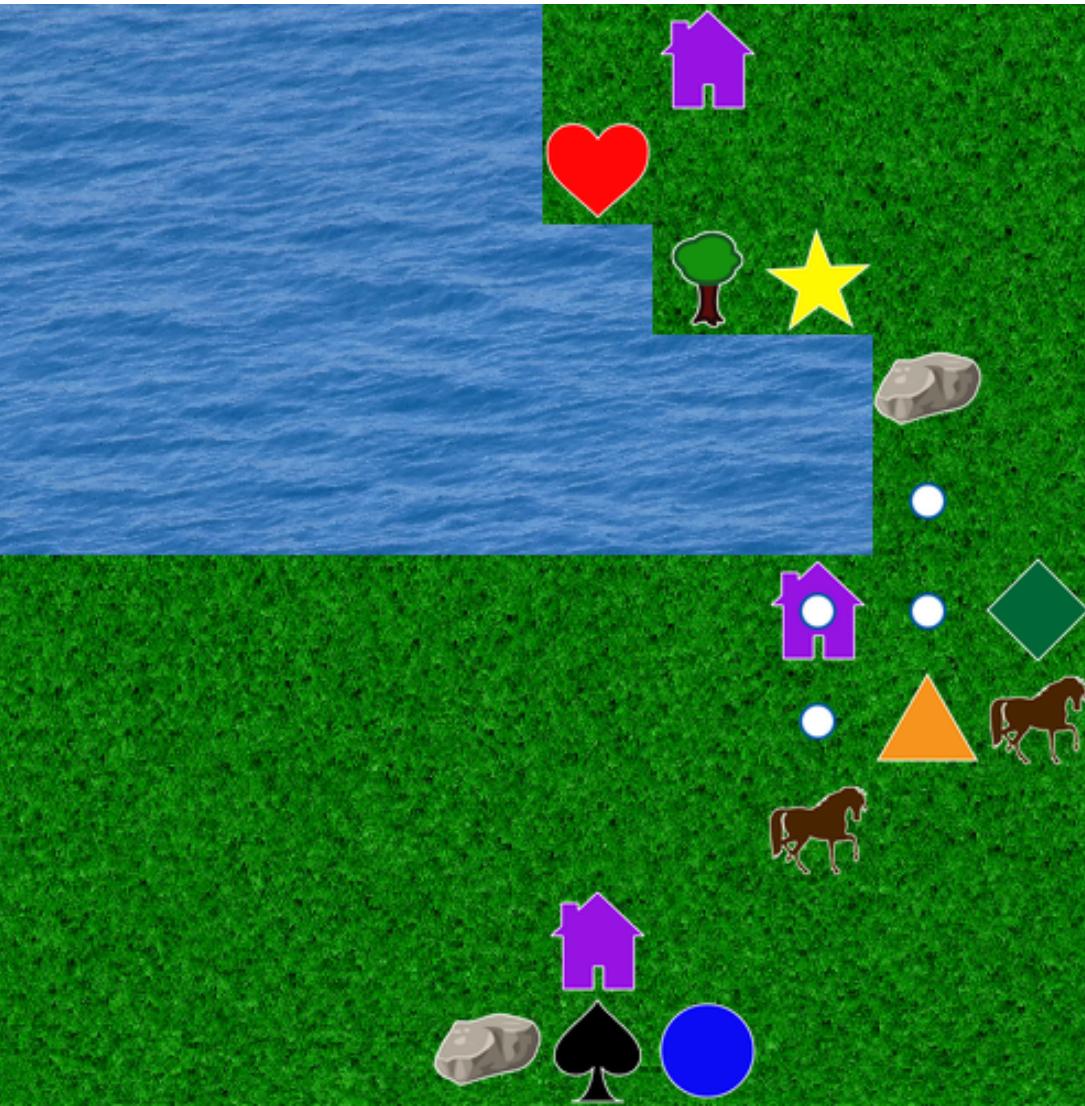
True description



reach cell on left of triangle

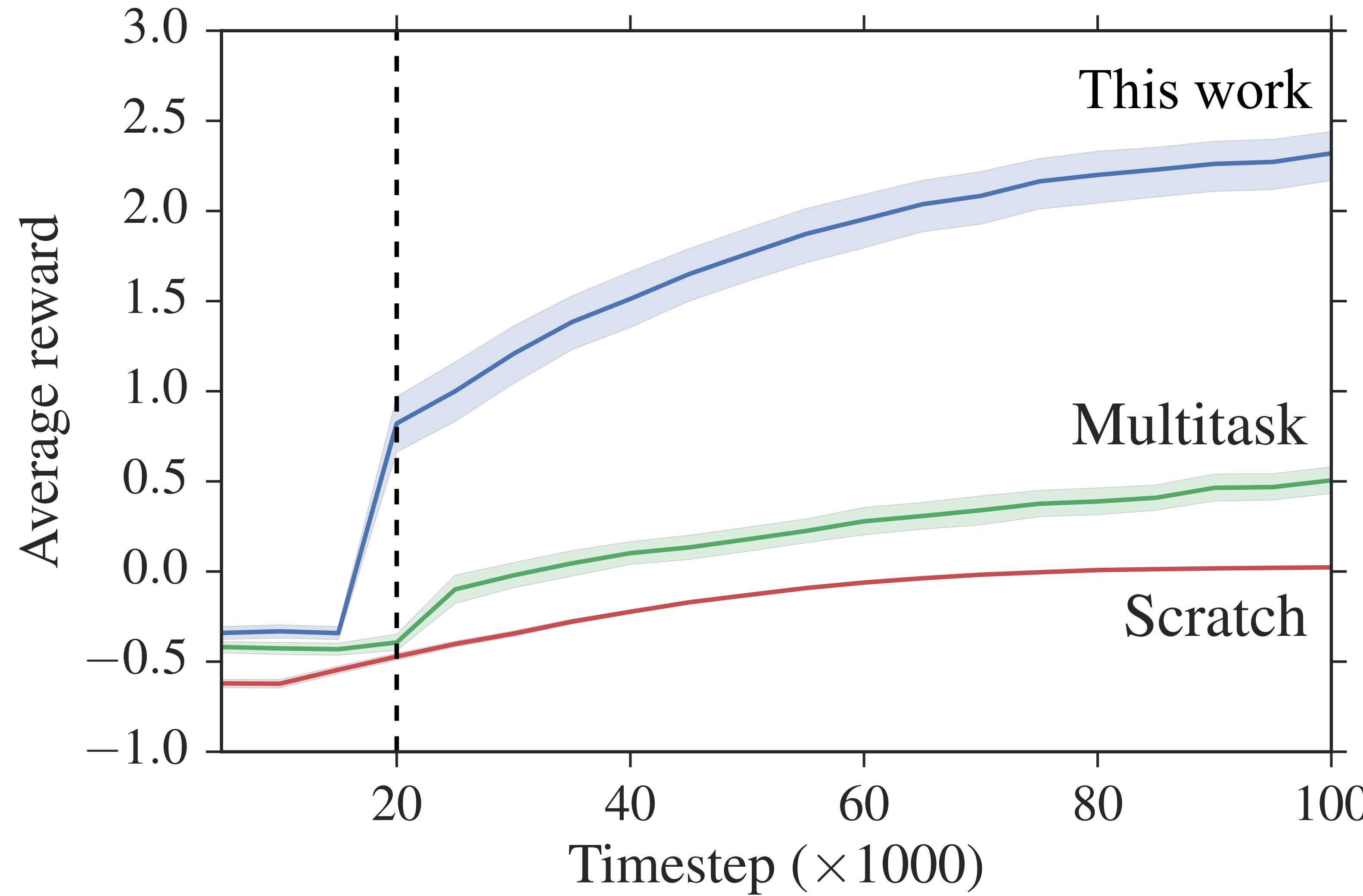
Pred description

reach square left of triangle



# Results: RL

---



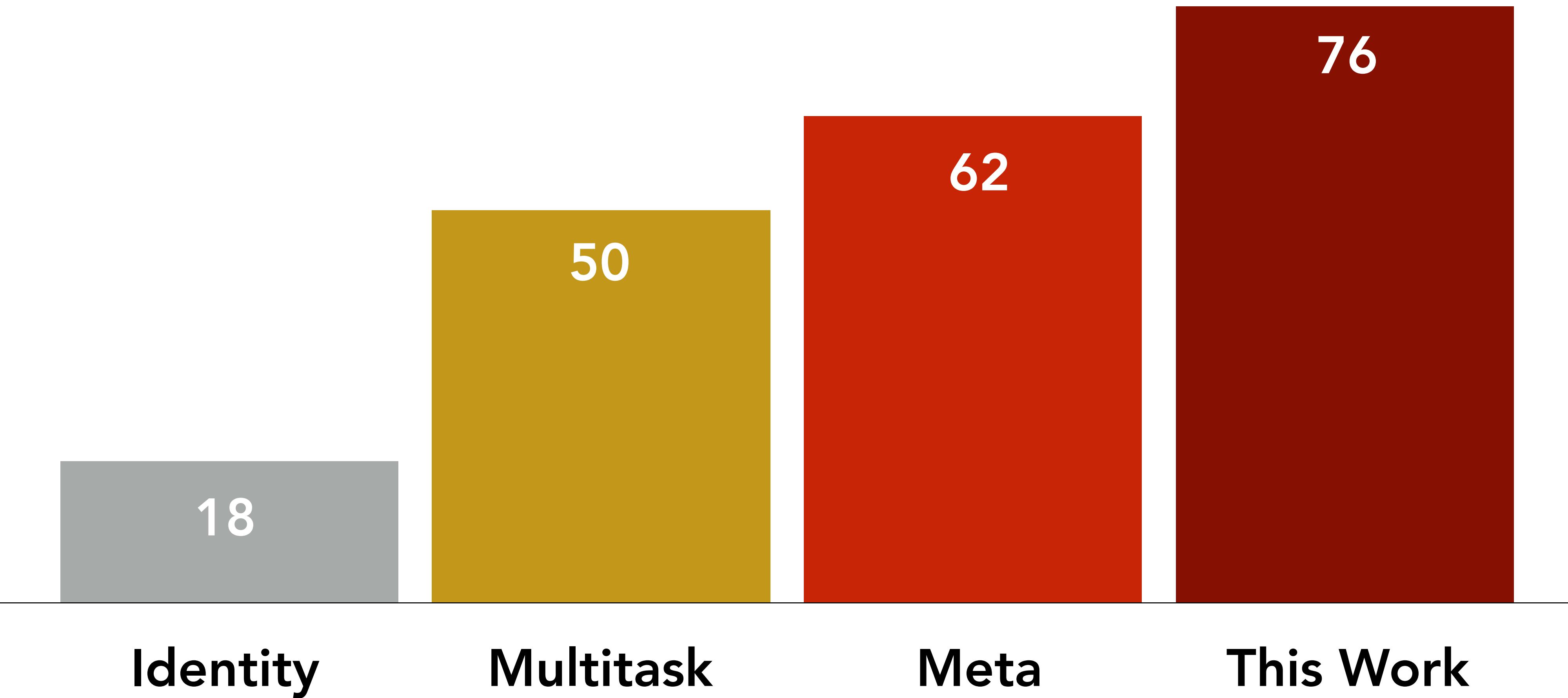


# Results

examples	true description	true output
emboldens	emboldecs	loocies
kisses	kisses	loonies
loneliness →	replace all n s with c	↑
vein	locelicess	loonies
dogtrot	veic	↓
	dogtrot	loocies
	change any n to a c	
	pred. description	pred. output

# Results: programming by demonstration

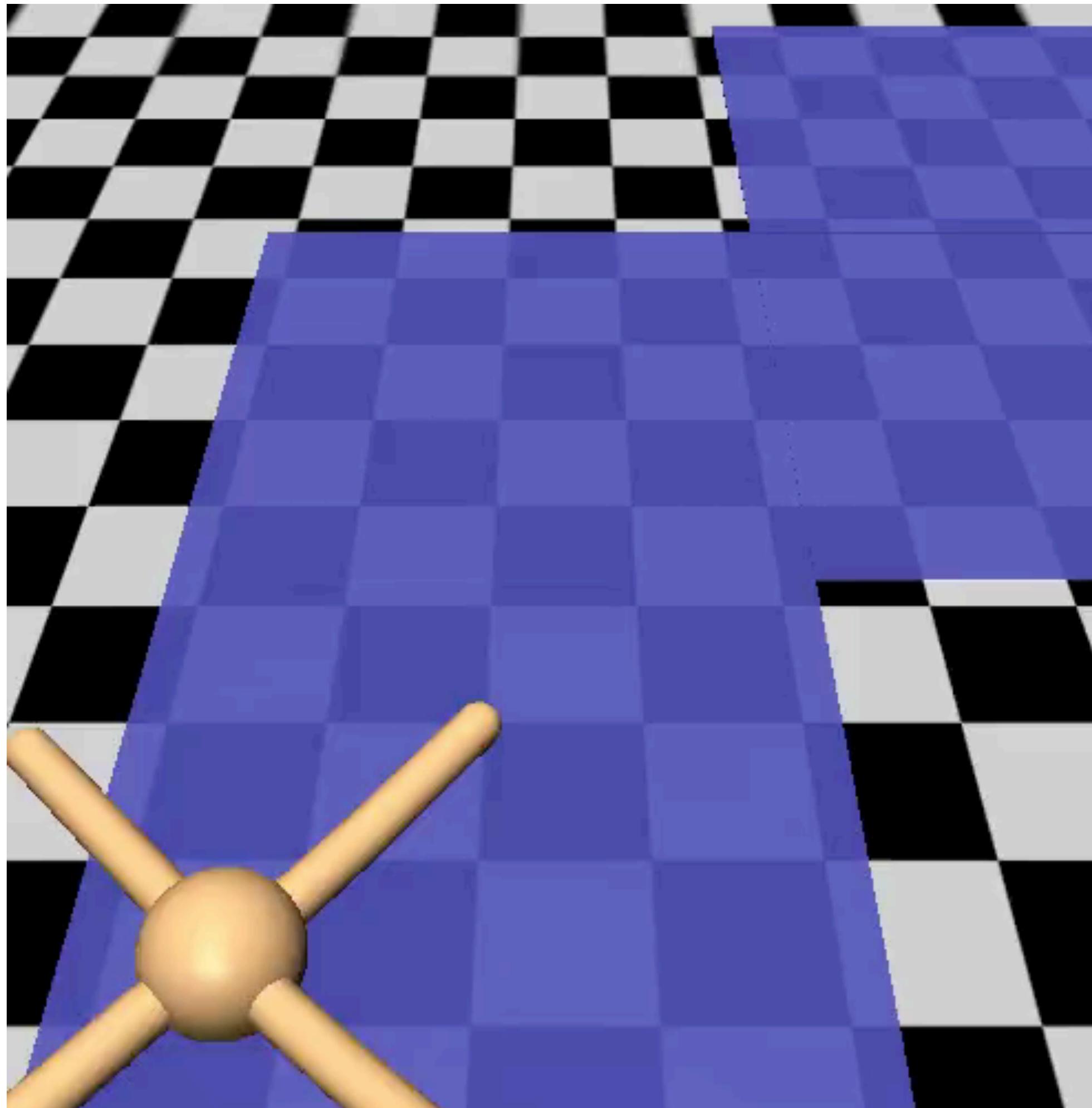
---



# Results: locomotion

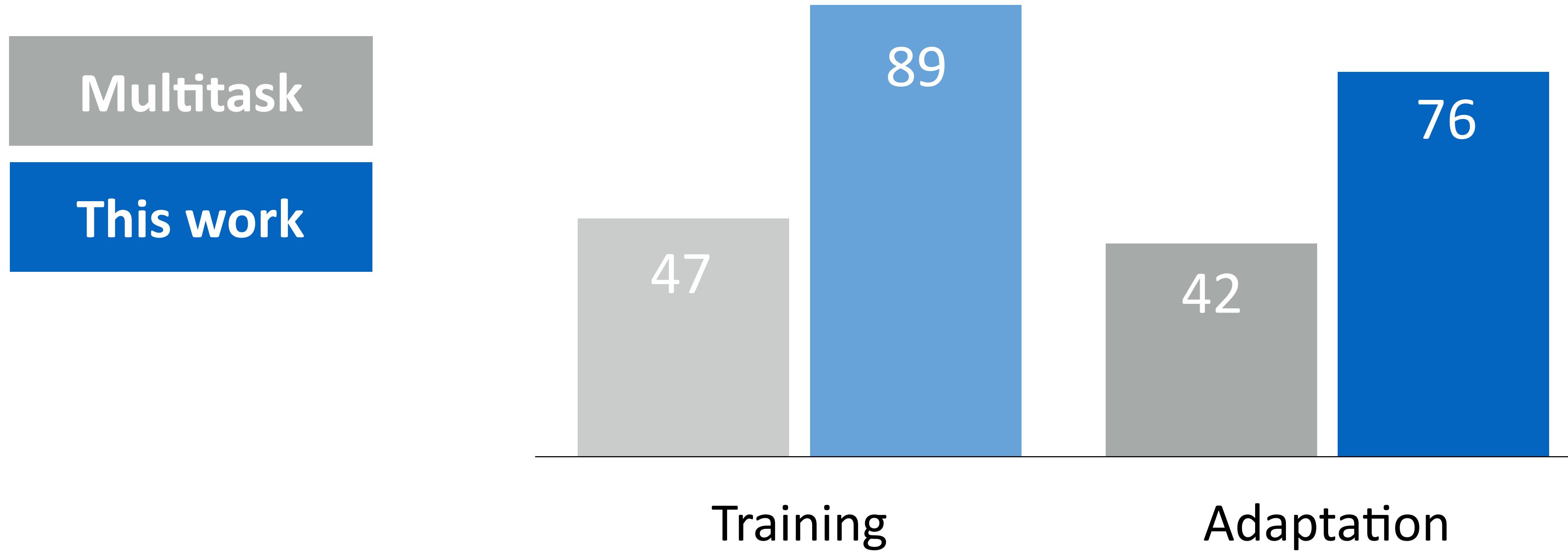
---

north, east,  
north



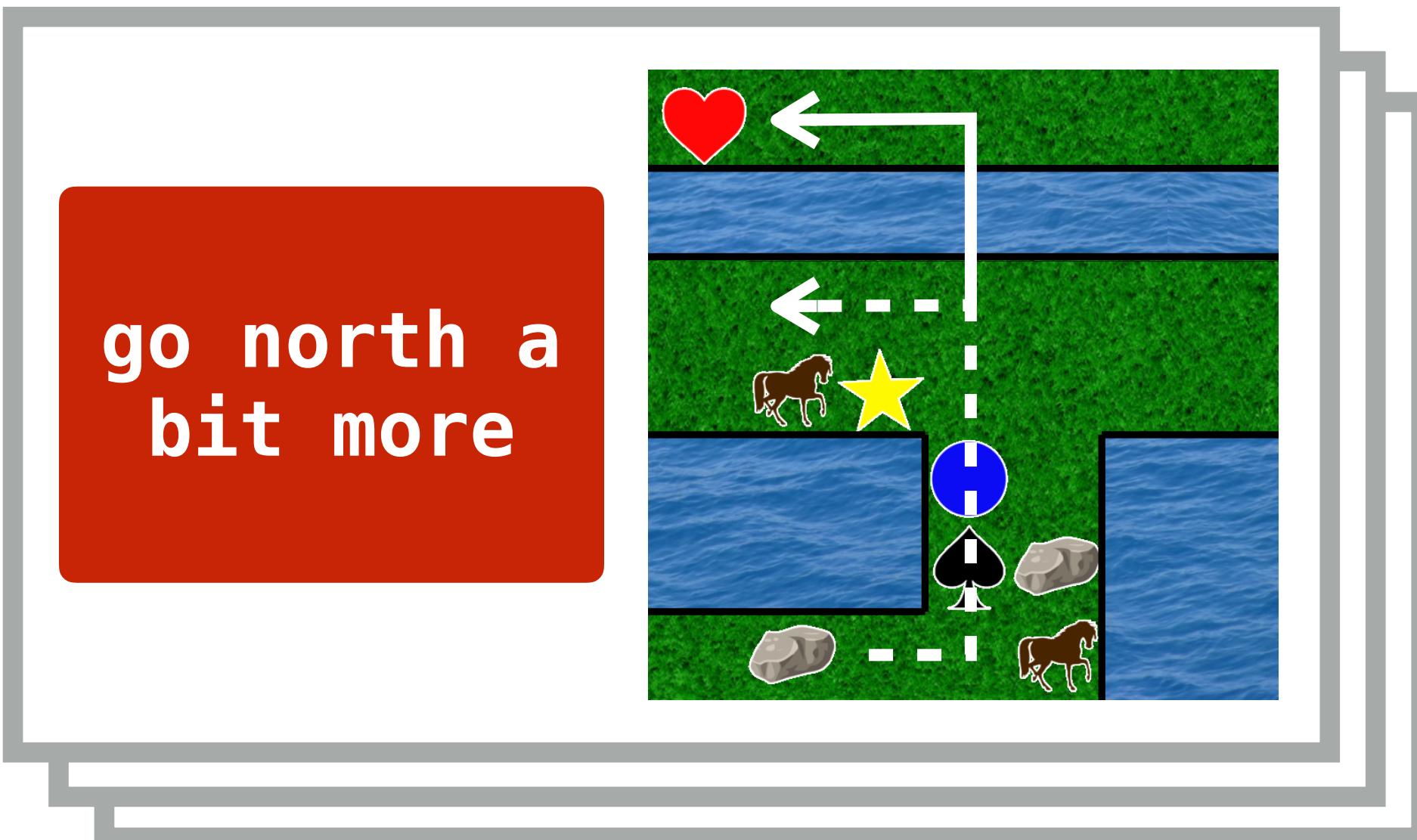
# Generalization

---

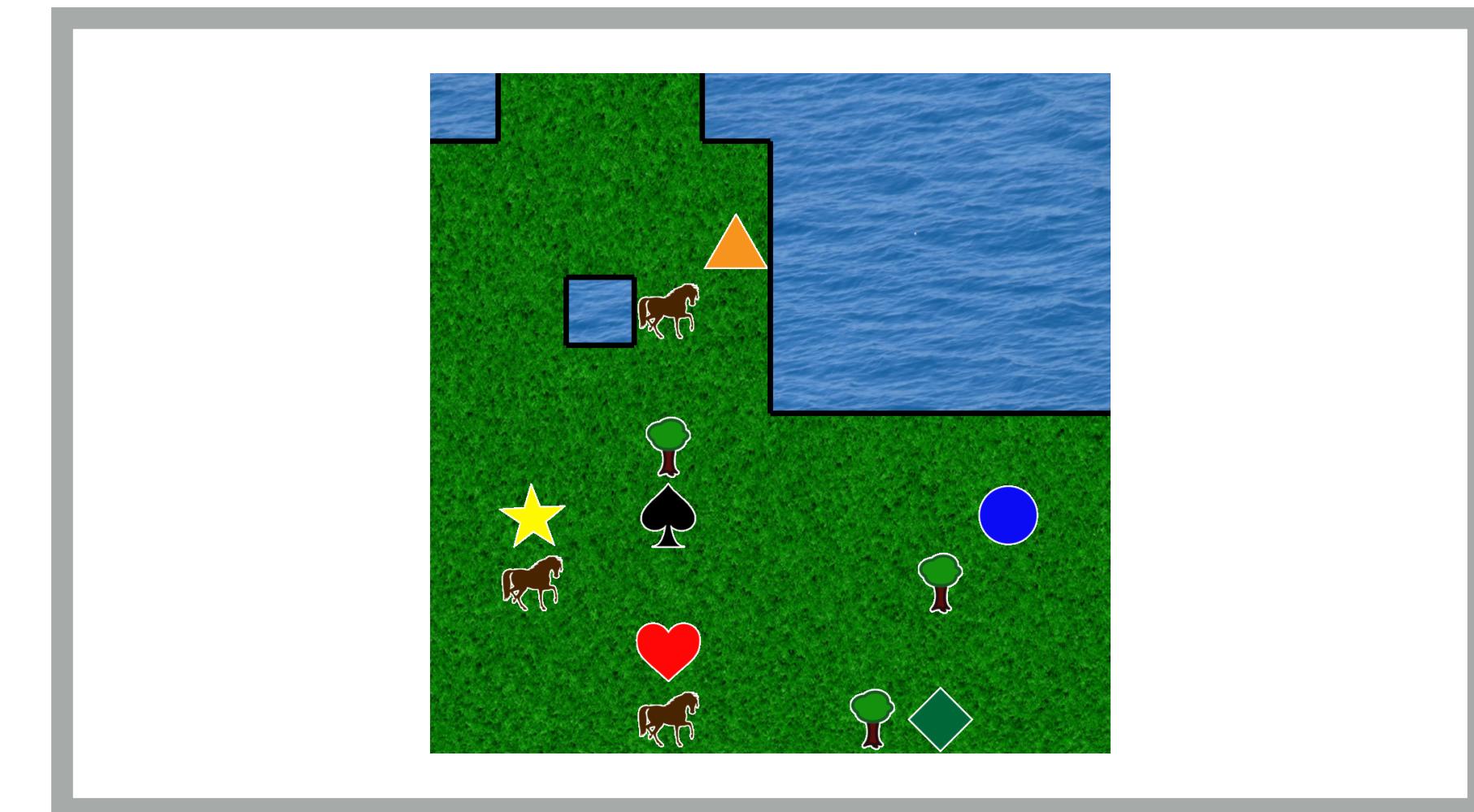


# Learning with corrections

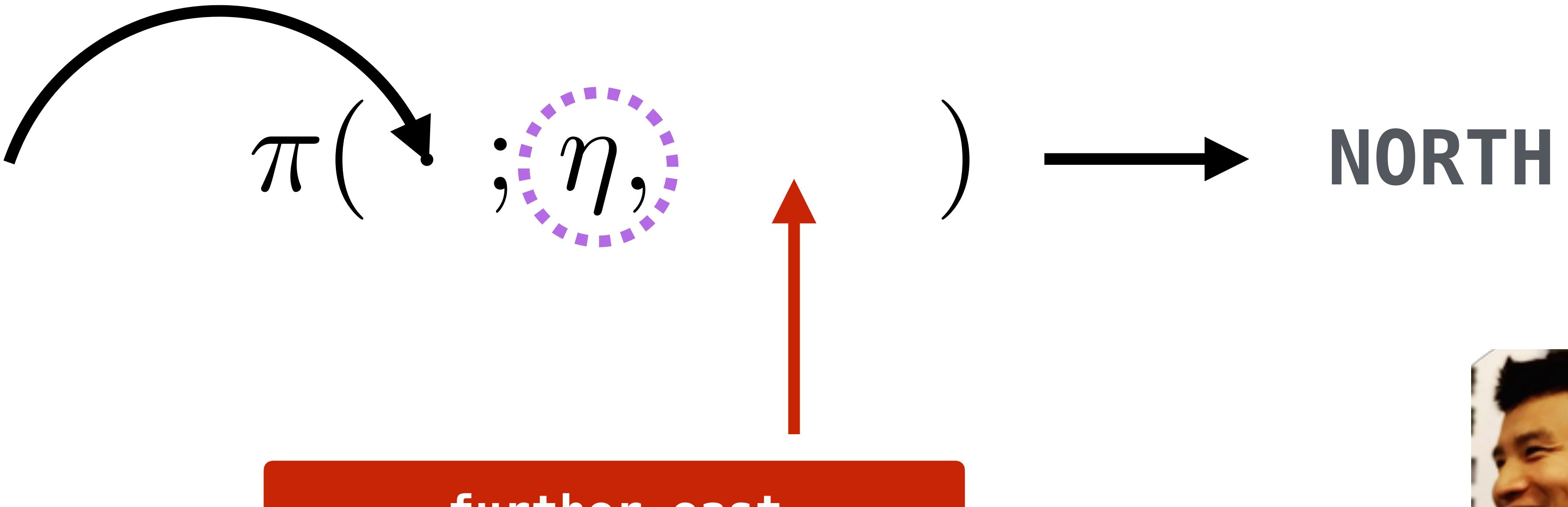
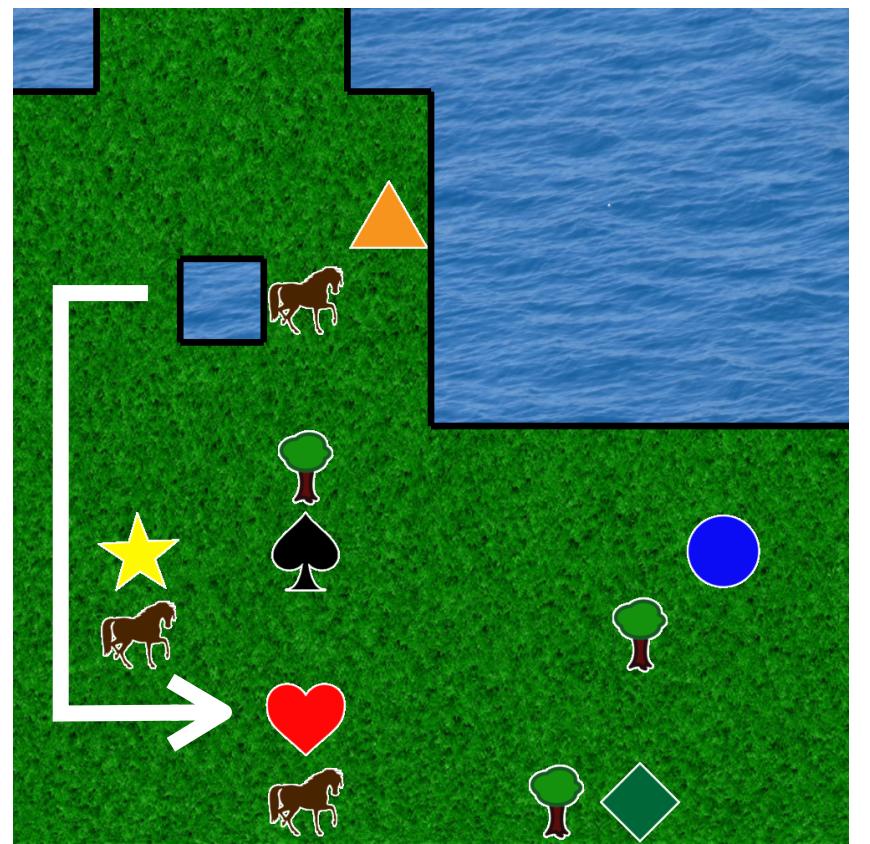
## Language learning



## Reinforcement learning



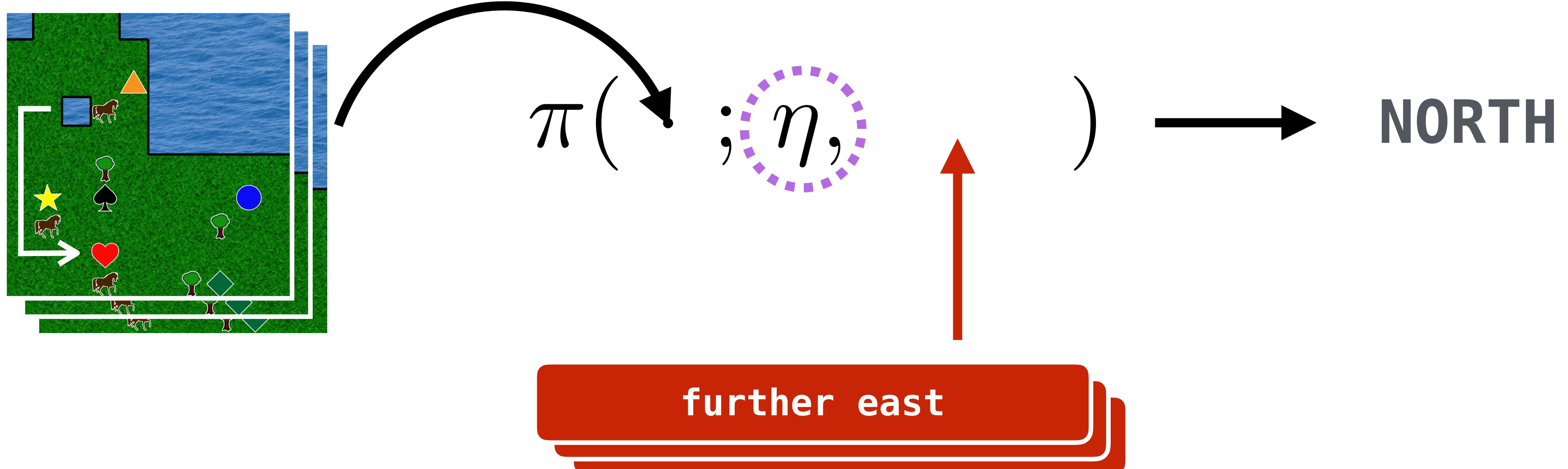
# Pretraining by learning to correct



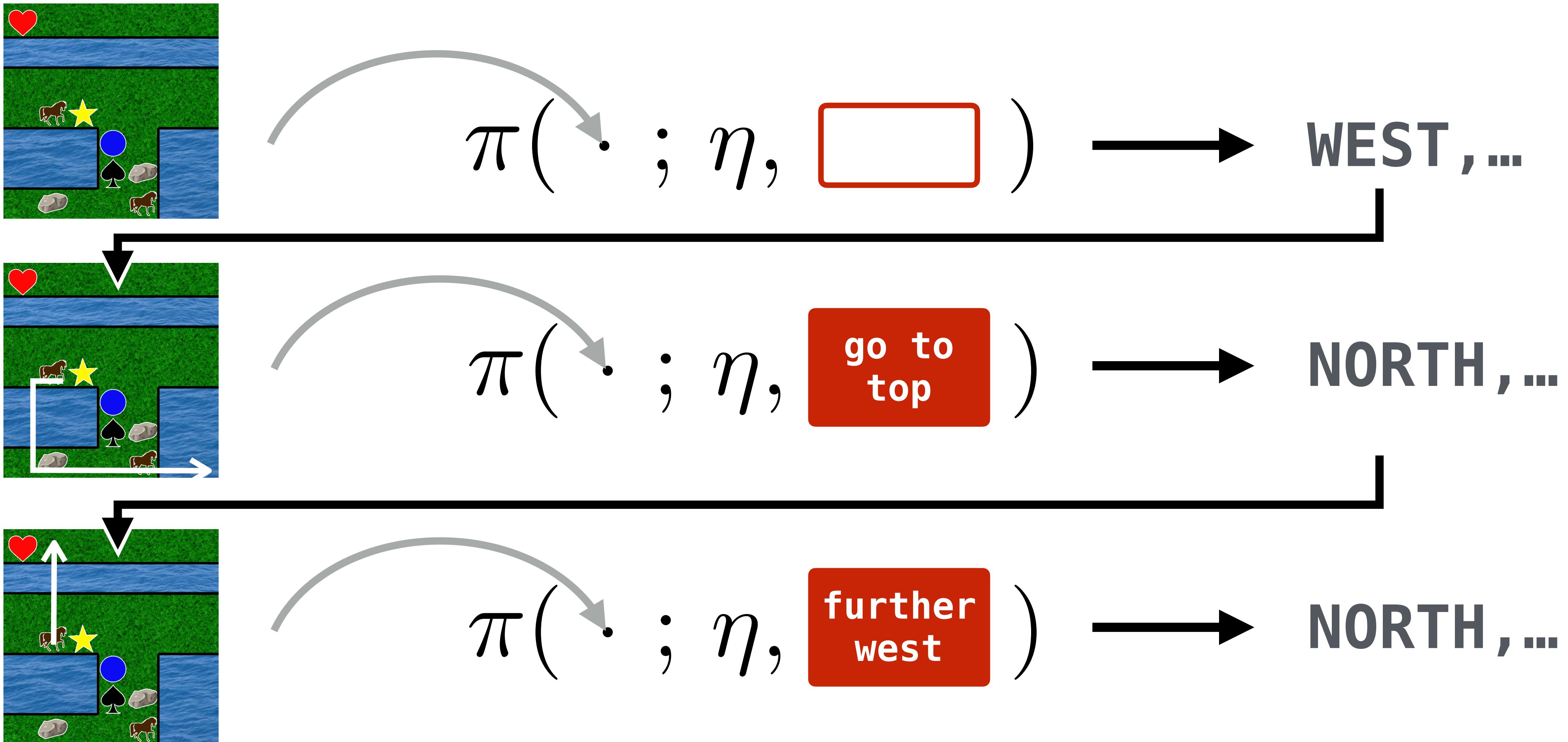
JD Co-Reyes

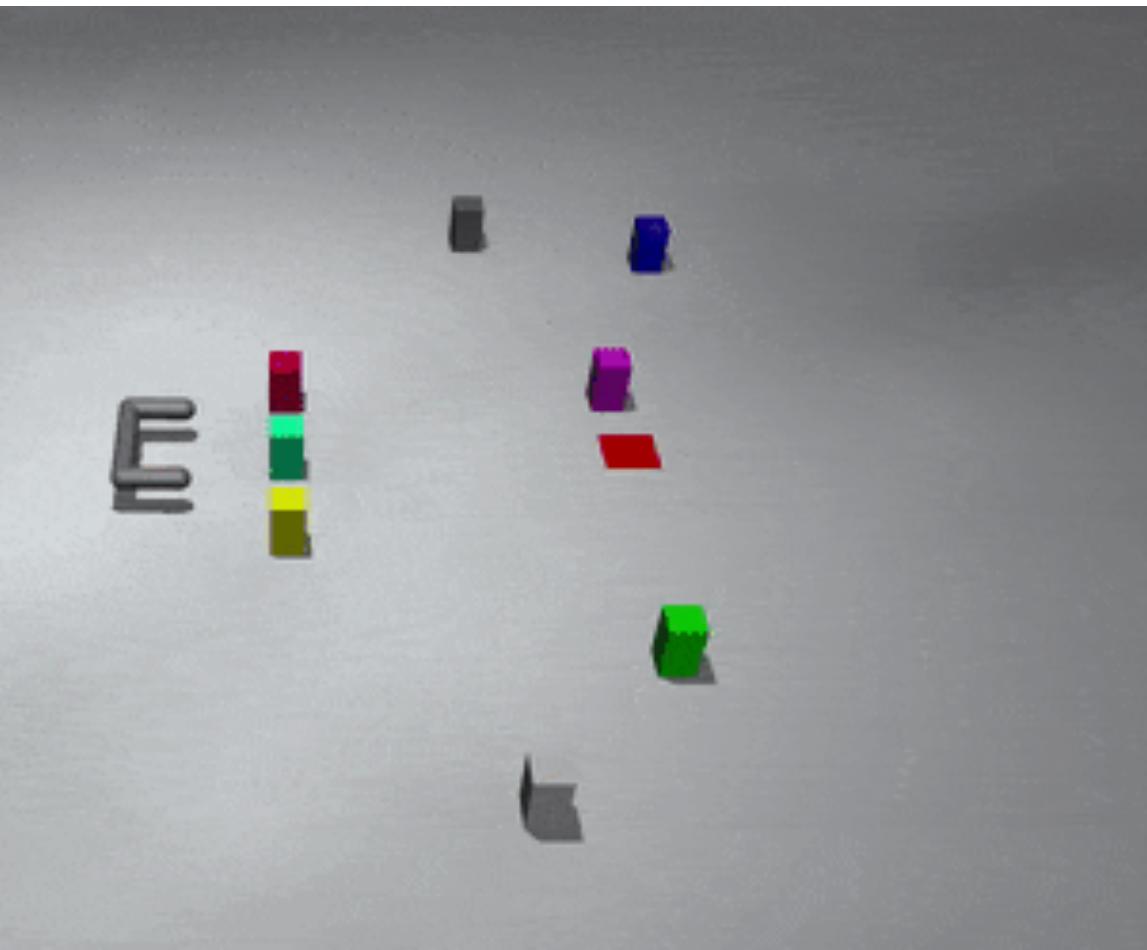
Guiding policies with language via  
meta-learning. ICLR 19.

# Pretraining by learning to correct

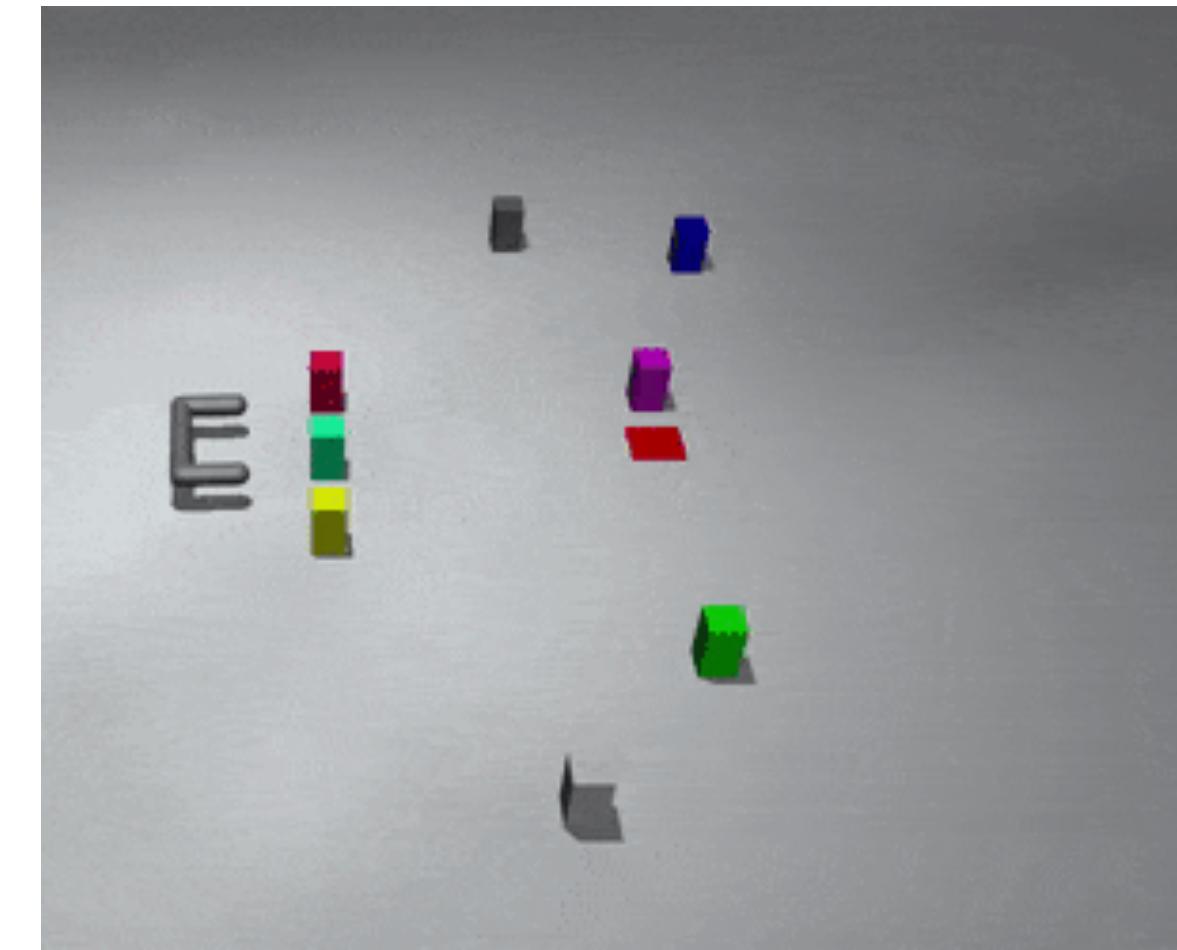


# Learning from corrections

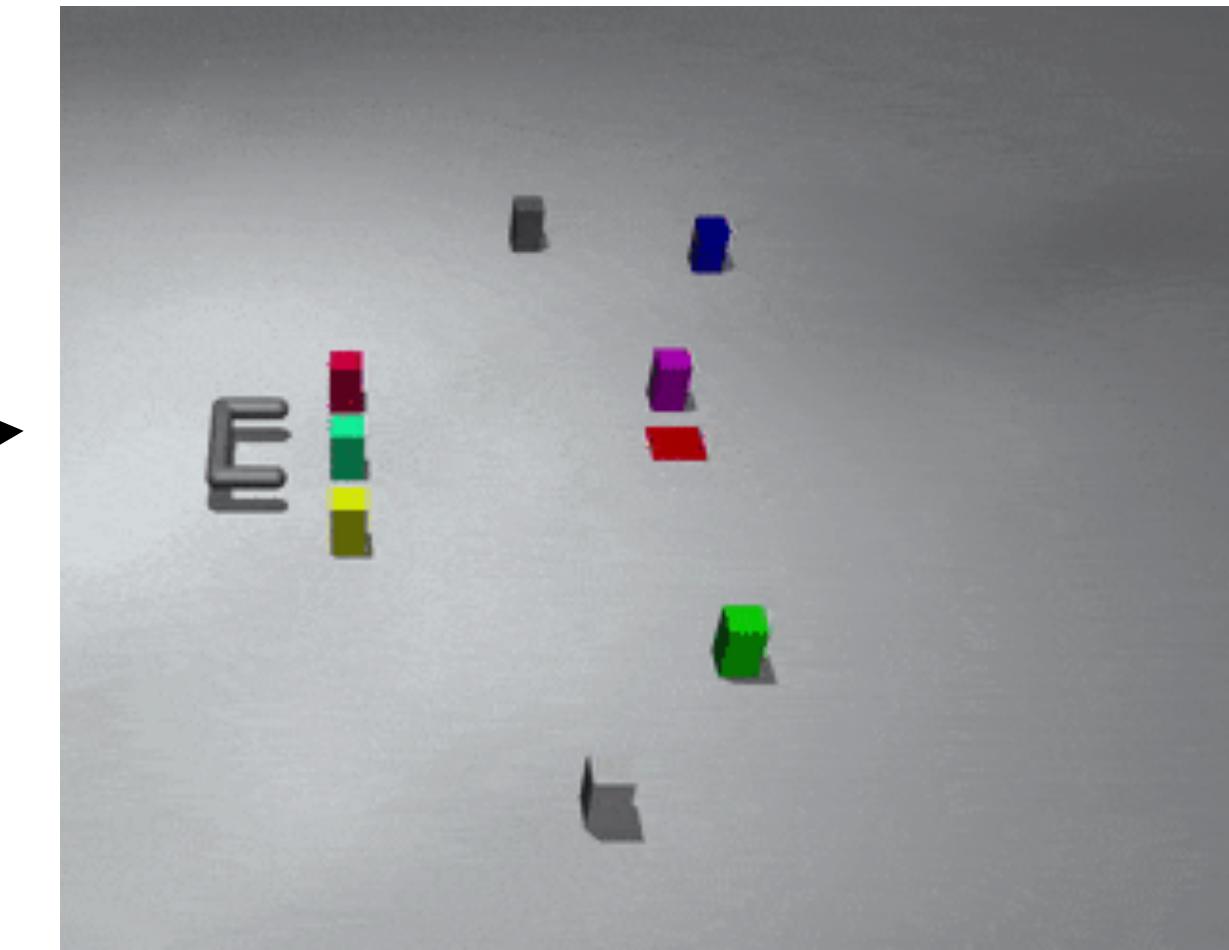




*Touch cyan block.*



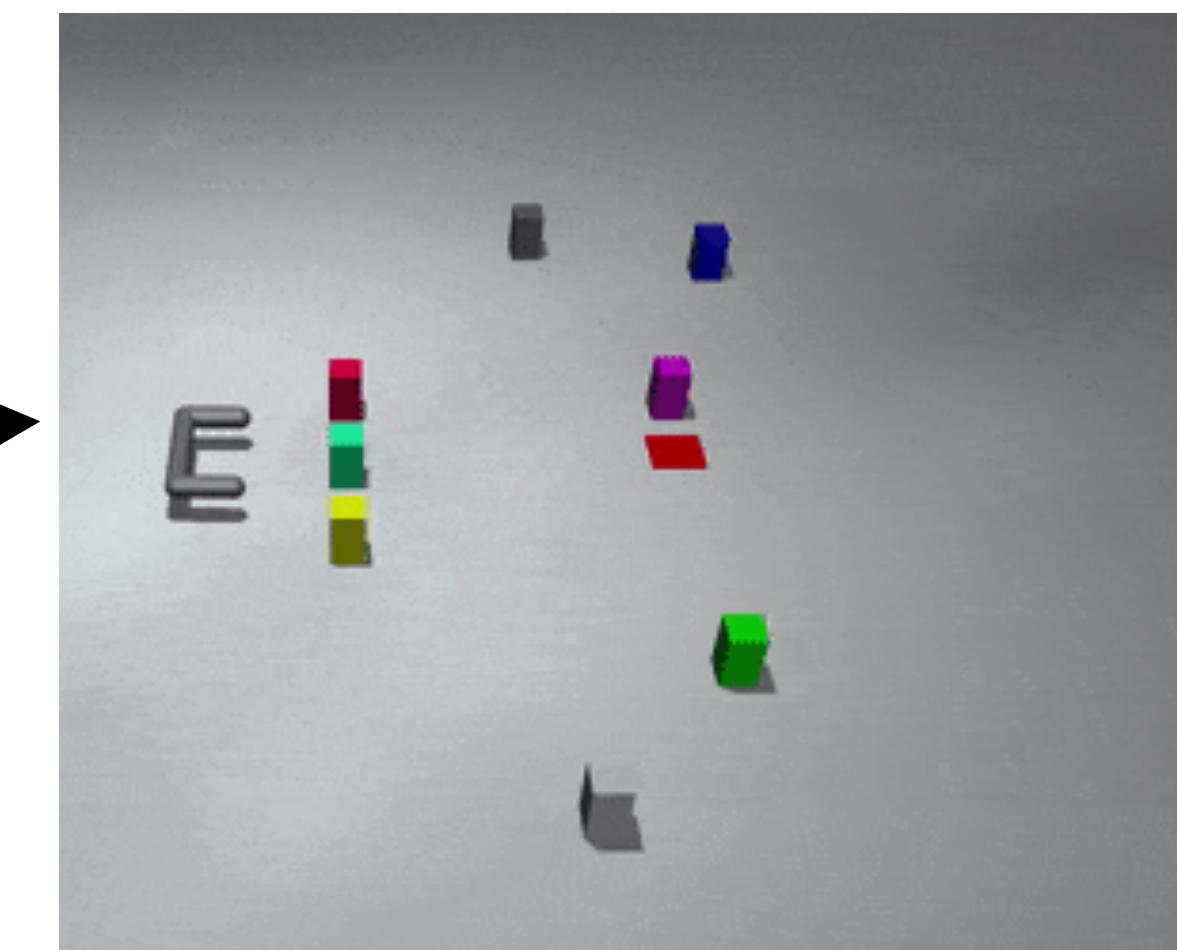
*Move closer to magenta block.*

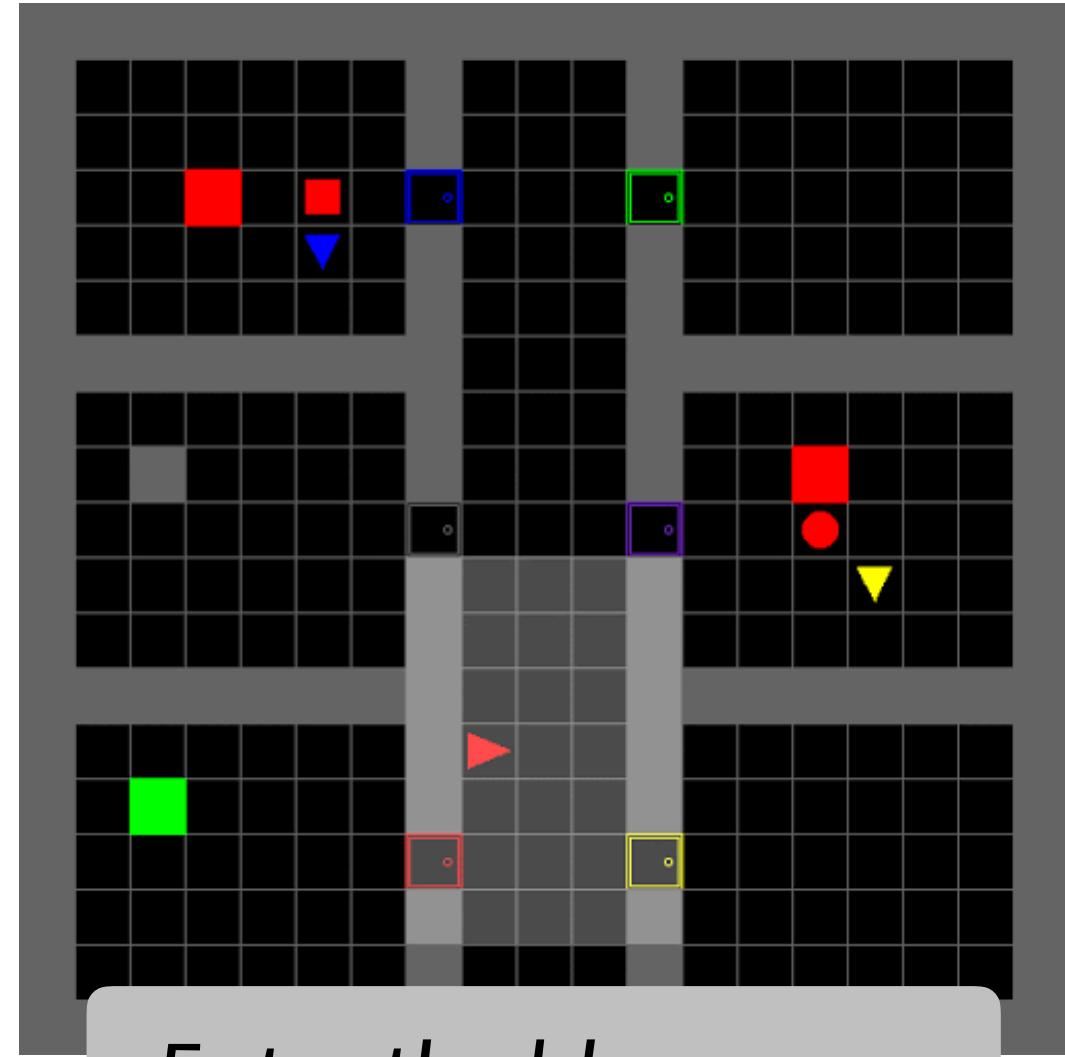


*Move a lot up.*

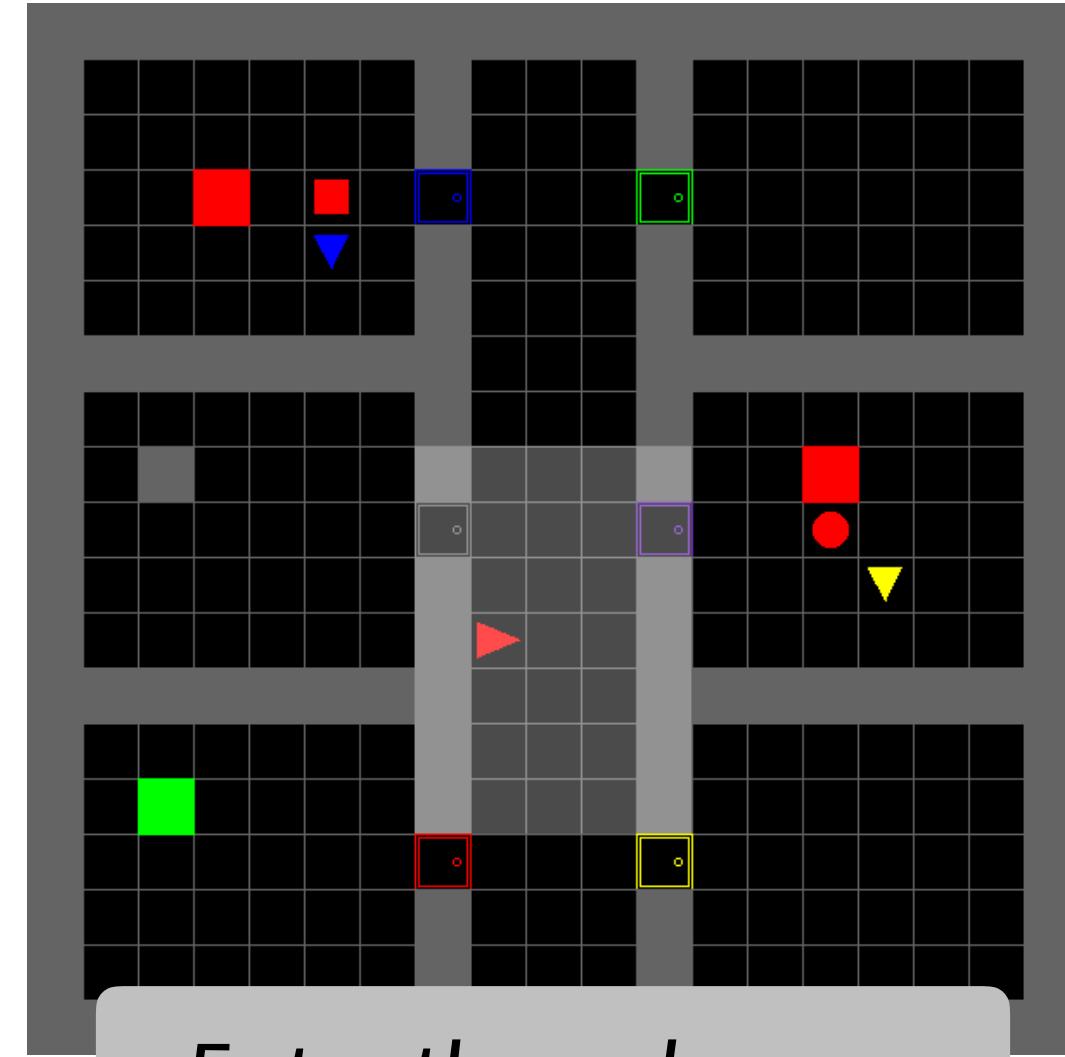


*Move a little up.*

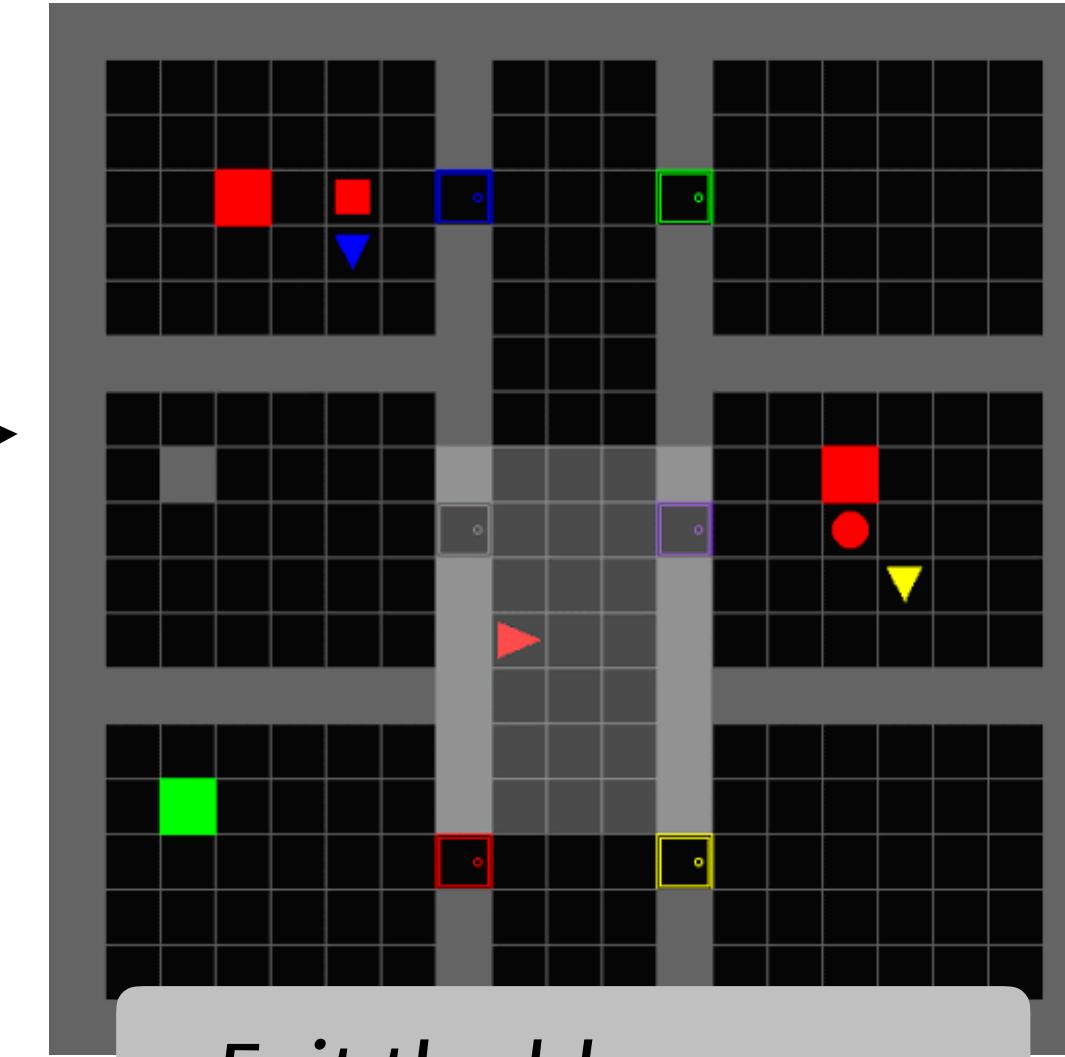




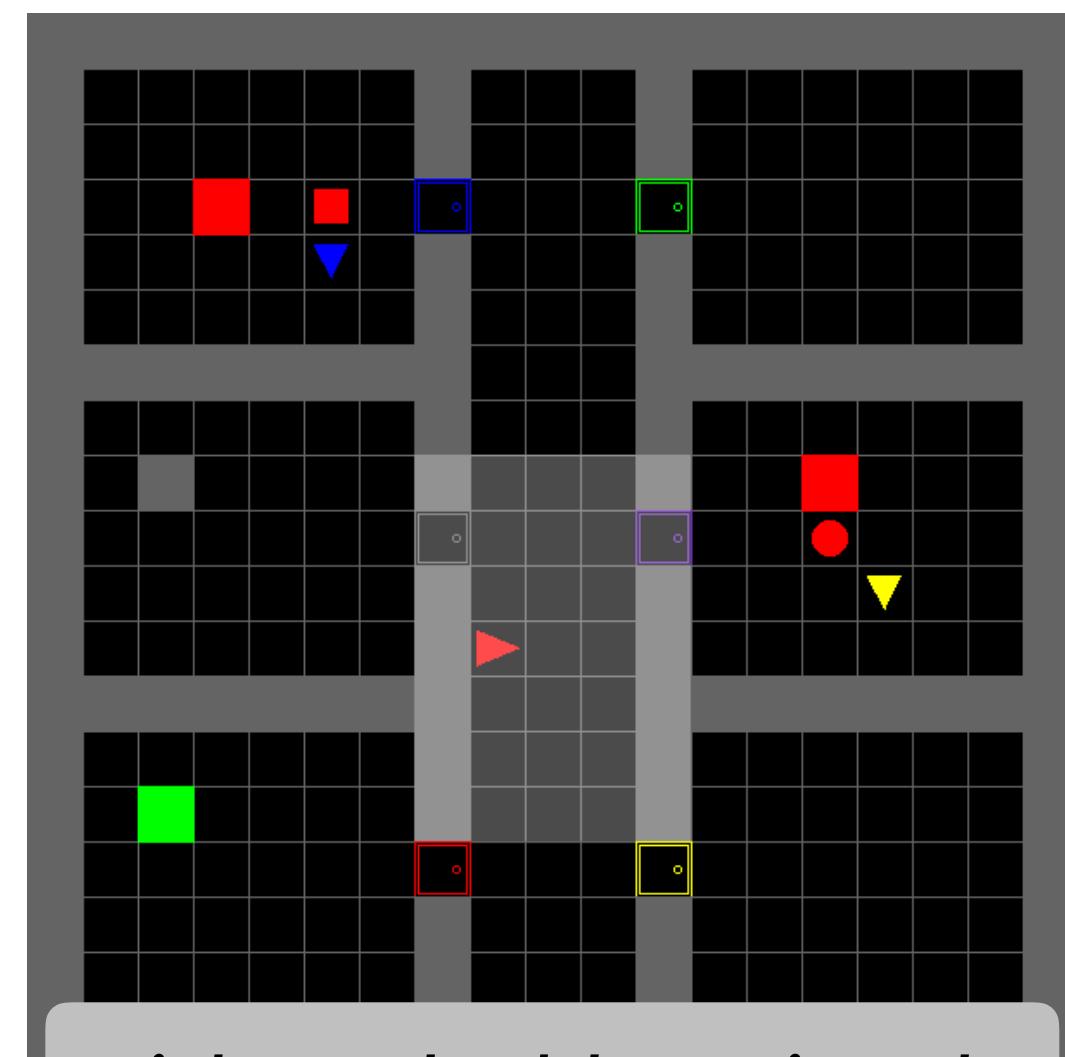
*Enter the blue room.*



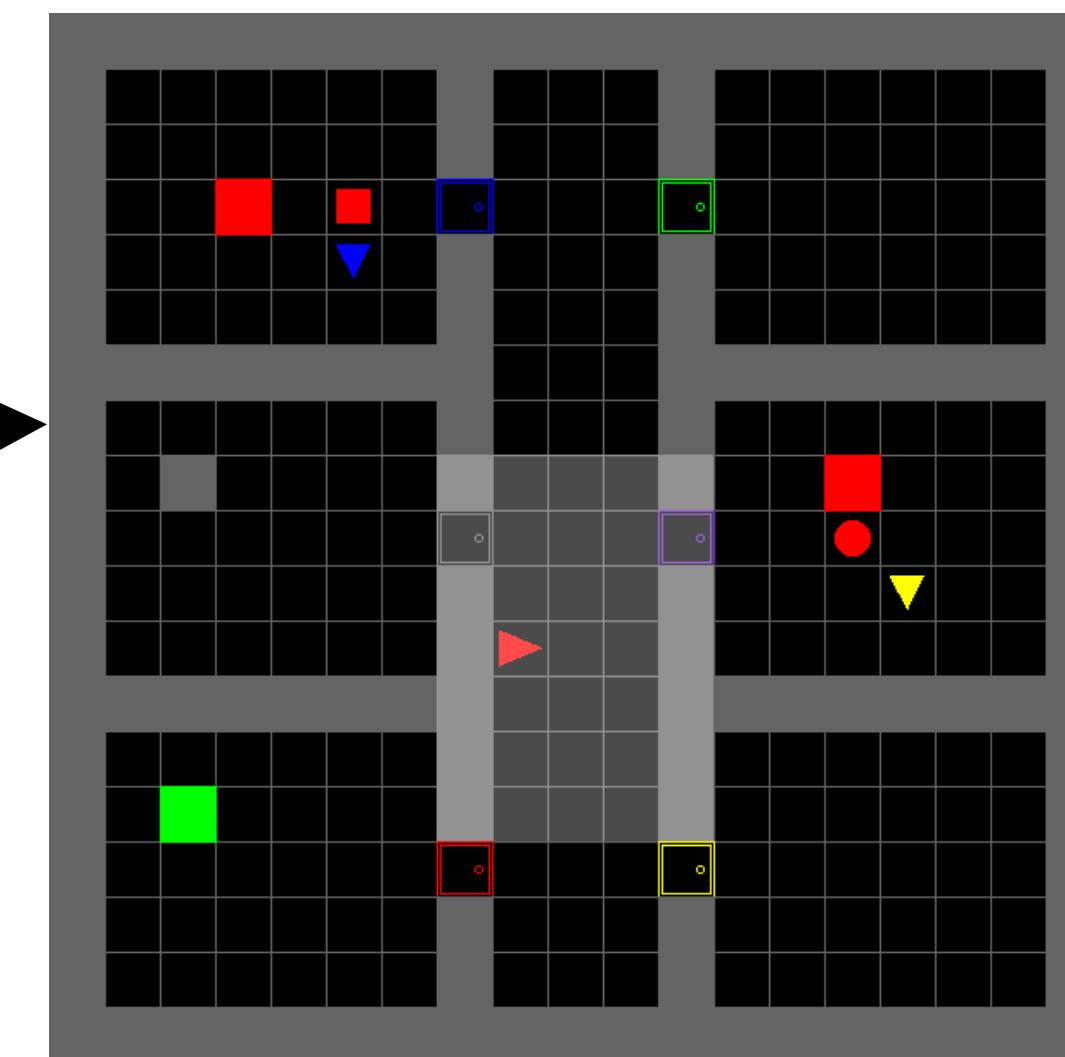
*Enter the red room.*



*Exit the blue room.*



*Pick up the blue triangle*



# Lesson

---

**Language is useful as side information,  
not just a goal specification.**

Use it with / instead of instructions as  
a representational bottleneck  
or interactive advice

So what comes next?

# What comes next?

---

Challenges for the field:

# What comes next?

---

Challenges for the field:

- **huge** datasets

# What comes next?

---

Challenges for the field:

- **huge** datasets
- with **fake** annotations

# What comes next?

---

Challenges for the field:

- **huge** datasets
- with **fake** annotations
- that look **very little like natural language**

# What comes next?

---

Challenges for the field:

- **huge** datasets → Learn to make do without an annotation for every rollout!
- with **fake** annotations
- that look **very little like natural language**

# What comes next?

---

Challenges for the field:

- **huge datasets** → Learn to make do without an annotation for every rollout!
- **with fake annotations** → Learn to generalize from fake strings to real ones!
- **that look very little like natural language**

# What comes next?

---

Challenges for the field:

- **huge datasets** → Learn to make do without an annotation for every rollout!
- **with fake annotations** → Learn to generalize from fake strings to real ones!
- **that look very little like natural language**  
→ Pay attention to human evals (or scope claims accordingly)!

# Learn more: Luketina et al.,

## *A survey of reinforcement learning informed by natural language*

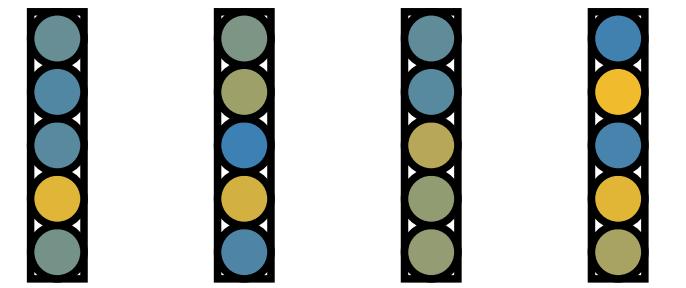
<https://arxiv.org/abs/1906.03926>

### Task-independent

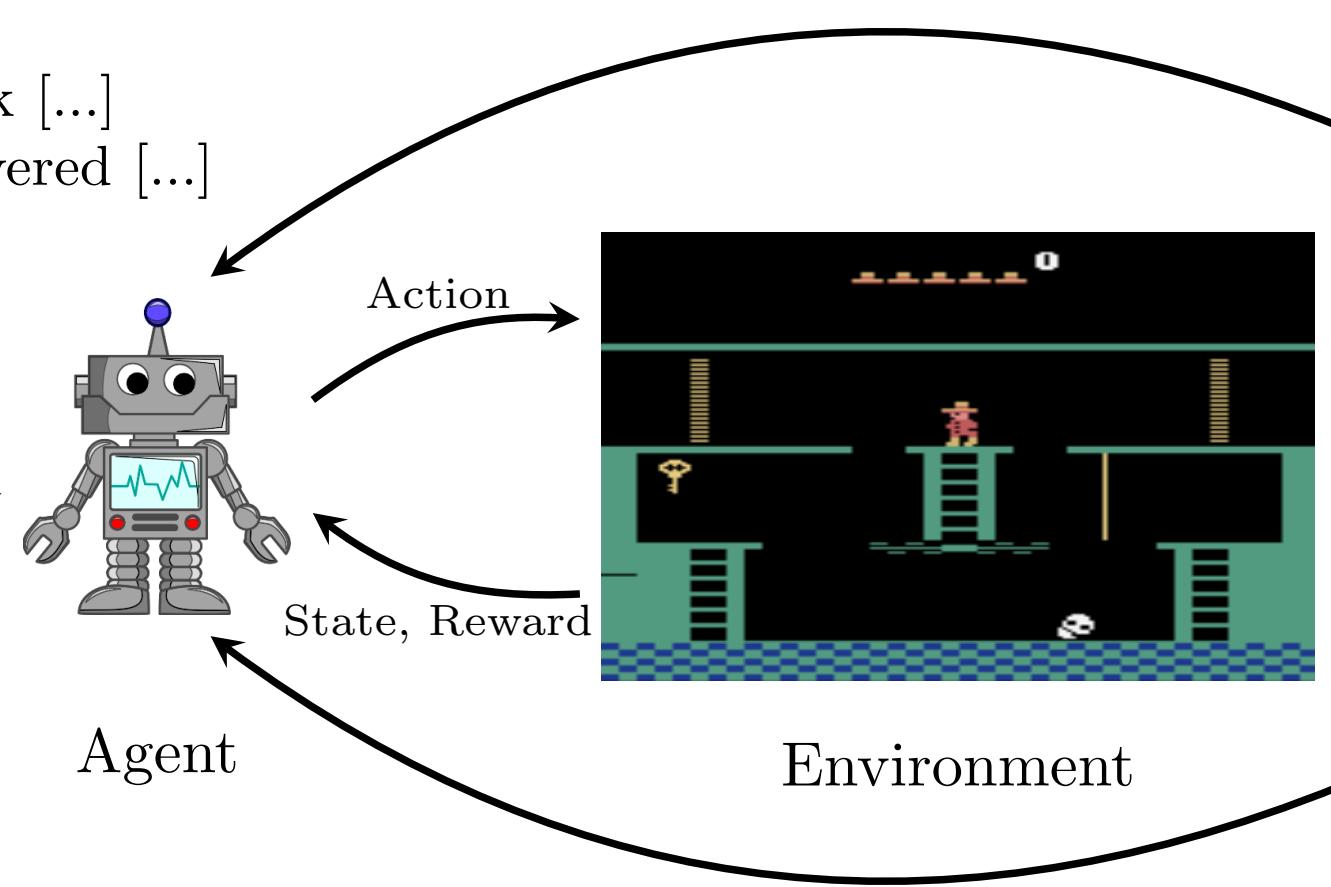
[...] having the correct  
[...] known lock and  
[...] unless the correct

**key** can open the lock [...]  
**key** device was discovered [...]  
**key** is inserted [...]

Pre-training



$v_{\text{key}}$   $v_{\text{skull}}$   $v_{\text{ladder}}$   $v_{\text{rope}}$



### Task-dependent

#### **Language-assisted**

**Key** Opens a door of the same color as the key.

**Skull** They come in two varieties, rolling skulls and bouncing skulls ... you must jump over rolling skulls and walk under bouncing skulls.

#### **Language-conditional**

Go down the ladder and walk right immediately to avoid falling off the conveyor belt, jump to the yellow rope and again to the platform on the right.