



# Unsupervised Reinforcement Learning

Pieter Abbeel  
UC Berkeley EECS

# Outline

---

- NN Architectures
- Representation Learning for Reinforcement Learning
- Active Pre-Training for Reinforcement Learning
- Human-in-the-loop Reinforcement Learning

# Outline

---

- *NN Architectures*
- Representation Learning for Reinforcement Learning
- Active Pre-Training for Reinforcement Learning
- Human-in-the-loop Reinforcement Learning

# Language Models (e.g. OpenAI's GPT-2)

Prompt: A train carriage containing controlled nuclear materials was stolen in Cincinnati today. Its whereabouts are unknown.

GPT-2: *The incident occurred on the downtown train line, which runs from Covington and Ashland stations. In an email to Ohio news outlets, the U.S. Department of Energy said it is working with the Federal Railroad Administration to find the thief.*

*'The theft of this nuclear material will have significant negative consequences on public and environmental health, our workforce and the economy of our nation,' said Tom Hicks, the U.S. Energy Secretary, in a statement. 'Our top priority is to secure the theft and ensure it doesn't happen again.'*

*The stolen material was taken from the University of Cincinnati's Research Triangle Park nuclear research site, according to a news release from Department officials.*

*The Nuclear Regulatory Commission did not immediately release any information. According to the release, the U.S. Department of Energy's Office of Nuclear Material Safety and Security is leading that team's investigation.*

*'The safety of people, the environment and the nation's nuclear stockpile is our highest priority,' Hicks said. 'We will get to the bottom of this and make no excuses.'*

# Text Generation I's GPT-2)

SYSTEM PROMPT  
(HUMAN-WRITTEN)

In a shocking find, a team of scientists has discovered a previously unknown species of lizard in a remote, previously unexplored area of the Amazon rainforest. The discovery is surprising to the researchers, who had been searching for new species for years.



MODEL COMPLETION  
(MACHINE-WRITTEN,  
10 TRIES)

The scientist named the new species Ovid's Unicorn. They had previously unknown

Now, after almost two years of research, the odd phenomenon is finally explained.

Dr. Jorge Pérez, an environmental scientist from La Paz, and several colleagues were climbing the Andes Mountains when they found a small, white, horned creature. Unlike other humans, Pérez noticed that the creature was not afraid of them. It was a natural fountain, surrounded by snow.

Pérez and the others then decided to climb the mountain. "It took us about three hours to reach the top, and we were exhausted by the time we reached the top. But it was worth it, because we found some crystals on top,"

Pérez and his friends were amazed by the creature. These creatures could be seen only at night, and they were too much to see them – they were covered in thick fur and had large horns.

AI text writing technology AI text writing technology

The Drum - 17 Feb 2019

This technology could 'absolutely devastate' the internet as we know it

NEWS.com.au - 17 Feb 2019

This AI is so good at writing that its creators won't let you use it

In-Depth - CNN - 18 Feb 2019

Lord of The Rings, Celebrity Gossip: This AI is So Good at Writing That ...

In-Depth - News18 - 18 Feb 2019

View all



When Is Technology Too Dangerous to Release to the Public?

Slate Magazine - 22 Feb 2019

If your knowledge of the model, called GPT-2, came solely on headlines ... U.K.

read, "Elon Musk-Founded OpenAI Builds Artificial Intelligence So ... had trained a

language model using text from 8 million webpages to predict ...

AI Weekly: Experts say OpenAI's controversial model is a potential ...

In-Depth - VentureBeat - 22 Feb 2019

View all



OpenAI's Text Model so Disruptive it's Deemed Too Dangerous To ...

Computer Business Review - 15 Feb 2019

OpenAI's Text Model so Disruptive it's Deemed Too Dangerous To Release ...

OpenAI has declined to release the full research due to concerns over ... We've

trained an unsupervised language model that can generate ...

New AI fake text generator may be too dangerous to release, say ...

Highly Cited - The Guardian - 14 Feb 2019

View all

bizarre creatures the scientists discovered in the Amazon rainforest. For example, they spoke some fairly regular English. Pérez said, "It's like a dialect or dialectic."

The scientists believe that the unicorns may have originated in South America. Unicorns were believed to be descendants of a creature that lived there before the arrival of humans in South America.

It is currently unclear, some believe that perhaps the first time a human and a unicorn met each other was during the Stone Age civilization. According to Pérez, such sightings seem to be quite common."

He added, "It is likely that the only reason we have never seen a unicorn is that they are indeed the descendants of mythical creatures."

"But they seem to be able to communicate with each other, which I believe is a sign of a highly advanced social organization," said the scientist.

# Benchmarks – GPT-2

Dataset	Metric	Our Result	Previous Record	Human
Winograd Schema Challenge	accuracy (+)	<b>70.70%</b>	63.7%	92%+
LAMBADA	accuracy (+)	<b>63.24%</b>	59.23%	95%+
LAMBADA	perplexity (-)	<b>8.6</b>	99	~1-2
Children's Book Test Common Nouns (validation accuracy)	accuracy (+)	<b>93.30%</b>	85.7%	96%
Children's Book Test Named Entities (validation accuracy)	accuracy (+)	<b>89.05%</b>	82.3%	92%
Penn Tree Bank	perplexity (-)	<b>35.76</b>	46.54	unknown
WikiText-2	perplexity (-)	<b>18.34</b>	39.14	unknown

# Benchmarks -- BERT

## GLUE Results

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

### MultiNLI

Premise: Hills and mountains are especially sanctified in Jainism.

Hypothesis: Jainism hates nature.

Label: Contradiction

### CoLa

Sentence: The wagon rumbled down the road.

Label: Acceptable

Sentence: The car honked down the road.

Label: Unacceptable

Might these pre-trained transformers  
be *even* more general?

---

# **Pretrained Transformers As Universal Computation Engines**

---

**Kevin Lu**

UC Berkeley

kz1@berkeley.edu

**Aditya Grover**

Facebook AI Research

adityagrover@fb.com

**Pieter Abbeel**

UC Berkeley

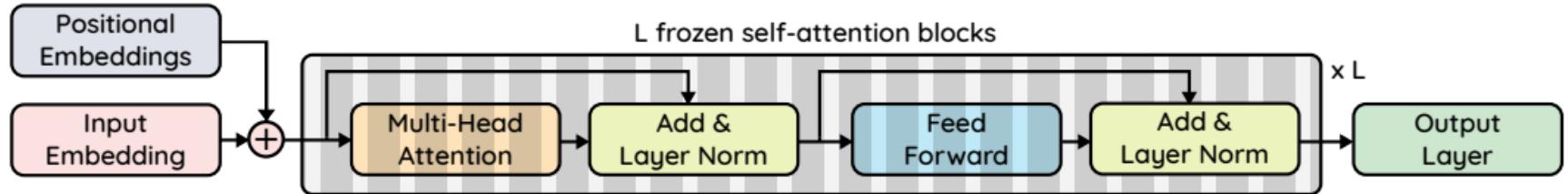
pabbeel@cs.berkeley.edu

**Igor Mordatch**

Google Brain

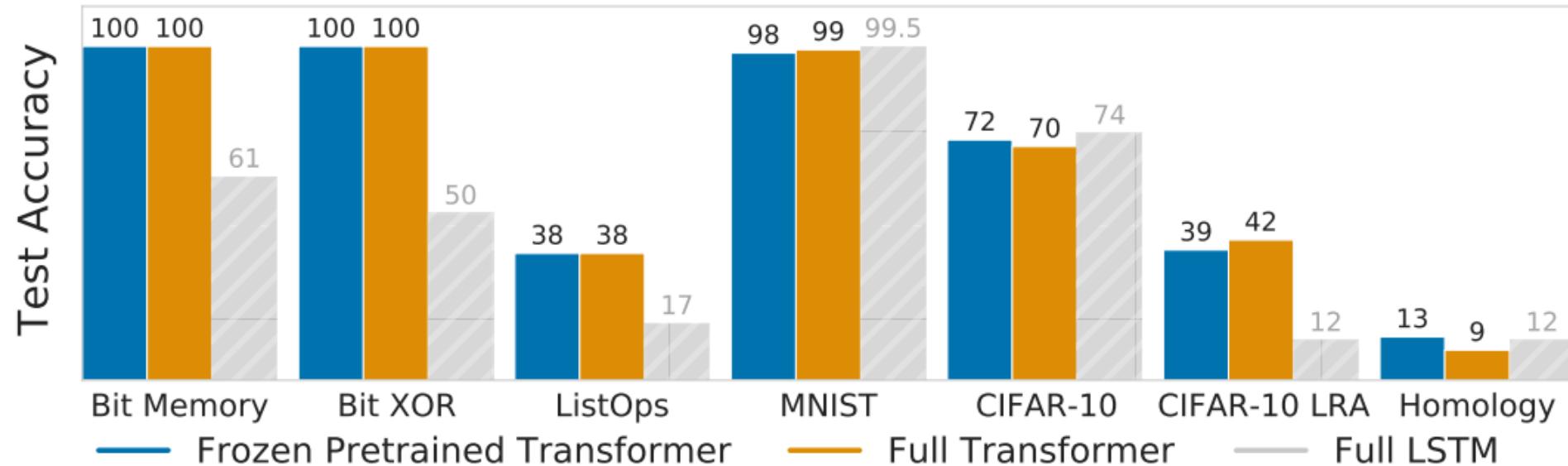
imordatch@google.com

# Pre-Trained Model + .1% finetune



- ***Pre-train:*** language corpus next-token prediction
- ***Minimally fine-tune:***
  - Bit memory
  - Bit XOR
  - ListOps
  - MNIST
  - CIFAR-10 and CIFAR-10 LRA
  - Remote homology detection

# Can pretrained LMs transfer to new modalities?



# What's the importance of the pretraining modality?

<b>Model</b>	<b>Bit Memory</b>	<b>XOR</b>	<b>ListOps</b>	<b>MNIST</b>	<b>C10</b>	<b>C10 LRA</b>	<b>Homology</b>
FPT	100%	100%	38.4%	98.0%	68.2%	38.6%	12.7%
Random	75.8%	100%	34.3%	91.7%	61.7%	36.1%	9.3%
Bit	100%	100%	35.4%	97.8%	62.6%	36.7%	7.8%
ViT	100%	100%	37.4%	97.8%	72.5%	43.0%	7.5%

Table 2: Test accuracy of language-pretrained (FPT) vs randomly initialized (Random) vs Bit Memory pretraining (Bit) vs pretrained Vision Transformer (ViT) models. The transformer is frozen.

# Does performance scale with model size?

<b>Model Size</b>	<b># Layers</b>	<b>Total Params</b>	<b>Trained Params</b>	<b>FPT</b>	<b>Random</b>
Small (Base)	12	117M	106K	68.2%	61.7%
Medium	24	345M	190K	69.8%	64.0%
Large	36	774M	300K	72.1%	65.7%

Table 6: Test accuracy of larger frozen transformer models on CIFAR-10.

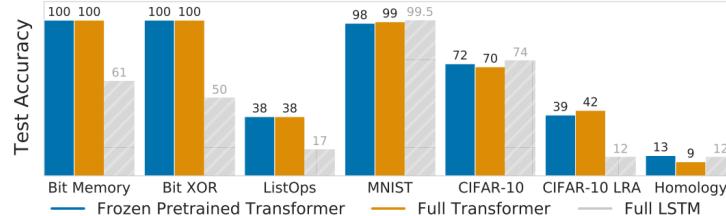
## Does the trend hold across other transformer models?

<b>Task</b>	<b>GPT-2 (FPT Default)</b>	<b>BERT</b>	<b>T5</b>	<b>Longformer</b>
ListOps	38.4%	38.3%	15.4%	17.0%
CIFAR-10	68.2%	68.8%	64.7%	66.8%

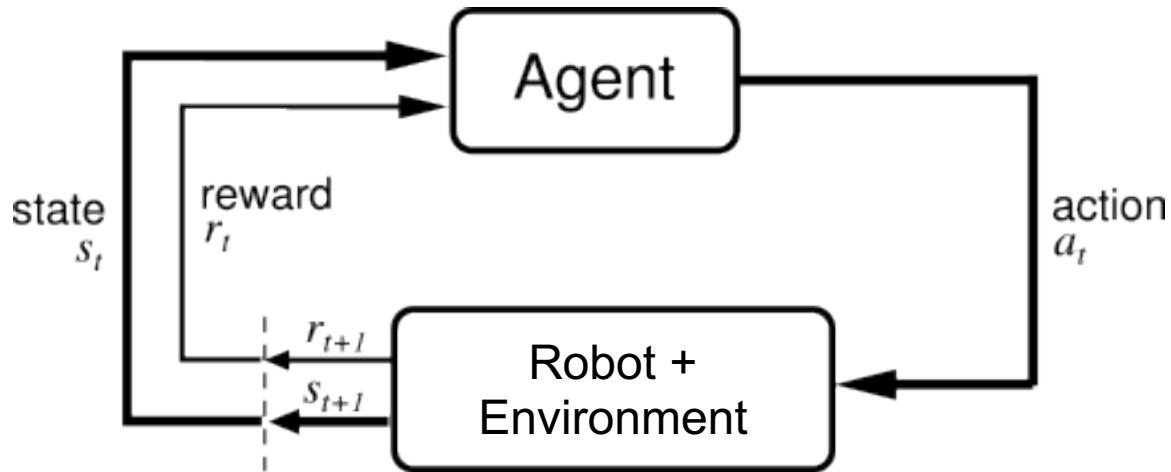
Table 9: Test accuracy for frozen pretrained transformer variants (base model sizes).

# Outline

- NN Architectures: Pretrained Transformers as Universal Computation Engines
- *Representation Learning for Reinforcement Learning*
- Active Pre-Training for Reinforcement Learning
- Human-in-the-loop Reinforcement Learning

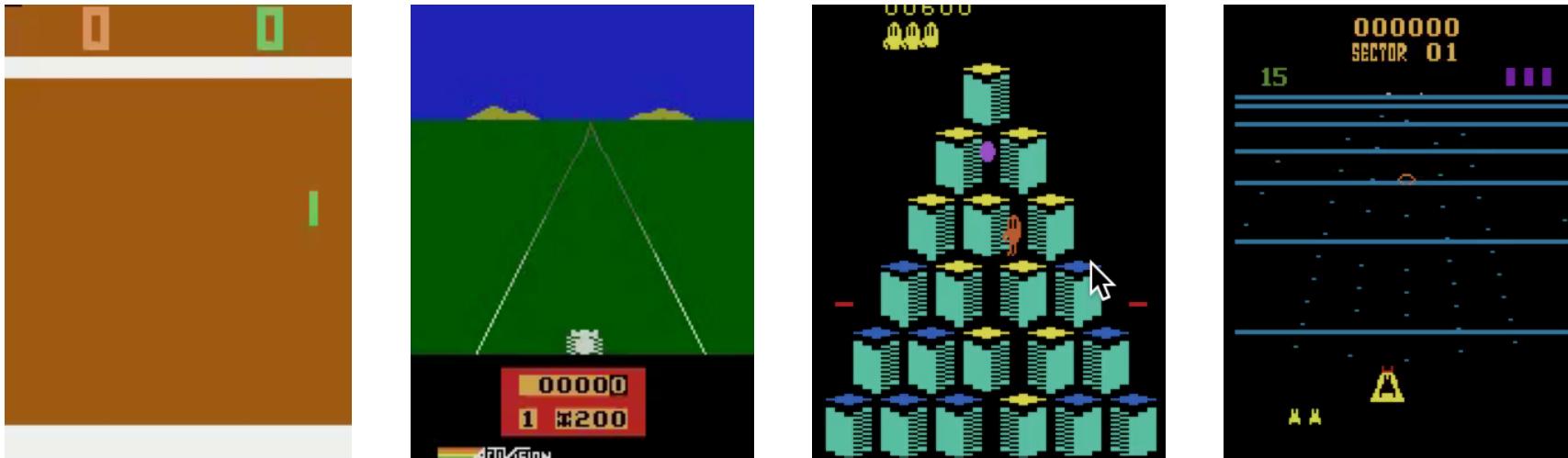


# Reinforcement Learning (RL)



$$\max_{\theta} \mathbb{E} \left[ \sum_{t=0}^H R(s_t) | \pi_{\theta} \right]$$

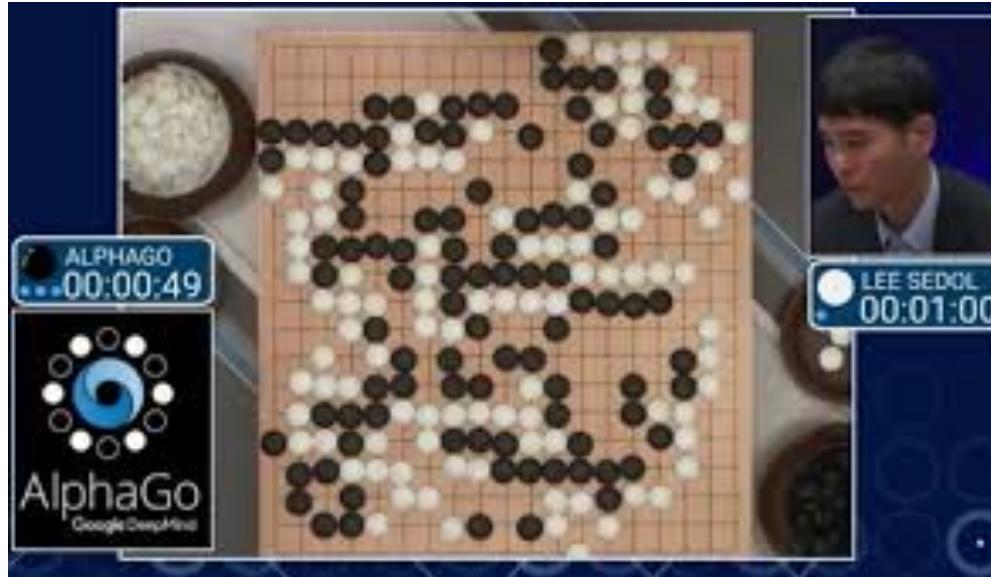
# Deep RL: Atari



DQN Mnih et al, NIPS 2013 / Nature 2015

**MCTS** Guo et al, NIPS 2014; **TRPO** Schulman, Levine, Moritz, Jordan, Abbeel, ICML 2015; **A3C** Mnih et al, ICML 2016; **Dueling DQN** Wang et al ICML 2016; **Double DQN** van Hasselt et al, AAAI 2016; **Prioritized Experience Replay** Schaul et al, ICLR 2016; **Bootstrapped DQN** Osband et al, 2016; **Q-Ensembles** Chen et al, 2017; **Rainbow** Hessel et al, 2017; **Accelerated Stooke and Abbeel, 2018;** ...

# Deep RL: Go



**AlphaGo** Silver et al, Nature 2015

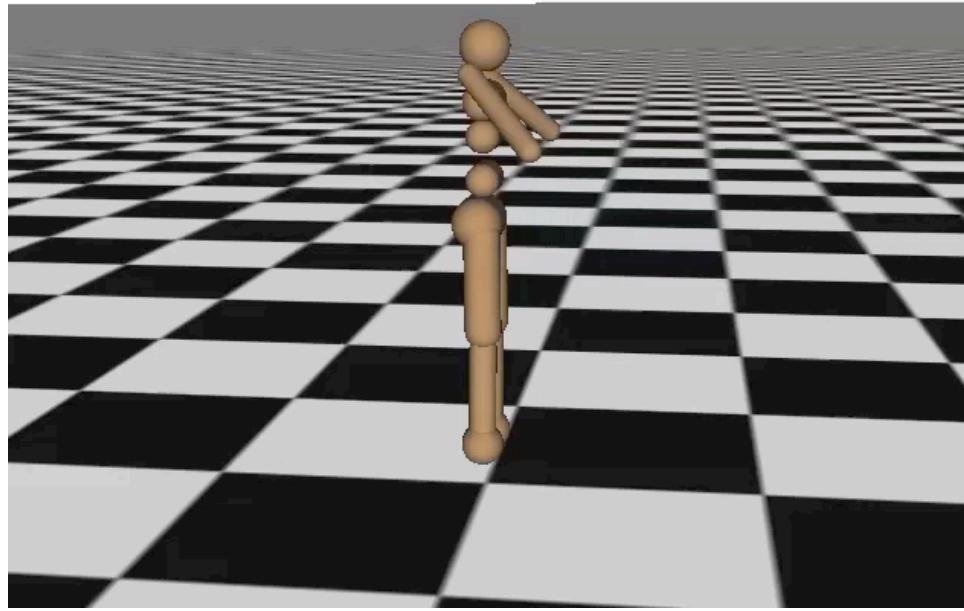
**AlphaGoZero** Silver et al, Nature 2017

**AlphaZero** Silver et al, 2017

Tian et al, 2016; Maddison et al, 2014; Clark et al, 2015

# Deep RL: Robot Locomotion

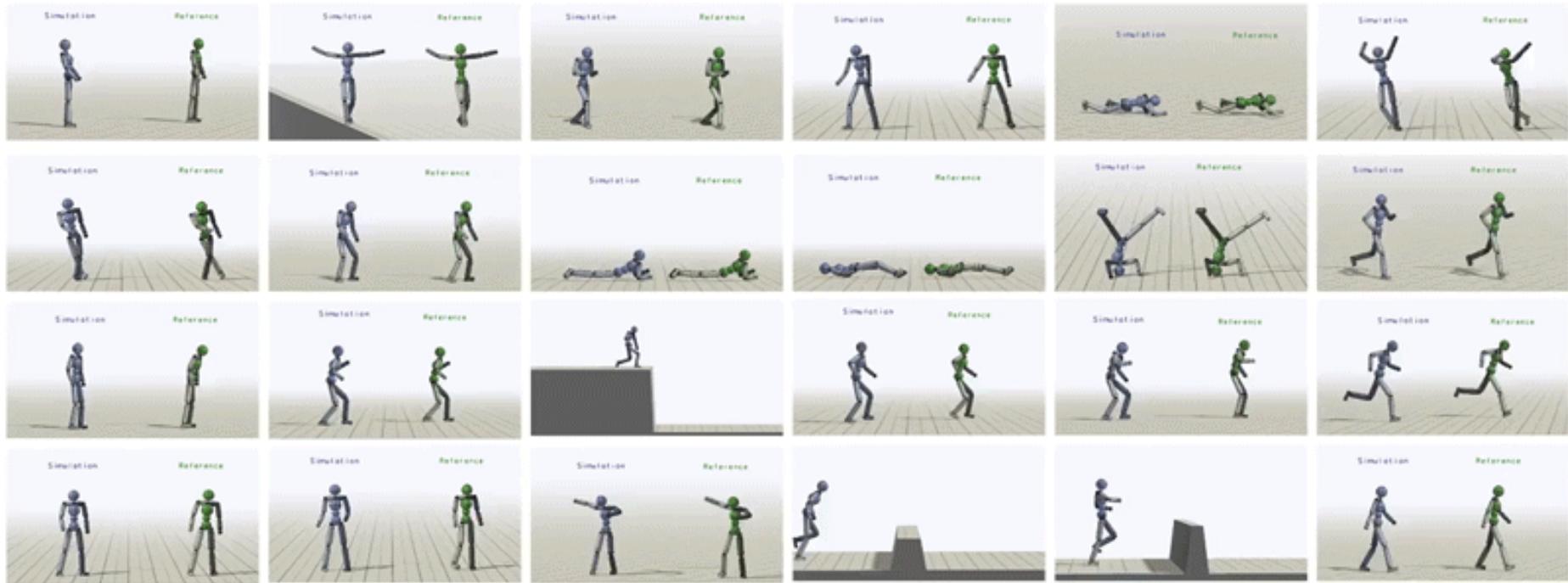
Iteration 0



^ **TRPO** Schulman et al, 2015 + **GAE** Schulman et al, 2016

See also: **DDPG** Lillicrap et al 2015; **SVG** Heess et al, 2015; **Q-Prop** Gu et al, 2016; **Scaling up ES** Salimans et al, 2017; **PPO** Schulman et al, 2017; **Parkour** Heess et al, 2017;

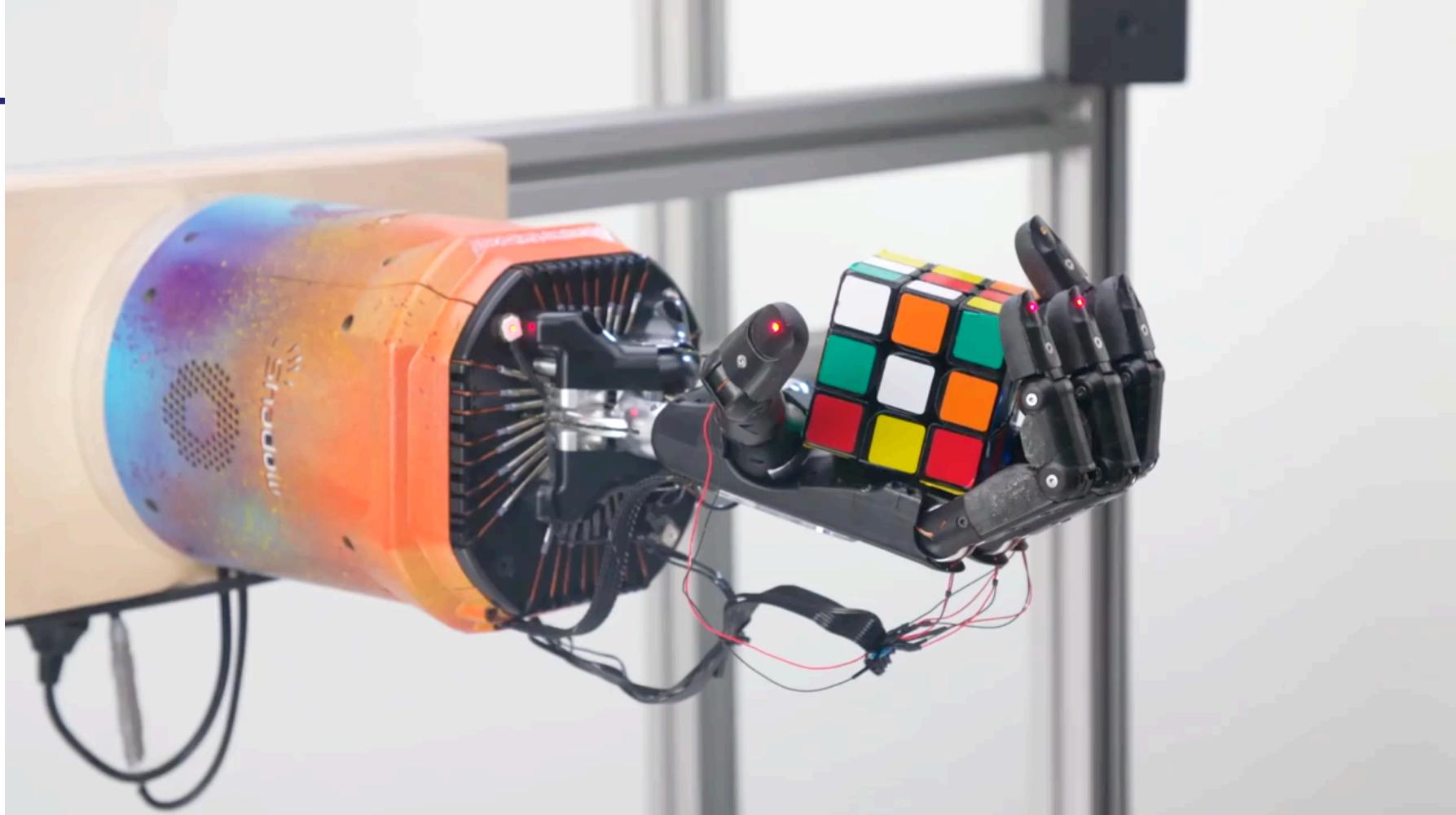
# Deep RL: Robot Locomotion



# BRETT: Berkeley Robot for the Elimination of Tedious Tasks

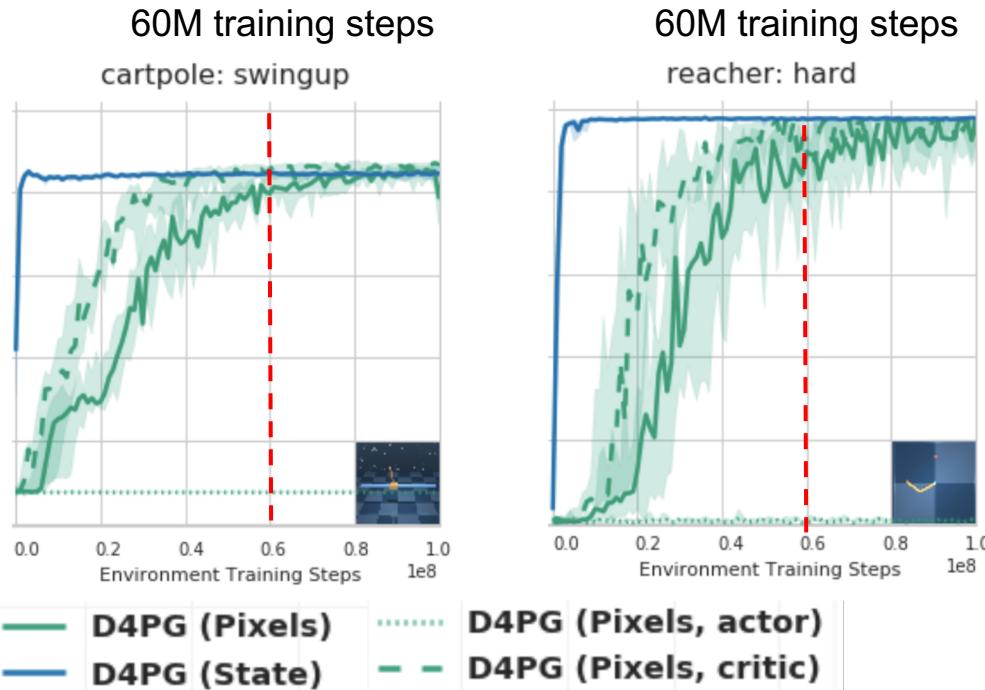
---





# Can visual RL match data-efficiency state RL?

- State-based D4PG (blue) vs pixel-based D4PG (green)



Pixel-based needs > 50M more training steps than state-based to solve same tasks

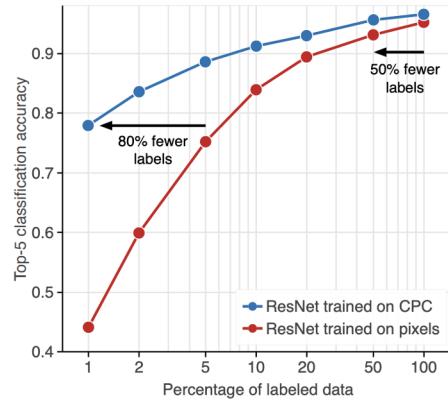
[Tassa et al., 2018] Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D.D.L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A. and Lillicrap, T. [DeepMind Control Suite](#), arxiv:1801.00690, 2018.



LeCake (Yann LeCun)

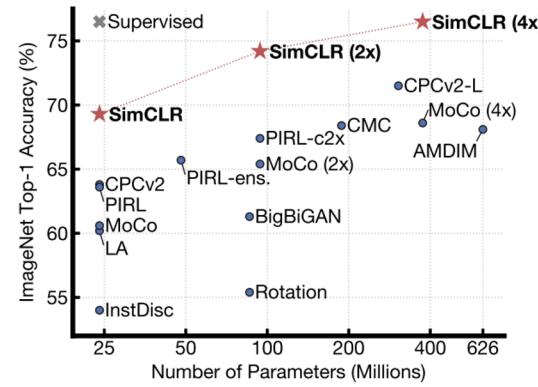
# Contrastive learning: SOTA in computer vision

CPCv2 **top-5** ImageNet accuracy as function of labels



[Henaff, Srinivas et al., 2019]

SimCLR **top-1** ImageNet accuracy as function of # of parameters



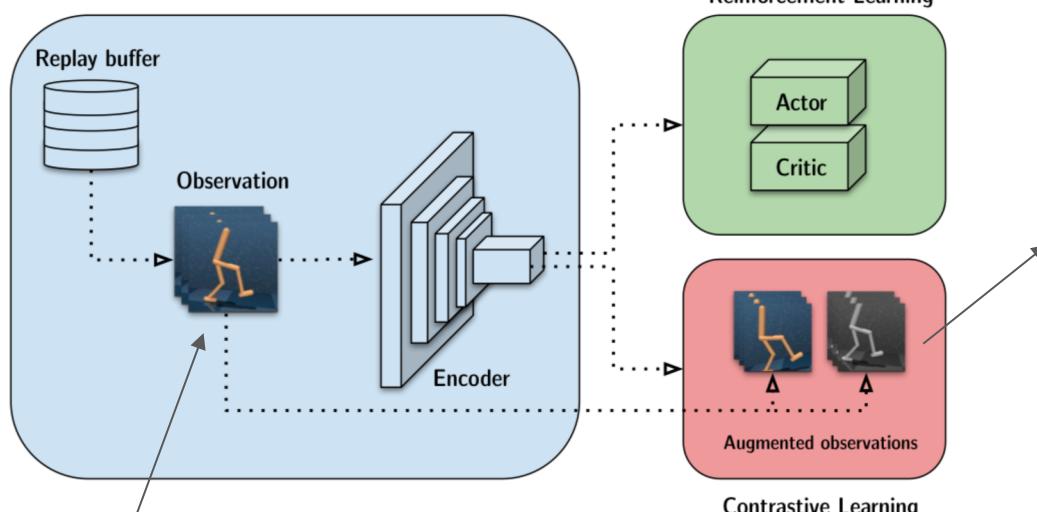
[Chen et al., 2020]

[Henaff et al., 2019] Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, Aaron van den Oord [Data-Efficient Image Recognition with Contrastive Coding](#) arxiv:1905.09272, 2019.

[Chen et al., 2020] Chen, T., Kornblith, S., Norouzi, M. and Hinton, G. [A Simple Framework for Contrastive Learning of Visual Representations](#) arxiv:2002.05709, 2020.

# Contrastive + RL

CURL

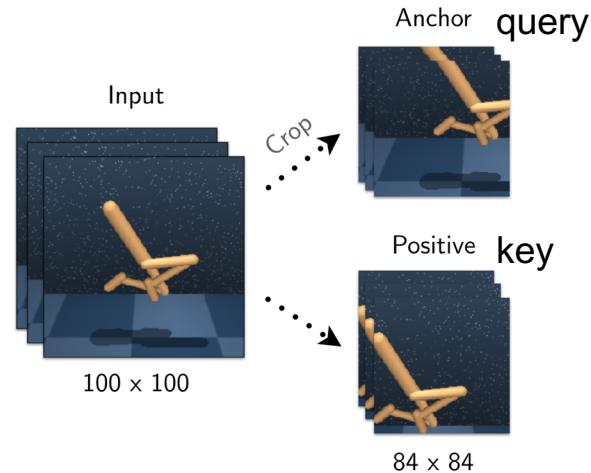


Observations are  
stacked frames

Need to define:

1. query / key pairs
2. similarity measure
3. architecture

# 1. Query / key pairs: random crop



## 2. Bilinear inner product with learned weight matrix

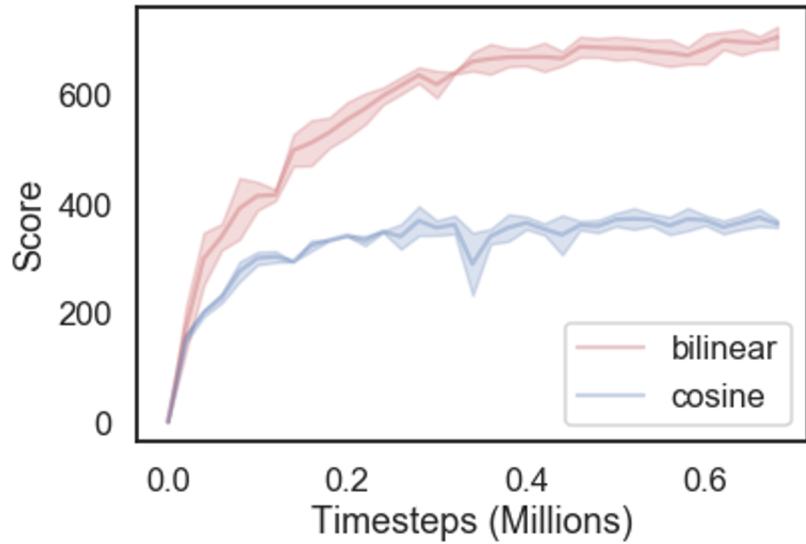
logits

$$\begin{bmatrix} q_0^T W k_0 & q_0^T W k_1 & \dots & q_0^T W k_j \\ q_1^T W k_0 & q_1^T W k_1 & \dots & q_1^T W k_j \\ \vdots & \vdots & \ddots & \vdots \\ q_j^T W k_0 & q_j^T W k_1 & \dots & q_j^T W k_j \end{bmatrix} \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

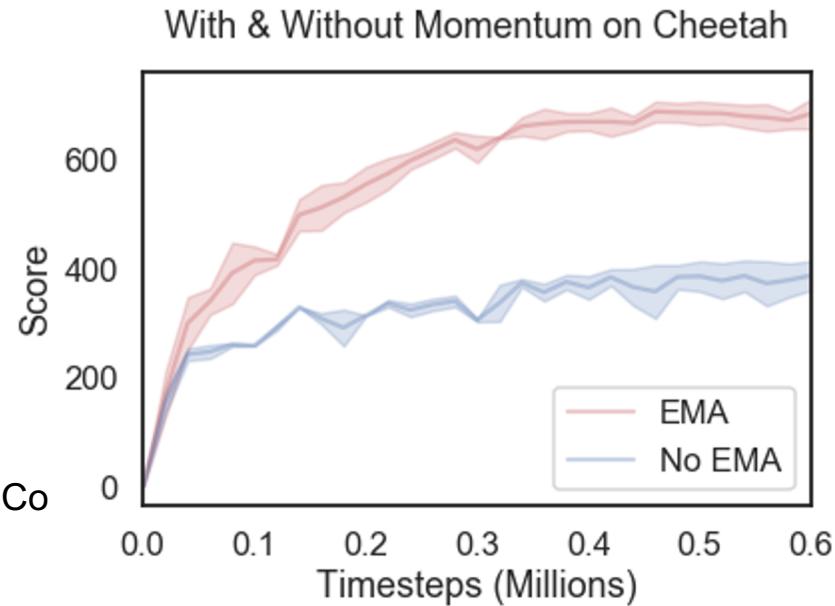
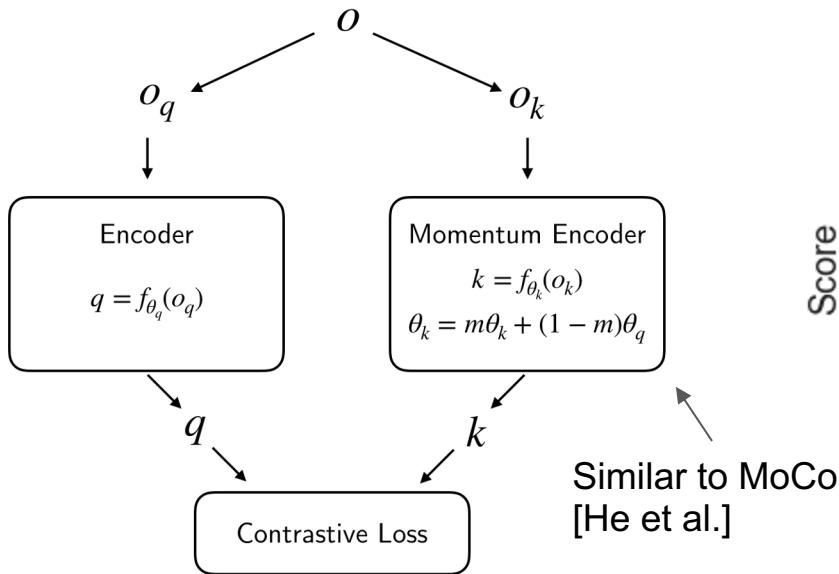
labels

$$\mathcal{L}_q = \frac{\exp(q^T W k_+)}{\exp(q^T W k_+) + \sum_{i=0}^{K-1} \exp(q^T W k_i)}$$

Comparing Similarity Measures on Cheetah

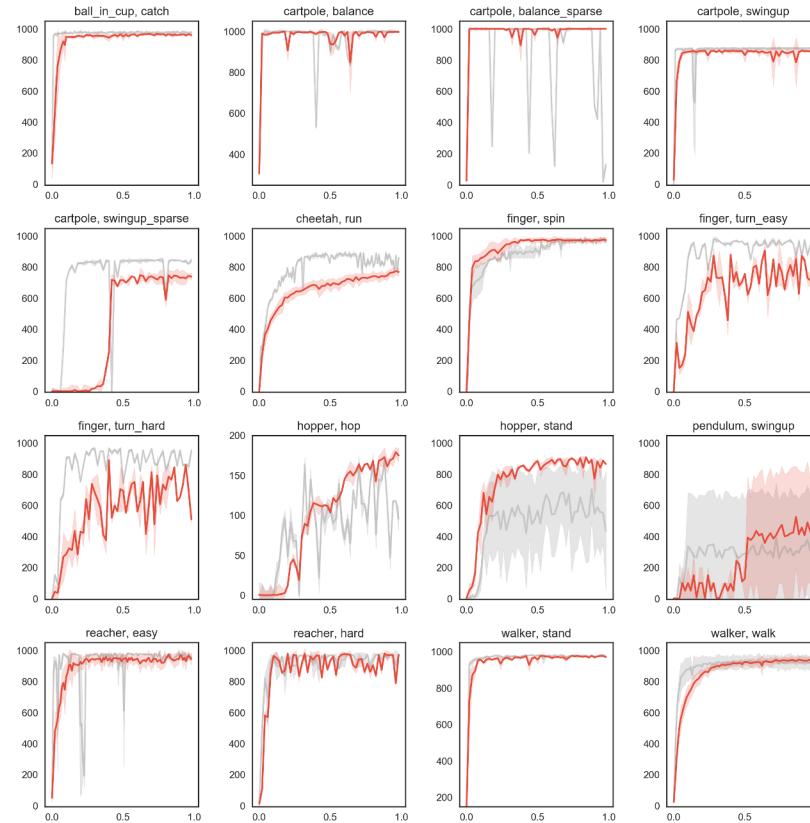


### 3. Keys encoded with momentum

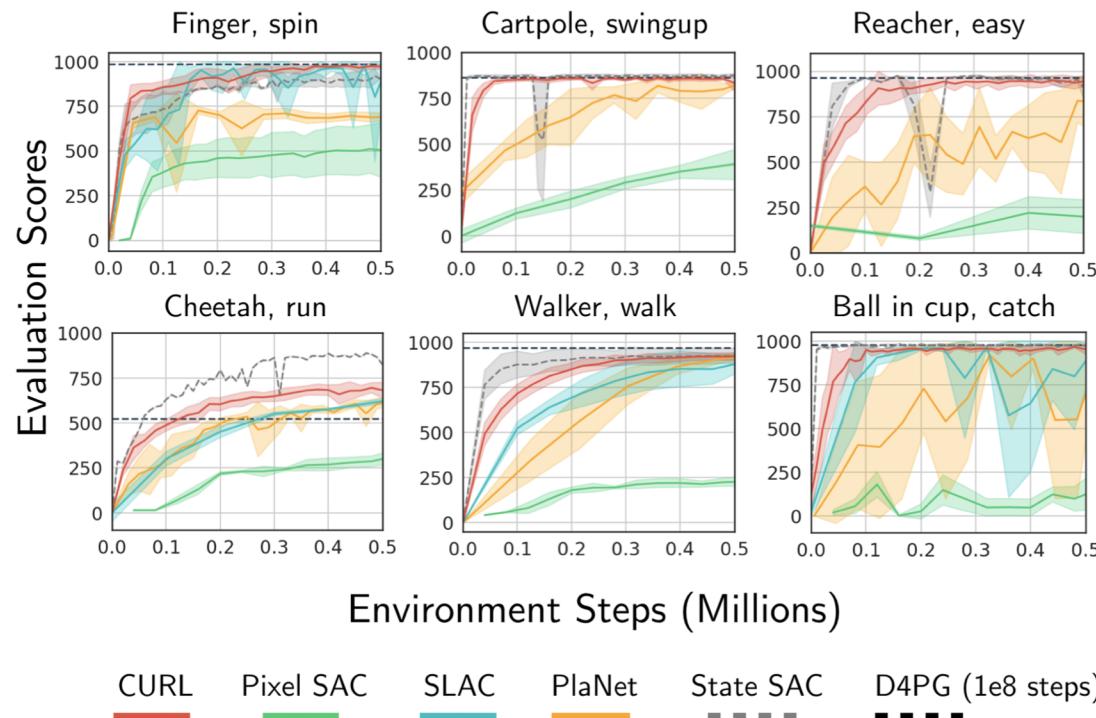


# CURL from pixels matches state-based SAC

**GRAY:** SAC State  
**RED:** CURL



# CURL Comparison: DeepMind Control Suite



# CURL Comparison: Atari

Atari performance benchmarked at 100K frames

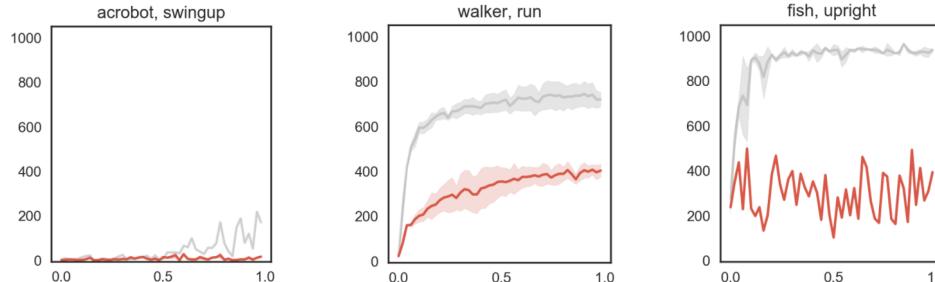
100K STEP SCORES	CURL RAINBOW	SIMPLE	RAINBOW	HUMAN	RANDOM
ALIEN	<b>1148.2</b>	616.9	318.7	6875	184.8
AMIDAR	<b>232</b>	74.3	32.5	1676	11.8
ASSAULT	473	<b>527.2</b>	231	1496	248.8
BATTLEZONE	<b>11208</b>	4031.2	3285.71	37800	2895
FREEWAY	<b>27</b>	16.7	0	29.6	0
FROSTBITE	<b>924</b>	236.9	60.2	4335	74
JAMESBOND	<b>400</b>	100.5	47.4	406.7	29.2
QBERT	<b>1352</b>	1288.8	123.46	13455	166.1
SEAQUEST	408	<b>683.3</b>	131.69	20182	61.1

# Hard Environments

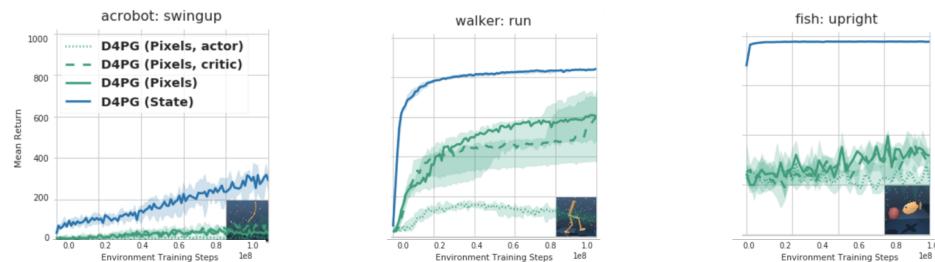
GRAY: SAC State

RED: CURL

Environment training steps 1 = 1M

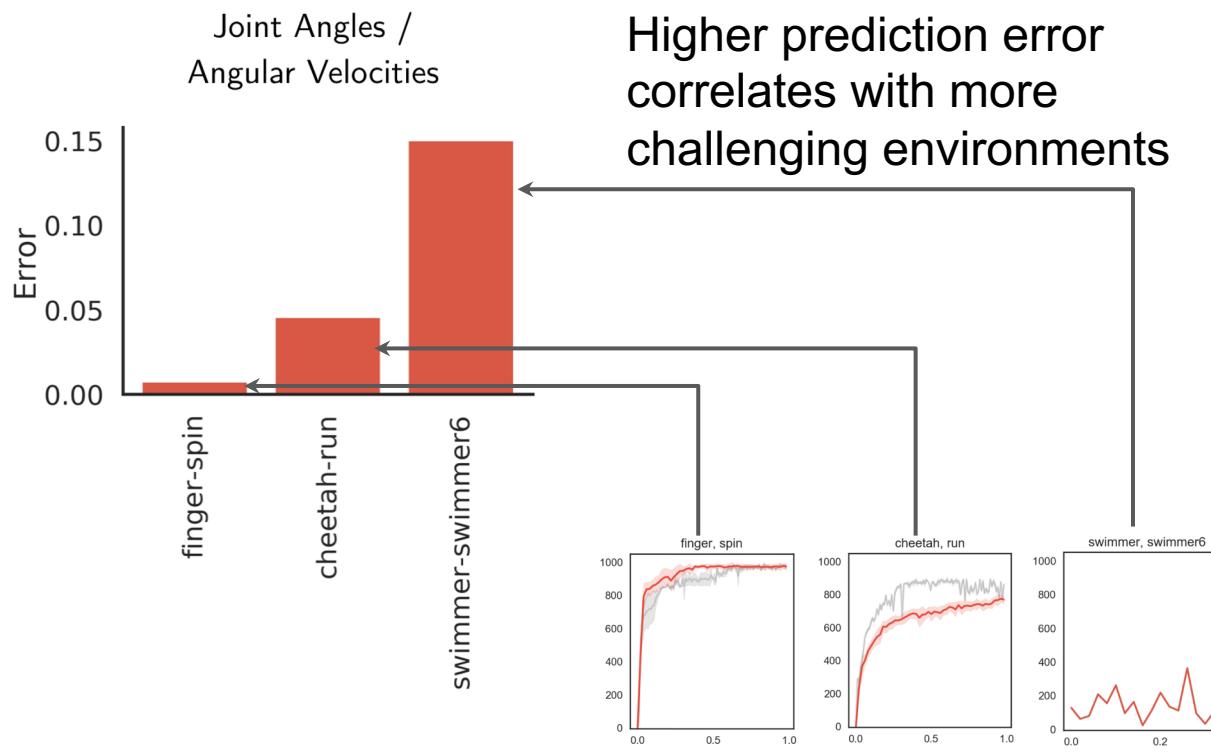


Environment training steps 1 = 100M



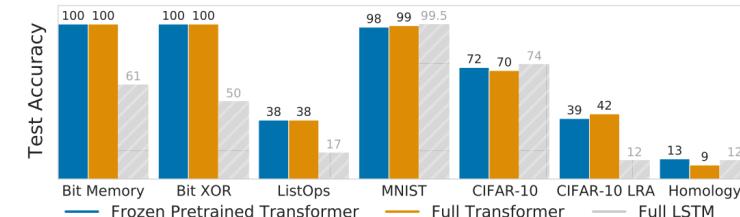
**Observation:**  
similar struggle  
asymptotically for  
pixel-based Deep RL

# Predicting state from pixels

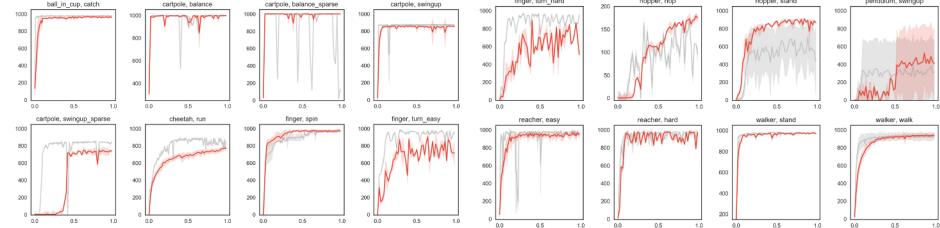


# Outline

- NN Architectures: Pretrained Transformers as Universal Computation Engines



- Representation Learning for Reinforcement Learning



- *Active Pre-Training for Reinforcement Learning*

- Human-in-the-loop Reinforcement Learning

# RL with Unsupervised Pretraining:

---

- Agent is allowed to train for a long period:
  - ***with*** access to the environment
  - ***without*** access to the ultimate reward

# RL with Unsupervised Pretraining

Algorithm	Objective	Visual Domain	Exploration	Off-policy	Pre-Trained model
MaxEnt	$\max H(s)$	✗	✓*	✗	$\pi(a s)$
CBB	$\max \mathbb{E}_s [c(s)]$	✗	✓	✓	$\pi(a s)$
MEPOL	$\max H(s)$	✗	✓*	✗	$\pi(a s)$
VISR	$\max -H(z s)$	✓	✗	✓	$\psi(s, z), \phi(s)$
DIAYN	$\max -H(z s) + H(a z, s)$	✗	✓*	✓	$\pi(a s, z)$
EDL	$\max H(s) - H(s z)$	✗	✓*	✓	$\pi(a s, z)$
DADS	$\max H(s) - H(s z)$	✗	✗	✓	$\pi(a s, z), q(s' s, z)$
APT	$\max H(s)$	✓	✓	✓	$\pi(a s), Q(s, a)$

$\psi(s, a)$ : successor feature,  $\phi(s)$ : state representation,  $c(s)$ : count-based bonus. \*: only in state-based RL.

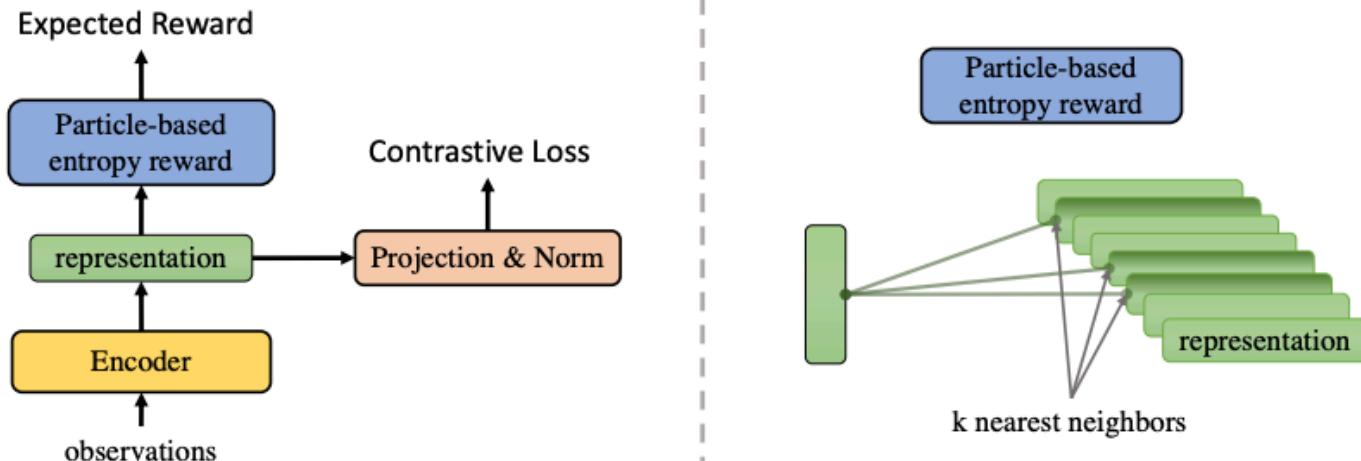
Table 1: Methods for pre-training RL in reward-free setting. DIAYN (Eysenbach et al., 2018), Count based bonus(CBB) (Bellemare et al., 2016), MEPOL (Mutti et al., 2020), VISR (Hansen et al., 2020), DADS (Sharma et al., 2019), EDL (Campos et al., 2020). Exploration: the method can explore efficiently. Visual: the method works well in visual RL. Off-policy: the method is off-policy RL.

# Our Approach (APT)

- During pre-training:

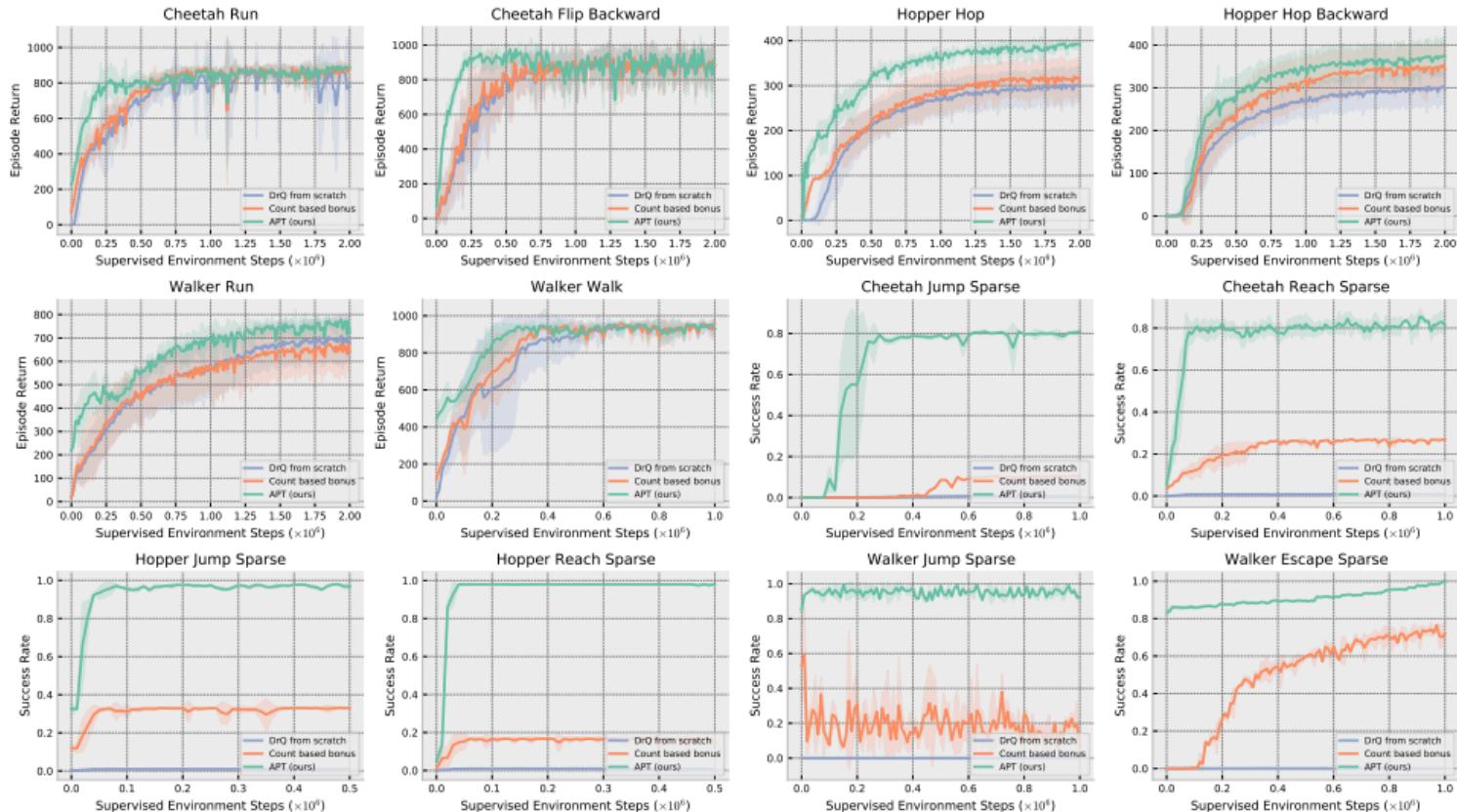
$$r(s, a, s') = \log \left( c + \frac{1}{k} \sum_{z^{(j)} \in N_k(z)} \|z - z^{(j)}\|_{n_z}^{n_z} \right)$$

where  $z = f_\theta(s)$



# Experiments: Deepmind Control Suite

DrQ from scratch  
Count based bonus  
APT (ours)



5MM pre-training steps

# Experiments: Deepmind Control Suite

Methods	DrQ	CURL	CBB	DIAYN*	MEPOL*	APT (ours)	SAC (state)
Cheetah Run	$660 \pm 96$	$641 \pm 51$	$113 \pm 69$	$81 \pm 29$	$156 \pm 128$	<b><math>671 \pm 89</math></b>	$772 \pm 60$
Hopper Hop	$315 \pm 103$	$289 \pm 87$	$327 \pm 121$	$171 \pm 45$	$61 \pm 29$	<b><math>398 \pm 45</math></b>	$285 \pm 95$
Walker Run	$713 \pm 139$	$676 \pm 127$	$381 \pm 171$	$137 \pm 43$	$275 \pm 71$	<b><math>789 \pm 58</math></b>	$801 \pm 76$
Cheetah Jump Sparse	.0 ± .0	.0 ± .0	.18 ± .08	.21 ± .09	.05 ± .02	<b>.79 ± .13</b>	.56 ± .23
Hopper Reach Sparse	.13 ± .05	.0 ± .0	.17 ± .31	.12 ± .08	.13 ± .06	<b>.94 ± .04</b>	.71 ± .18
Walker Turnover Sparse	.0 ± .0	.0 ± .0	.67 ± .14	.58 ± .22	.43 ± .12	<b>.97 ± .03</b>	.69 ± .16

# Experiments: Atari

Algorithm	26 Game Subset			Full 57 Games		
	Mdn	M	> H	Mdn	M	> H
<i>// Fully-Supervised Training</i>						
SimPLE @100K	14.39	44.30	2	–	–	–
OTR @100K	20.40	26.42	1	–	–	–
PPO @500K	20.93	43.74	7	–	–	–
DQN @10M	27.80	52.95	7	8.61	27.55	7
SPR @100K	41.50	70.40	7	–	–	–
CURL @100K	17.50	38.10	2	–	–	–
DrQ @100K	28.42	35.70	2	–	–	–
<i>// Unsupervised Pre-Training w/ Supervised Fine-Tuning</i>						
DIAYN @100K	1.34	25.39	2	2.95	23.90	6
VISR @100K	9.50	<b>128.07</b>	7	6.81	102.31	11
GPI VISR @100K	6.59	111.23	7	8.99	<b>109.16</b>	<b>12</b>
CBB@100K	1.23	21.94	3	–	–	–
MEPOL@100K	0.34	17.94	2	–	–	–
APT(ours)@100K	<b>43.50</b>	62.06	7	<b>31.55</b>	43.59	<b>12</b>

Take-aways:

- APT and VISR are strongest unsupervised pre training methods
- VISR stronger mean
- APT stronger median
- qualitatively: APT better on hard exploration games

# VISR <> APT

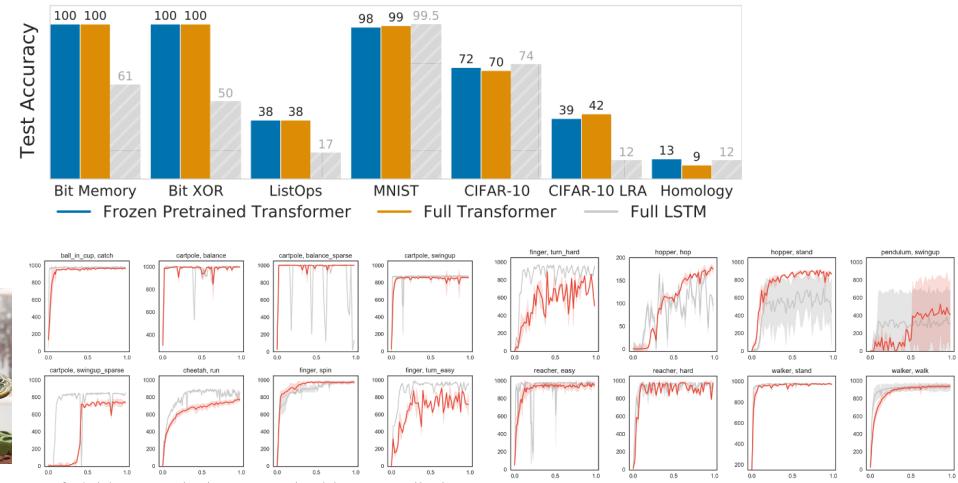
Game	Random	Human	VISR	APT
Alien	227.8	7127.7	364.4	<b>2614.8</b>
Amidar	5.8	1719.5	186.0	<b>211.5</b>
Assault	222.4	742.0	<b>1209.1</b>	891.5
Asterix	210.0	8503.3	<b>6216.7</b>	185.5
Asteroids	7191	47388.7	<b>4443.3</b>	678.7
Atlantis	12850.0	29028.1	<b>140542.8</b>	40231.0
Bank Heist	14.2	753.1	71.3	<b>416.7</b>
Battle Zone	2360.0	37187.5	<b>7072.7</b>	7065.1
Beam Rider	363.9	16826.5	1741.9	<b>3487.2</b>
Berzerk	123.7	2630.4	490.0	<b>493.4</b>
Bowling	23.1	160.7	<b>21.2</b>	-56.5
Boxing	0.1	12.1	13.4	<b>21.3</b>
Breakout	1.7	30.5	<b>17.9</b>	10.9
Centipede	2090.9	12017.1	<b>7184.9</b>	6233.9
Chopper Command	811.0	7387.8	<b>800.8</b>	317.0
Crazy Climber	10780.5	23829.4	<b>49373.9</b>	44128.0
Defender	2874.5	18688.9	<b>15876.1</b>	5927.9
Demon Attack	10780.5	35829.4	<b>8994.9</b>	6871.8
Double Dunk	-18.6	-16.4	-22.6	<b>-17.2</b>
Enduro	0.0	860.5	-3.1	<b>-0.3</b>
Fishing Derby	-91.7	-38.7	-93.9	<b>-5.6</b>
Freeway	0.0	29.6	-12.1	<b>29.9</b>
Frostbite	65.2	4334.7	230.9	<b>1796.1</b>
Gopher	257.6	2412.5	498.6	<b>2190.4</b>
Gravitar	173.0	3351.4	328.1	<b>542.0</b>
Hero	1027.0	30826.4	663.5	<b>6789.1</b>
Ice Hockey	-11.2	0.9	<b>-18.1</b>	-30.1
Jamesbond	29.0	302.8	<b>484.4</b>	356.1
Kangaroo	52.0	3035.0	<b>1761.9</b>	412.0
Krull	1598.0	2665.5	<b>3142.5</b>	2312.0
Kung Fu Master	258.5	22736.3	16754.9	<b>17357.0</b>
Montezuma Revenge	0.0	4753.3	0.0	<b>0.2</b>
Ms Pacman	307.3	6951.6	558.5	<b>2527.1</b>
Name This Game	2292.3	8049.0	<b>2605.8</b>	1387.2
Phoenix	761.4	7242.6	<b>7162.2</b>	3874.2
Pitfall	-229.4	6463.7	-370.8	<b>-12.8</b>
Pong	-20.7	14.6	-26.2	<b>-8.0</b>
Private Eye	24.9	69571.3	<b>98.3</b>	96.1
Qbert	163.9	13455.0	666.3	<b>17671.2</b>
Riverraid	1338.5	17118.0	<b>5422.2</b>	4671.0
Road Runner	11.5	7845.0	<b>6146.7</b>	4782.1
Robotank	2.2	11.9	10.0	<b>13.7</b>
Seasekt	68.4	42054.7	706.6	<b>2116.7</b>
Skiing	-17098.1	-4336.9	<b>-19692.5</b>	-38434.1
Solaris	1236.3	12326.7	<b>1921.5</b>	841.8
Space Invaders	148.0	1668.7	<b>9741.0</b>	3687.2
Star Gunner	664.0	10250.0	<b>25827.5</b>	8717.0
Surround	-10.0	6.5	15.5	<b>2.5</b>
Tennis	-23.8	-8.3	0.7	<b>1.2</b>
Time Pilot	3568.0	5229.2	<b>4503.6</b>	2567.0
Tutankham	11.4	167.6	50.7	<b>124.6</b>
Up N Down	533.4	11693.2	<b>10037.6</b>	8289.4
Venture	0.0	1187.5	-1.7	<b>231.0</b>
Video Pinball	0.0	17667.9	<b>35120.3</b>	2817.1
Wizard Of Wor	563.5	4756.5	853.3	<b>1265.0</b>
Yars Revenge	3092.9	54576.9	<b>5543.5</b>	1871.5
Zaxxon	32.5	9173.3	897.5	<b>3231.0</b>
Mean Human-Norm'd	0.000	1.000	68.42	<b>81.24</b>
Median Human-Norm'd	0.000	1.000	9.41	<b>66.89</b>
#Superhuman	0	N/A	11	<b>12</b>

# Outline

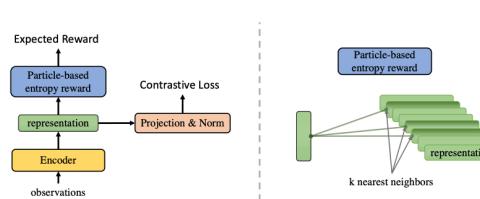
- NN Architectures: Pretrained Transformers as Universal Computation Engines



- Representation Learning for Reinforcement Learning

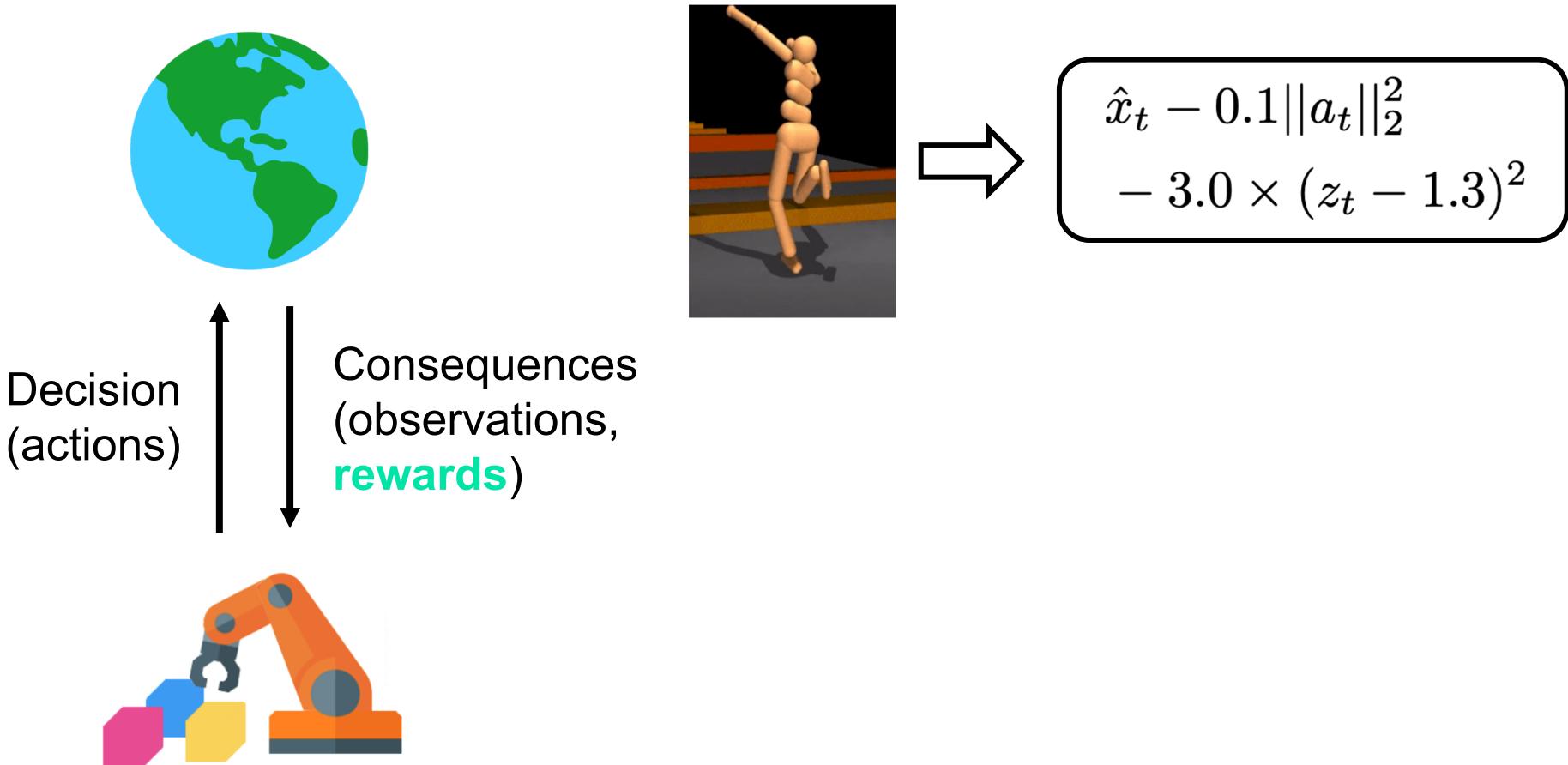


- Active Pre-Training for Reinforcement Learning

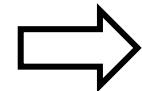
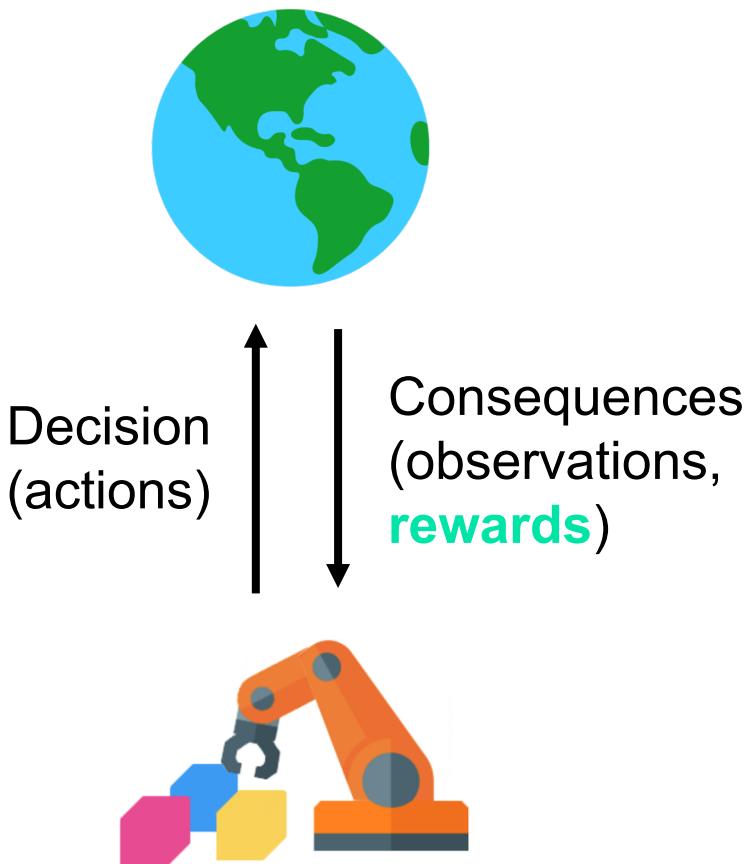


- *Human-in-the-loop Reinforcement Learning*

# Challenge: Designing Suitable Reward



# Challenge: Designing Suitable Reward

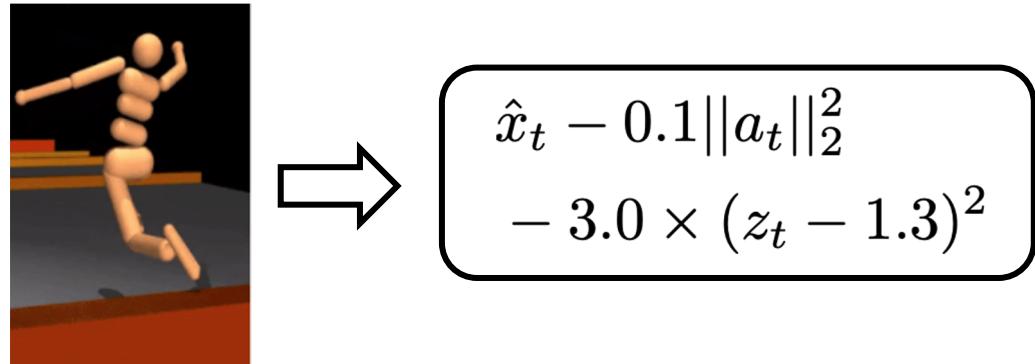
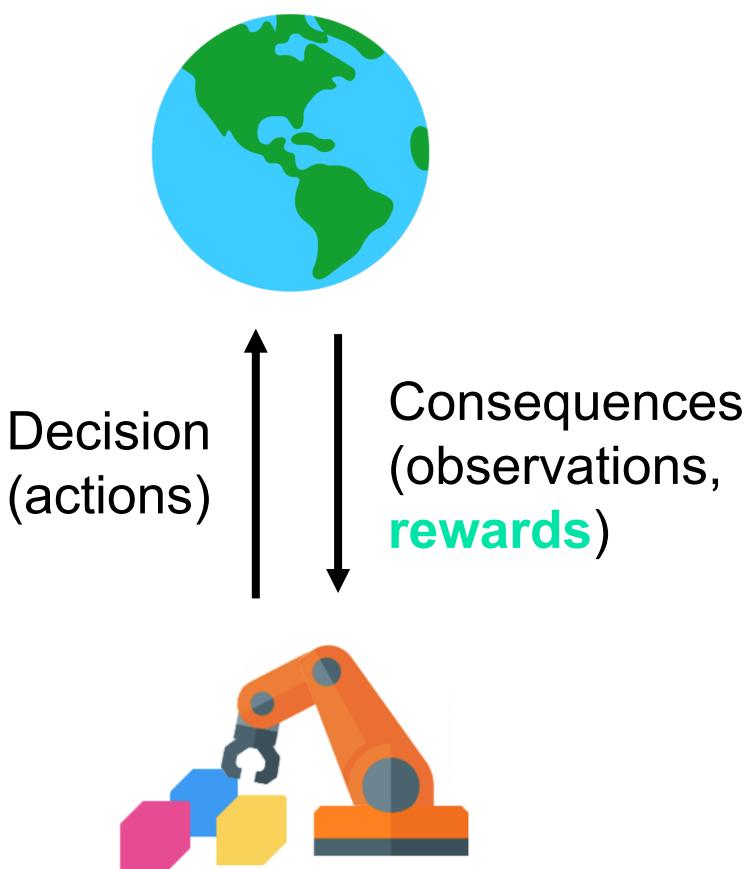


$$\begin{aligned}\hat{x}_t - 0.1\|a_t\|_2^2 \\ - 3.0 \times (z_t - 1.3)^2\end{aligned}$$



Hard tasks to define a reward (e.g. cooking)

# Challenge: Designing Suitable Reward



$$\hat{x}_t - 0.1\|a_t\|_2^2 - 3.0 \times (z_t - 1.3)^2$$



Hard tasks to define a reward (e.g. cooking)



Reward exploitation

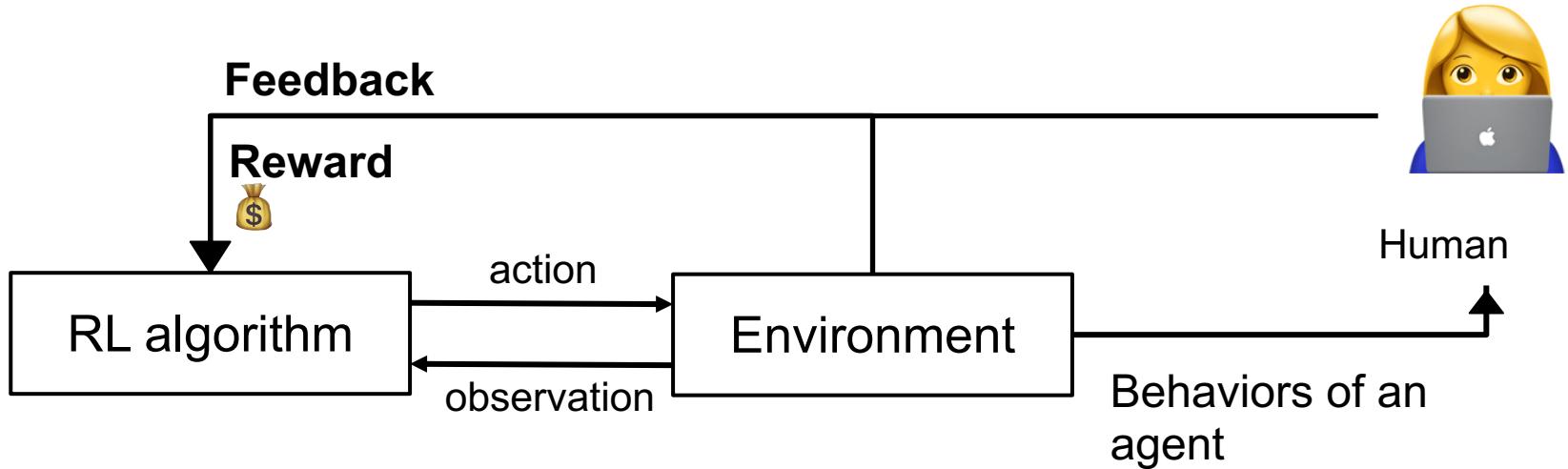
<https://openai.com/blog/faulty-reward-functions>

# What is Alternative Solution?

---

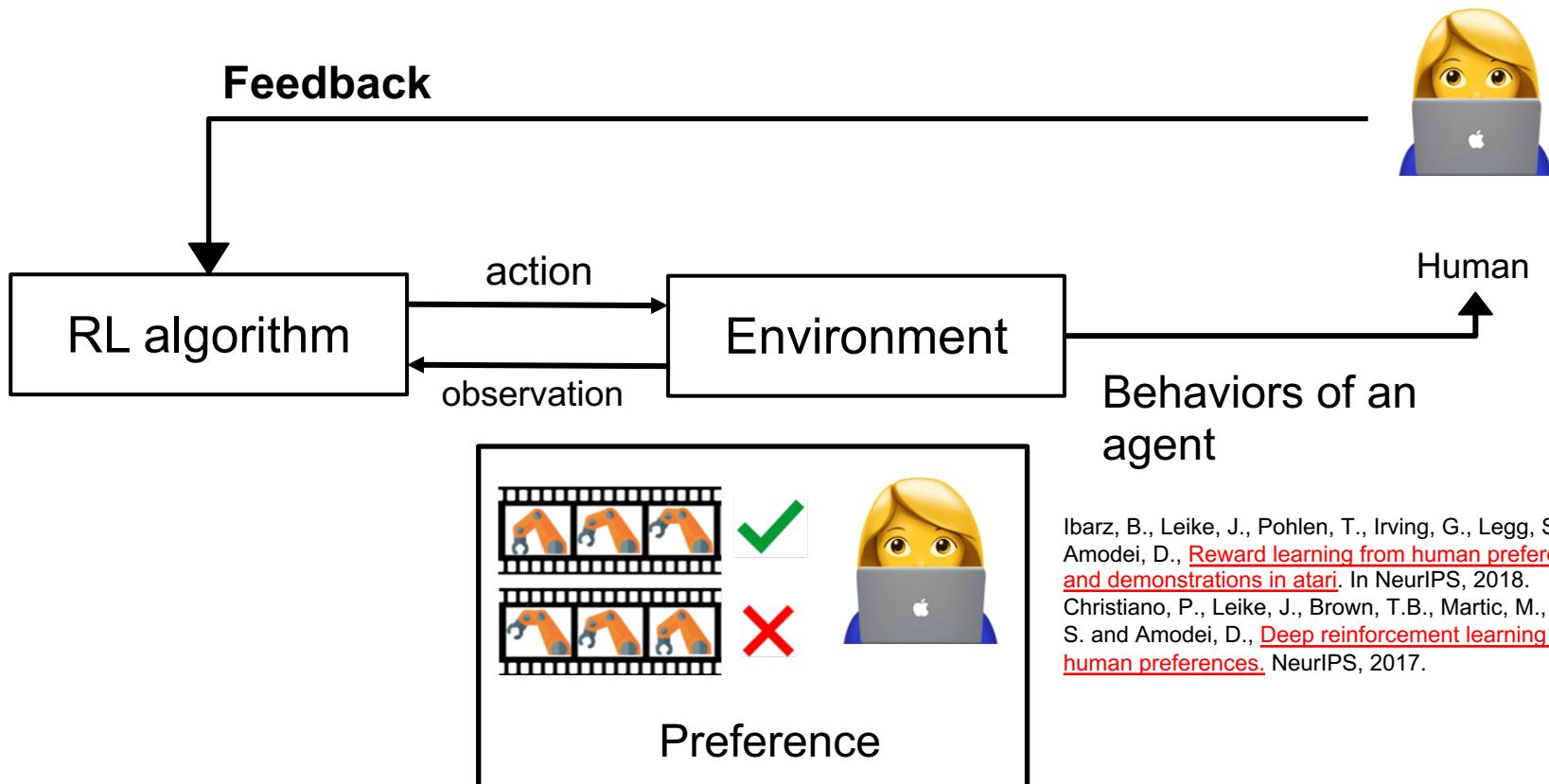
# What is Alternative Solution?

- Putting (non-expert) humans into the agent learning loop!



# What is Alternative Solution?

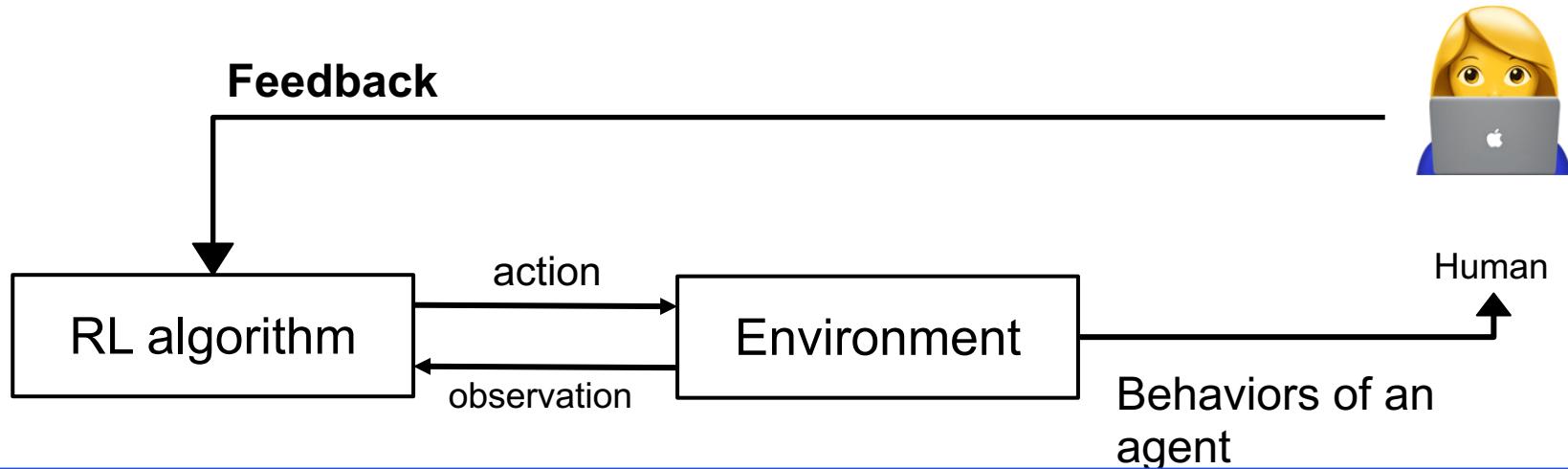
- Putting (non-expert) humans into the agent learning loop!



Ibarz, B., Leike, J., Pohlen, T., Irving, G., Legg, S. and Amodei, D., [Reward learning from human preferences and demonstrations in atari](#). In NeurIPS, 2018.  
Christiano, P., Leike, J., Brown, T.B., Martic, M., Legg, S. and Amodei, D., [Deep reinforcement learning from human preferences](#). NeurIPS, 2017.

# What is Alternative Solution?

- Putting (non-expert) humans into the agent learning loop!



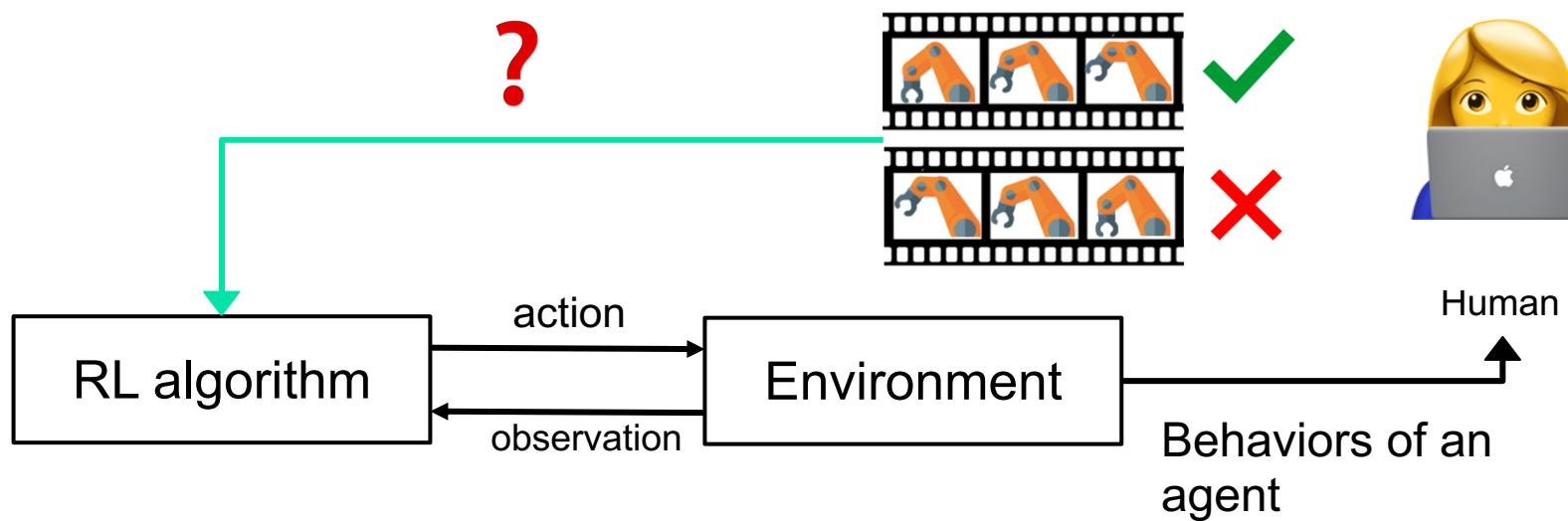
Human has ability to interactively guide agents according to their progress

- Can teach more hard tasks, where we can't easily define the reward
- Can avoid reward exploitation

# Research Questions



Q1. How to train RL agents using preferences?



# Research Questions



Q1. How to train RL agents using preferences?



Q2. How to design a feedback-efficient algorithm?



Large burden can be placed on the human!

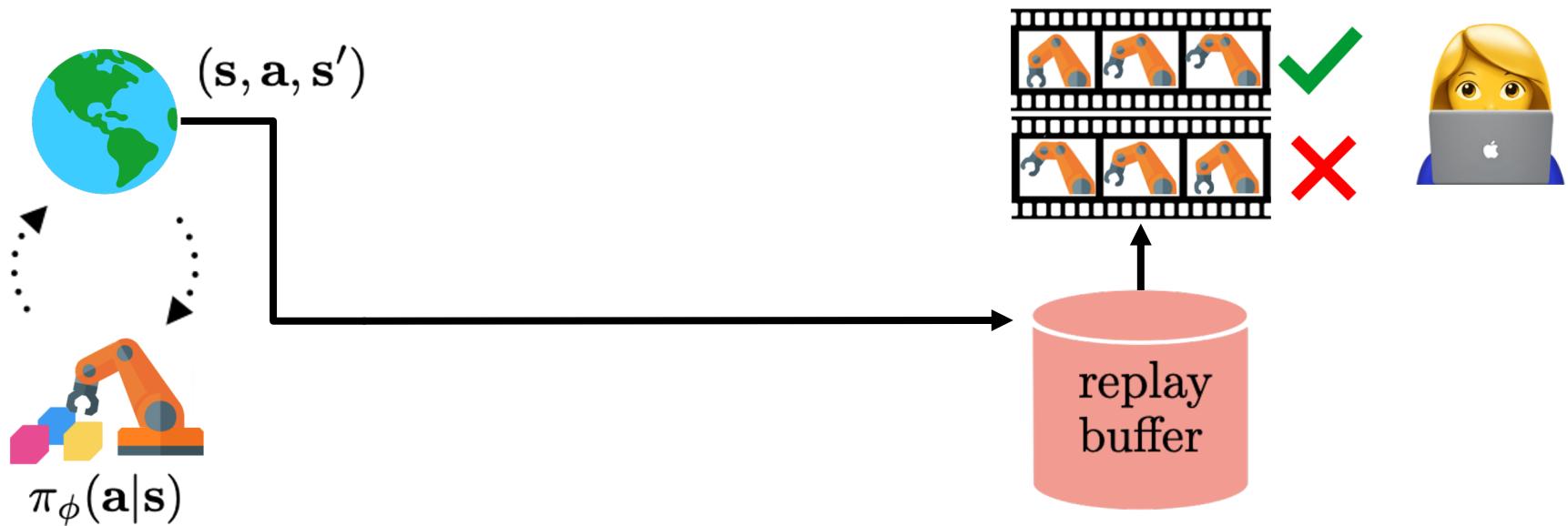
# Overall Framework

- Step 1. Collect samples via interactions with environment



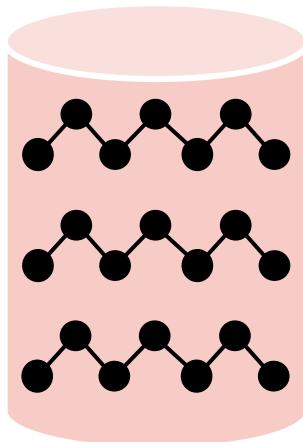
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences



# Overall Framework

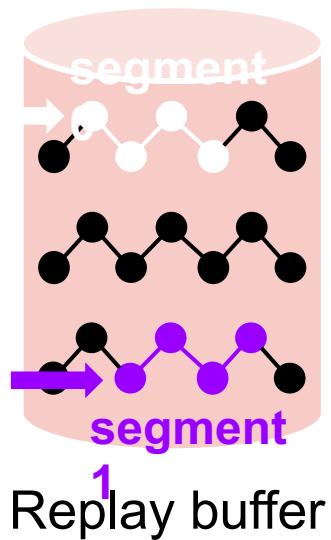
- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
  - How to select two segments? Entropy sampling



Replay buffer

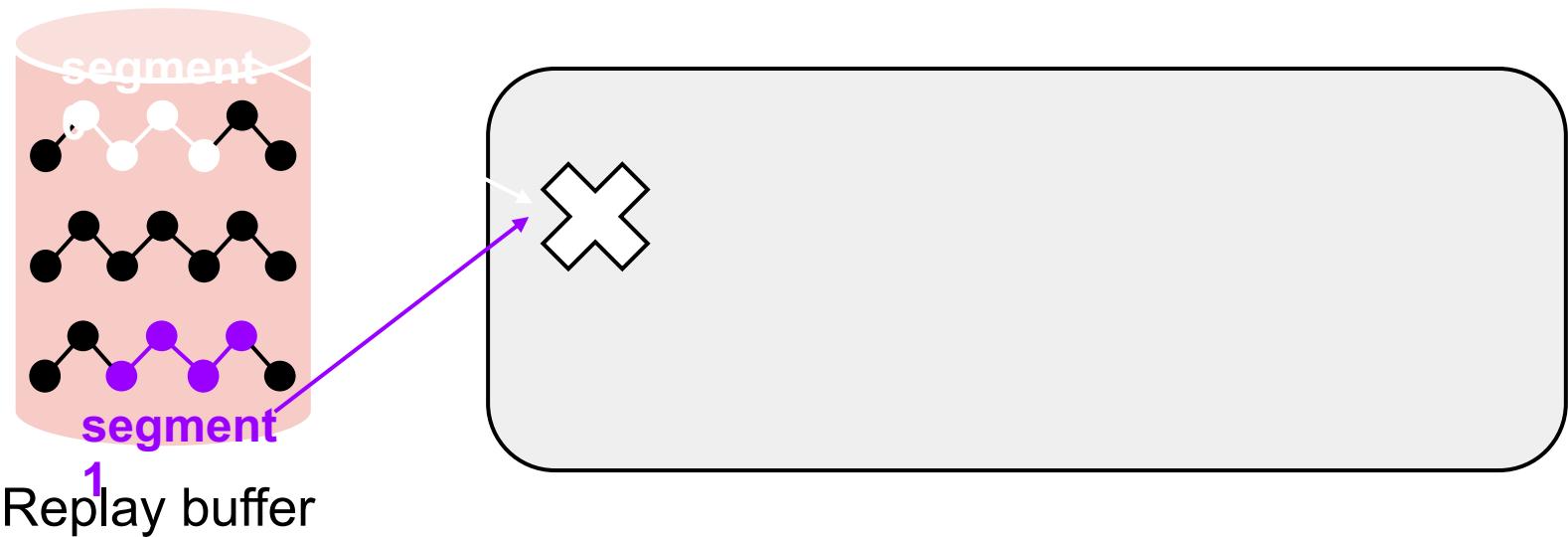
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
  - How to select two segments? Entropy sampling



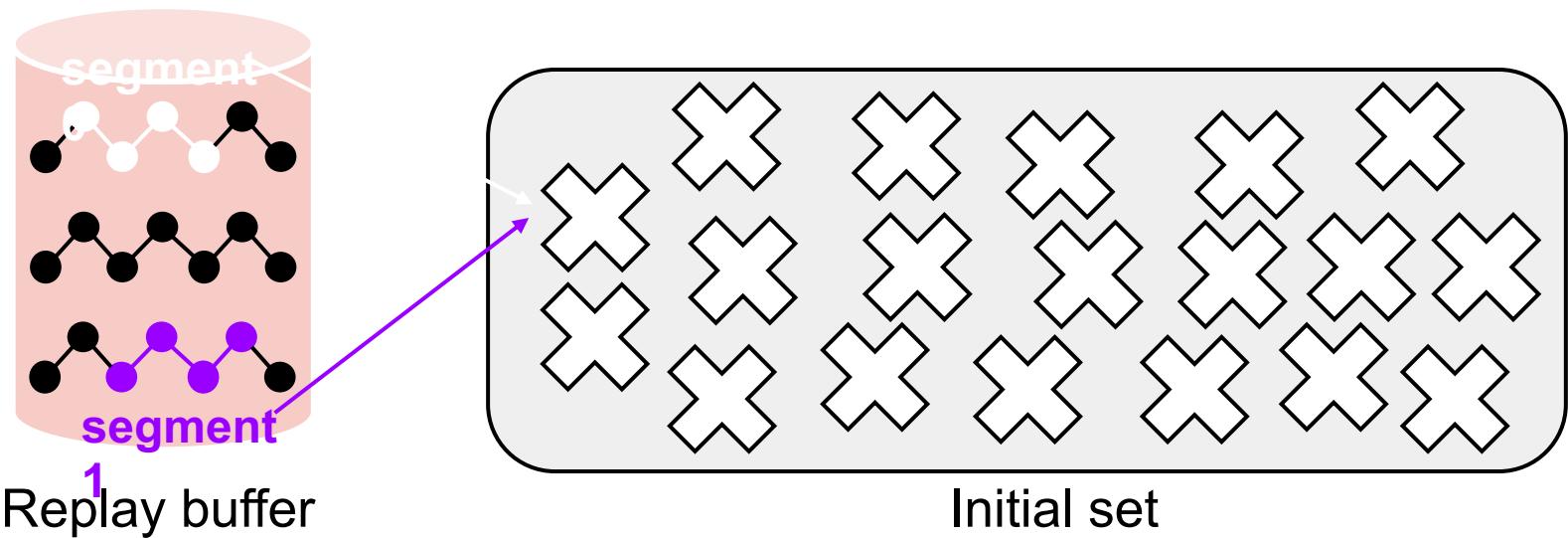
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
  - How to select two segments? Entropy sampling



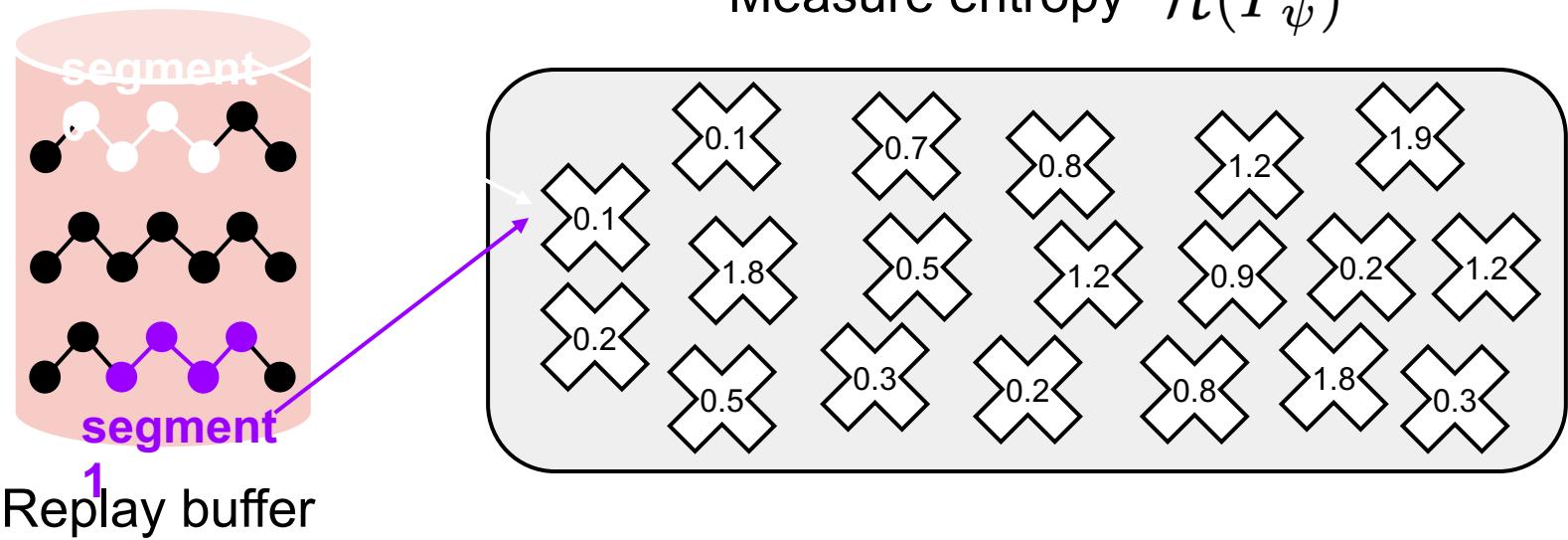
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
  - How to select two segments? Entropy sampling



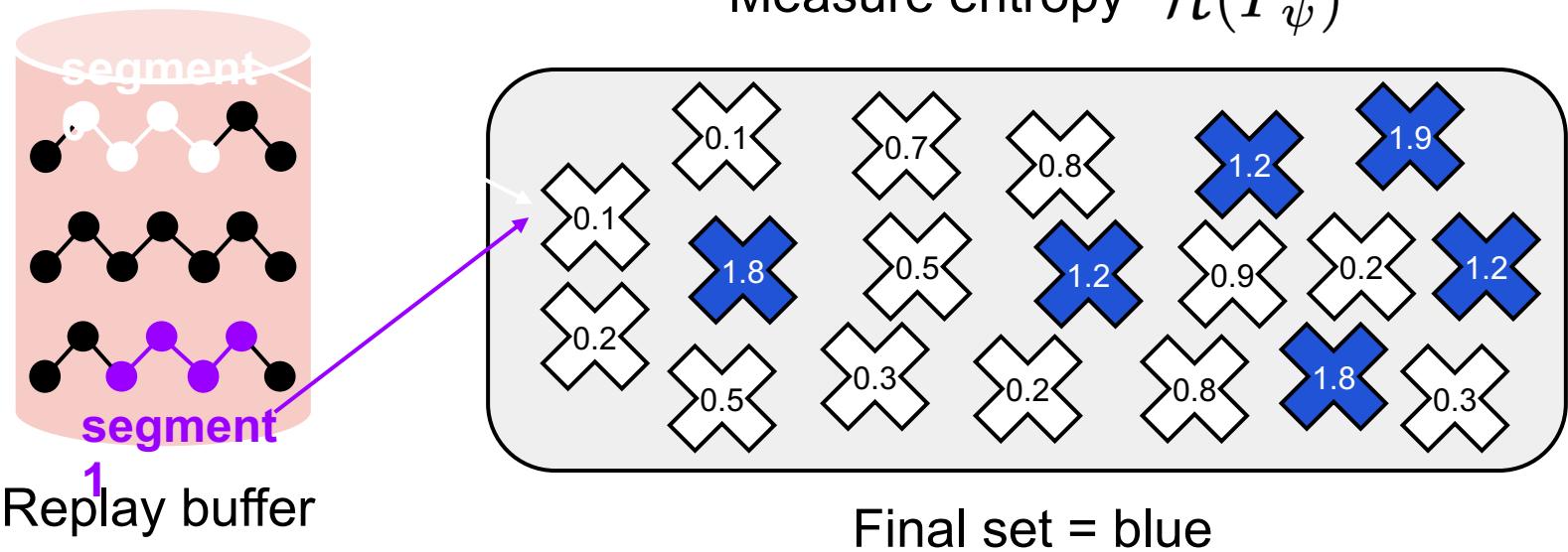
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
  - How to select two segments? Entropy sampling



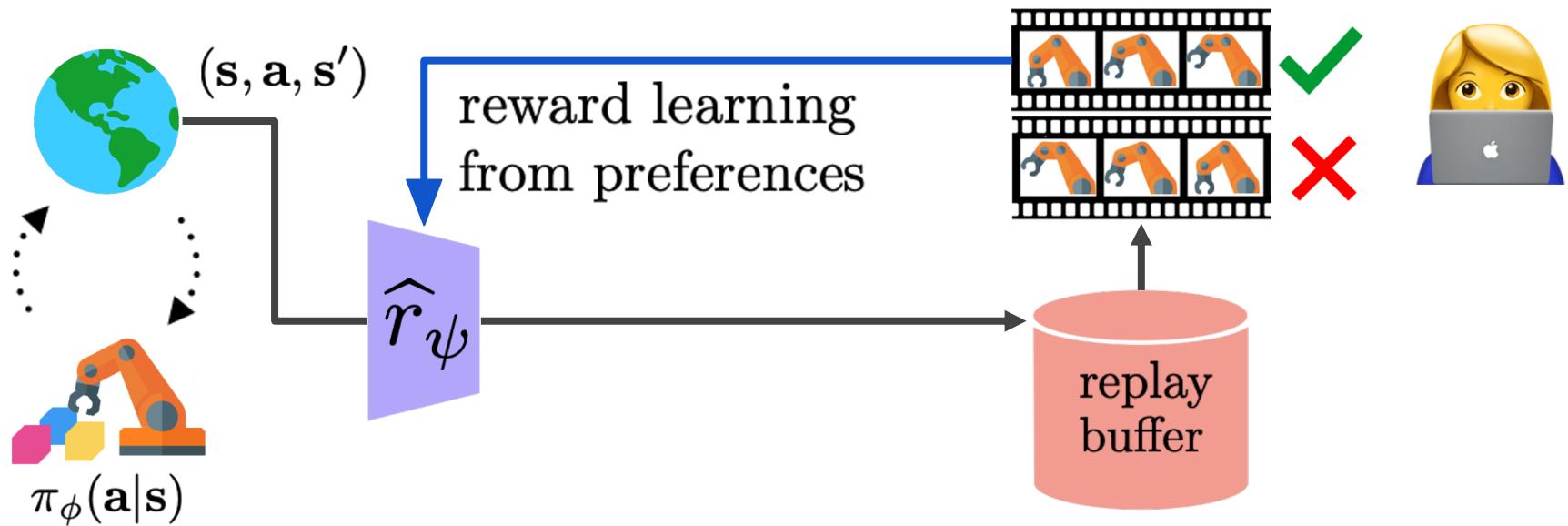
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
  - How to select two segments? Entropy sampling



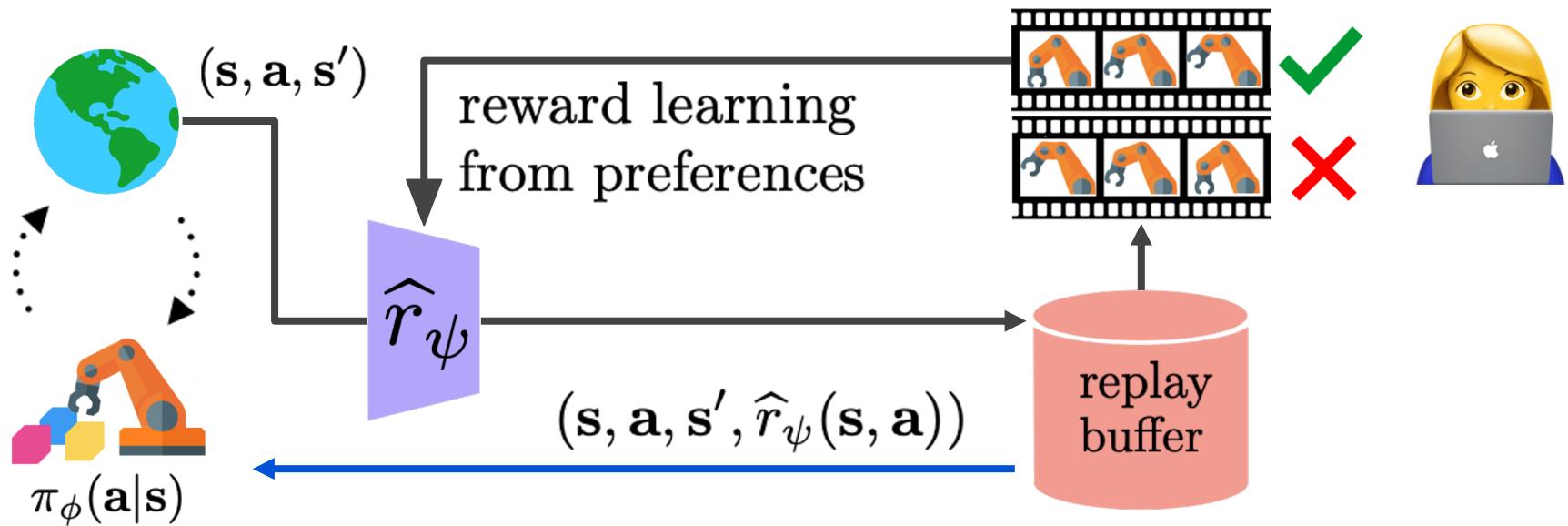
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
- Step 3. Optimize a reward model using cross entropy loss



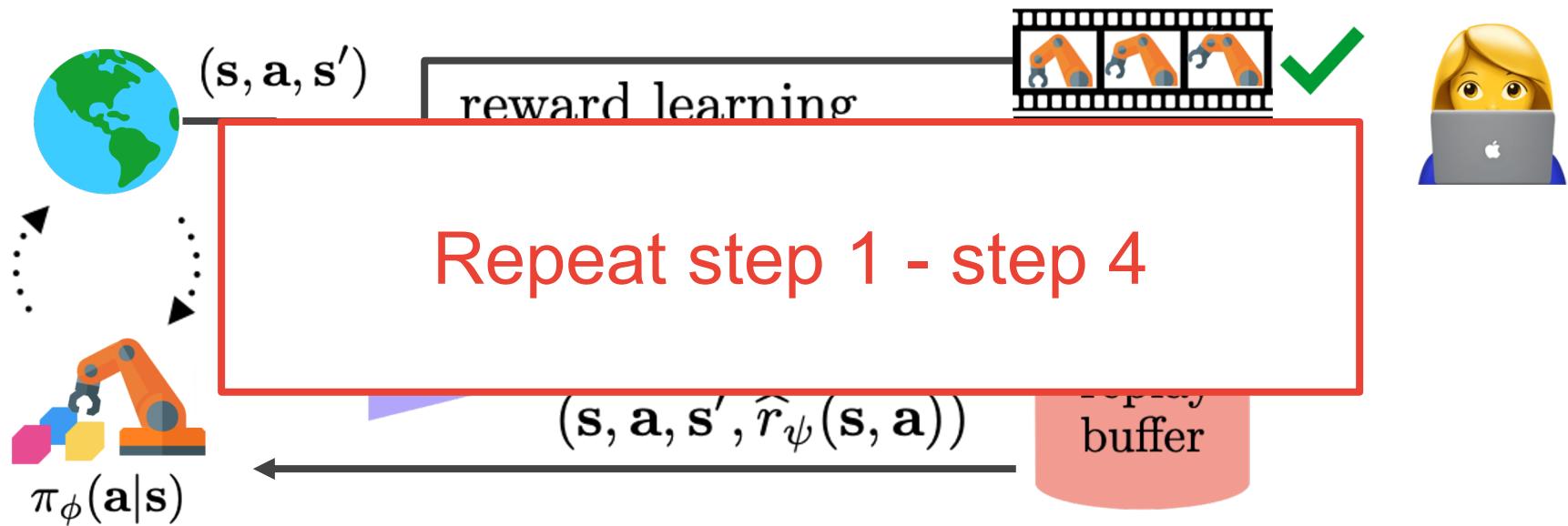
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
- Step 3. Optimize a reward model using cross entropy loss
- Step 4. Optimize a policy using off-policy algorithms



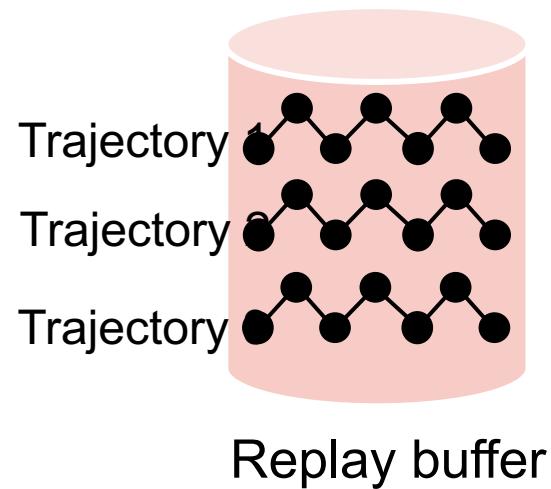
# Overall Framework

- Step 1. Collect samples via interactions with environment
- Step 2. Collect human preferences
- Step 3. Optimize a reward model using cross entropy loss
- Step 4. Optimize a policy using off-policy algorithms



# Learning Reward from Preferences

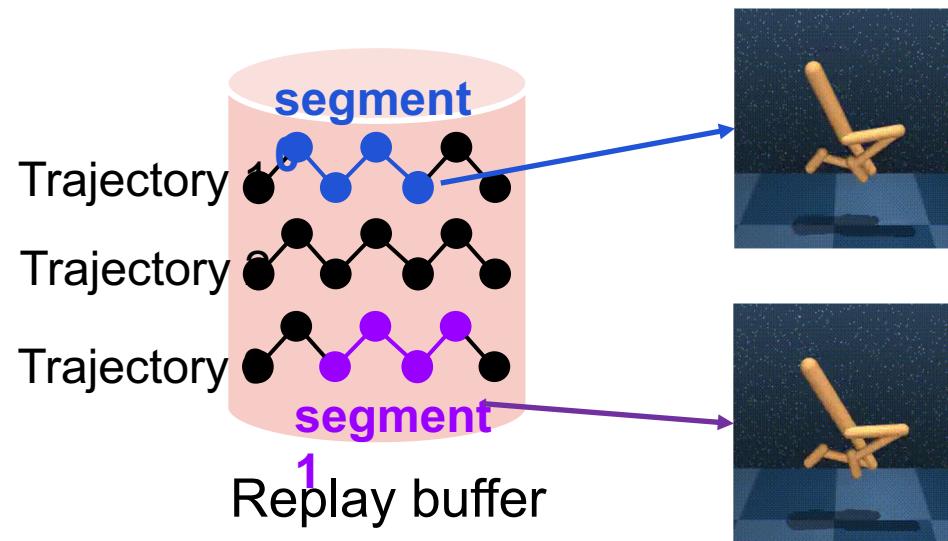
- Setup



# Learning Reward from Preferences

- Setup
  - Select two segments.  $(\sigma^0, \sigma^1)$ , i.e. behaviors, from buffer

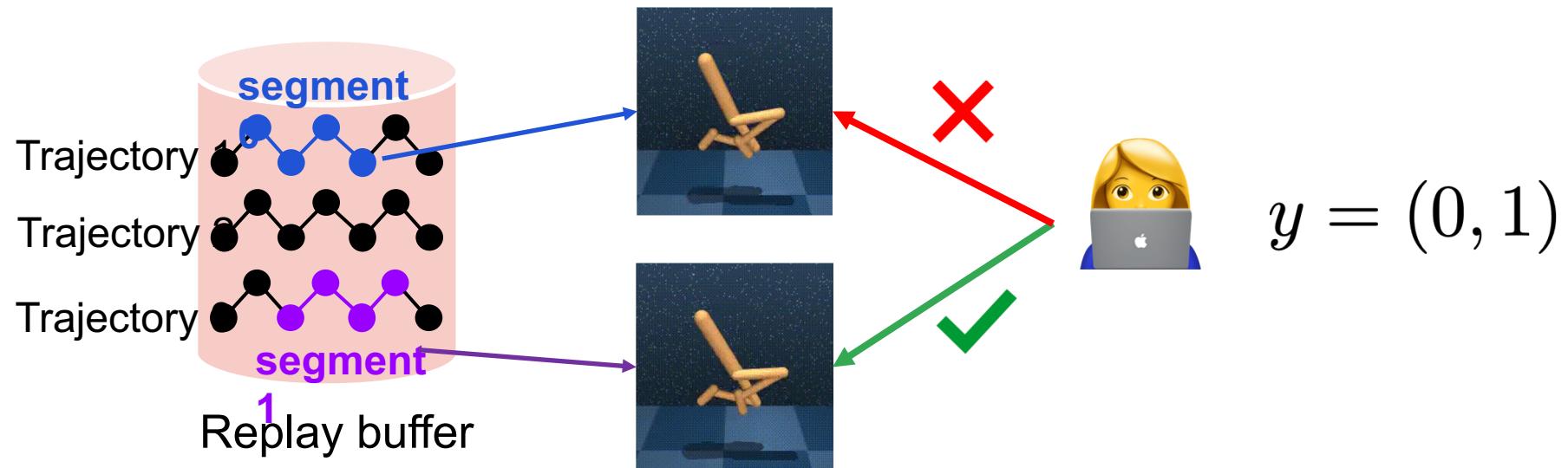
$$\sigma^i = \{(\mathbf{s}_t^I, \mathbf{a}_t^i), \dots, (\mathbf{s}_{t+H}^I, \mathbf{a}_{t+H}^i)\}$$



# Learning Reward from Preferences

- Setup
  - Select two segments  $(\sigma^0, \sigma^1)$ , i.e. behaviors, from buffer
  - Human provides a preference as a binary label:

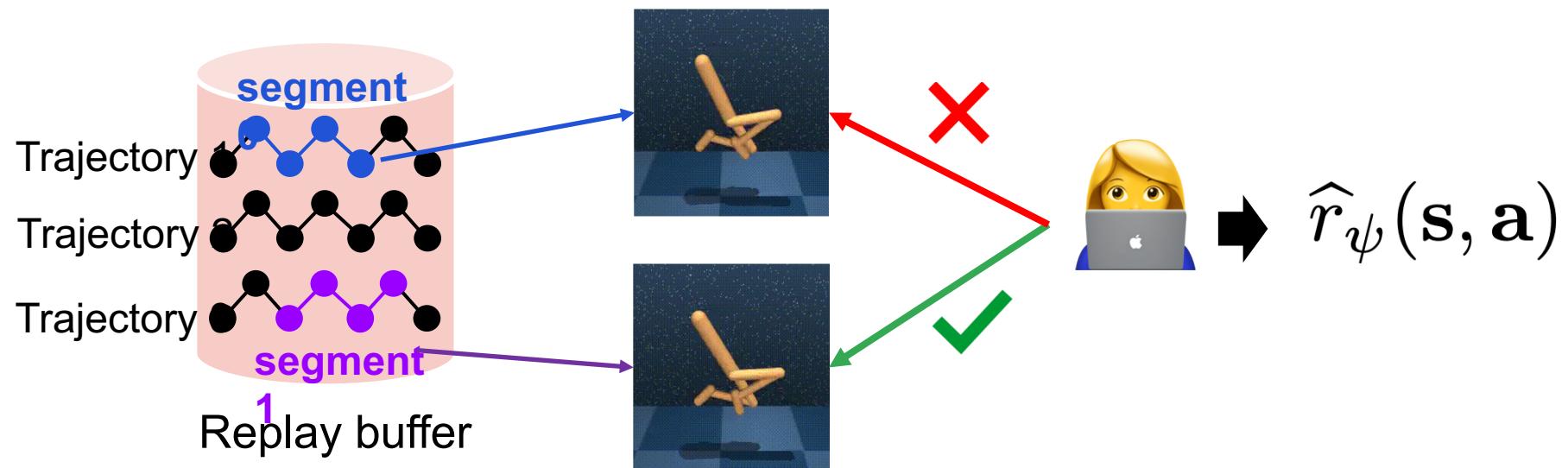
$$y \in \{(1, 0), (0, 1)\}$$



# Learning Reward from Preferences

- Setup
  - Select two segments ( $\sigma^0, \sigma^1$ ), i.e. behaviors, from buffer
  - Human provides a preference as a binary label
  - Goal: learning a reward model

$$\hat{r} : S \times A \rightarrow \mathbb{R}$$



# Learning Reward from Preferences

---

- Fitting a reward model [1]
  - Main idea: formulate this problem as a binary classification!

# Learning Reward from Preferences

- Fitting a reward model [1]
  - Main idea: formulate this problem as a binary classification!
  - By following the Bradley-Terry model [2], we can model a **preference predictor** as follows:

$$P_\psi[\sigma^1 \succ \sigma^0] = \frac{\exp \sum_t \hat{r}(\mathbf{s}_t^1, \mathbf{a}_t^1)}{\sum_{i \in \{0,1\}} \exp \sum_t \hat{r}(\mathbf{s}_t^i, \mathbf{a}_t^i)}$$

[1] Christiano, P., Leike, J., Brown, T.B., Martic, M., Legg, S. and Amodei, D., Deep reinforcement learning from human preferences. NeurIPS, 2017.

[2] Bradley, R.A. and Terry, M.E., Rank analysis of incomplete block designs: I. The method of paired comparisons. Biometrika, 39(3/4), pp.324-345, 1952.

# Learning Reward from Preferences

- Fitting a reward model [1]
  - Main idea: formulate this problem as a binary classification!
  - By following the Bradley-Terry model [2], we can model a **preference predictor** as follows:

$$P_{\psi}[\sigma^1 \succ \sigma^0] = \frac{\exp \sum_t \hat{r}(\mathbf{s}_t^1, \mathbf{a}_t^1)}{\sum_{i \in \{0,1\}} \exp \sum_t \hat{r}(\mathbf{s}_t^i, \mathbf{a}_t^i)}$$

Sum of rewards  
over segment 1

↓

Event that segment 1 is preferable to segment 0

[1] Christiano, P., Leike, J., Brown, T.B., Martic, M., Legg, S. and Amodei, D., Deep reinforcement learning from human preferences. NeurIPS, 2017.

[2] Bradley, R.A. and Terry, M.E., Rank analysis of incomplete block designs: I. The method of paired comparisons. Biometrika, 39(3/4), pp.324-345, 1952.

# Learning Reward from Preferences

- Fitting a reward model [1]
  - Main idea: formulate this problem as a binary classification!
  - By following the Bradley-Terry model [2], we can model a **preference predictor** as follows:

$$P_\psi[\sigma^1 \succ \sigma^0] = \frac{\exp \sum_t \hat{r}(\mathbf{s}_t^1, \mathbf{a}_t^1)}{\sum_{i \in \{0,1\}} \exp \sum_t \hat{r}(\mathbf{s}_t^i, \mathbf{a}_t^i)}$$

- We can learn a reward model by optimizing cross entropy:

$$\mathcal{L}^{\text{Reward}} = -\mathbb{E}_{(\sigma^0, \sigma^1, y)} \left[ y(0) \log P_\psi[\sigma^0 \succ \sigma^1] + y(1) \log P_\psi[\sigma^1 \succ \sigma^0] \right].$$

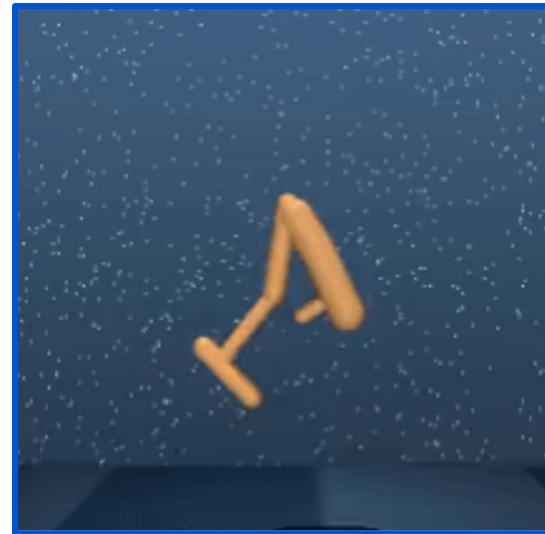
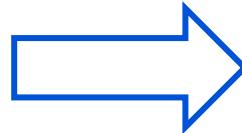
# Unsupervised Pre-training: APT

---

- Obtaining a good initial state space coverage is important!
  - Human can't convey much meaningful information to the agent

# Unsupervised Pre-training: APT

- Obtaining a good initial state space coverage is important!
  - Human can't convey much meaningful information to the agent

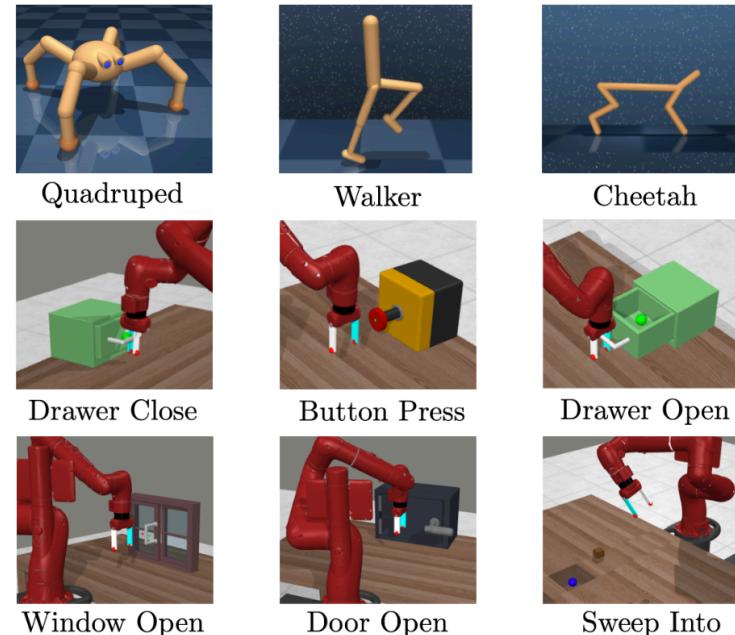


Behavior from random exploration policy

Behavior from pre-trained policy

# Experiments

- Goal
  - Solving benchmark tasks **without** observing the true reward
- Tasks
  - Locomotion tasks from DMControl [1]
  - Robotic manipulation tasks from Meta-world [2]
- We will solve these tasks using as few queries as possible!
  - How to get preferences?

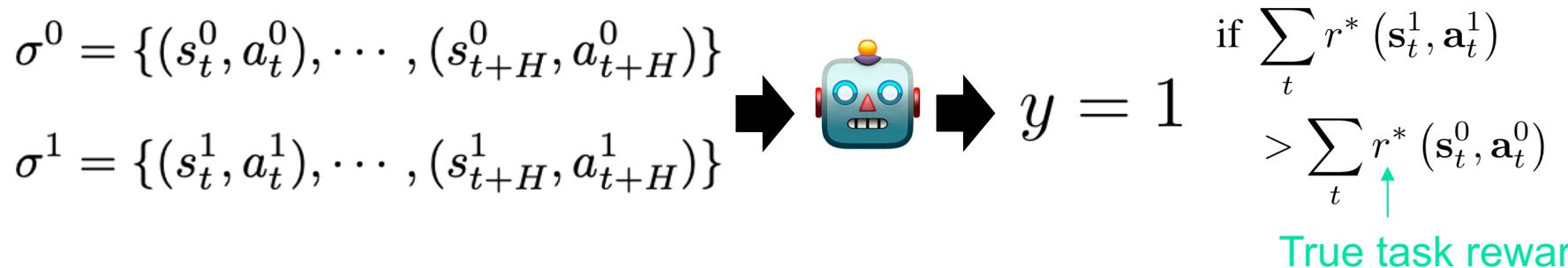


[1] Tassa, Y. et al. [Deepmind control suite](#). arXiv preprint arXiv:1801.00690, 2018.

[2] Yu, T. et al. [Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning](#). In CoRL, 2020.

# Experiments

- We generate preferences using a scripted teacher [1, 2]:



Preferences are immediately generated  
→ more rapid experiments

We can evaluate the agent quantitatively by  
measuring the true average return

[1] Ibarz, B., Leike, J., Pohlen, T., Irving, G., Legg, S. and Amodei, D., [Reward learning from human preferences and demonstrations in atari](#). In NeurIPS, 2018.

[2] Christiano, P., Leike, J., Brown, T.B., Martic, M., Legg, S. and Amodei, D., [Deep reinforcement learning from human preferences](#). NeurIPS, 2017.

# Tested Algorithms

---

- **PEBBLE (ours)**



- Unsupervised pre-training (**yellow**)
  - Collect initial samples via unsupervised pre-training

# Tested Algorithms

- **PEBBLE (ours)**



- Unsupervised pre-training (**yellow**)
  - Collect initial samples via unsupervised pre-training
- Learn a reward (**red**)
  - Get feedback (preference btw two behaviors) from a teacher
  - Learn a reward model

# Tested Algorithms

- **PEBBLE (ours)**



- Unsupervised pre-training (**yellow**)
  - Collect initial samples via unsupervised pre-training
- Learn a reward (**red**)
  - Get feedback (preference btw two behaviors) from a teacher
  - Learn a reward model
- Interaction (**blue**)
  - Interact with environment (e.g. 20 episodes)
  - Update policy using off-policy RL (SAC [1])

[1] Haarnoja, T., Zhou, A., Abbeel, P. and Levine, S., Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In ICML, 2018.

# Tested Algorithms

- **PEBBLE (ours)**



- **Preference-PPO [1]**



- No unsupervised pre-training → random exploration (green)
- Other parts (red & blue) are same except uses on-policy RL (PPO [2]) for updating policy

[1] Christiano, P., Leike, J., Brown, T.B., Martic, M., Legg, S. and Amodei, D., Deep reinforcement learning from human preferences. NeurIPS, 2017.

[2] Schulman, J., Wolski, F., Dhariwal, P., Radford, A. and Klimov, O., Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.

# Tested Algorithms

- **PEBBLE (ours)**



- **Preference-PPO [1]**



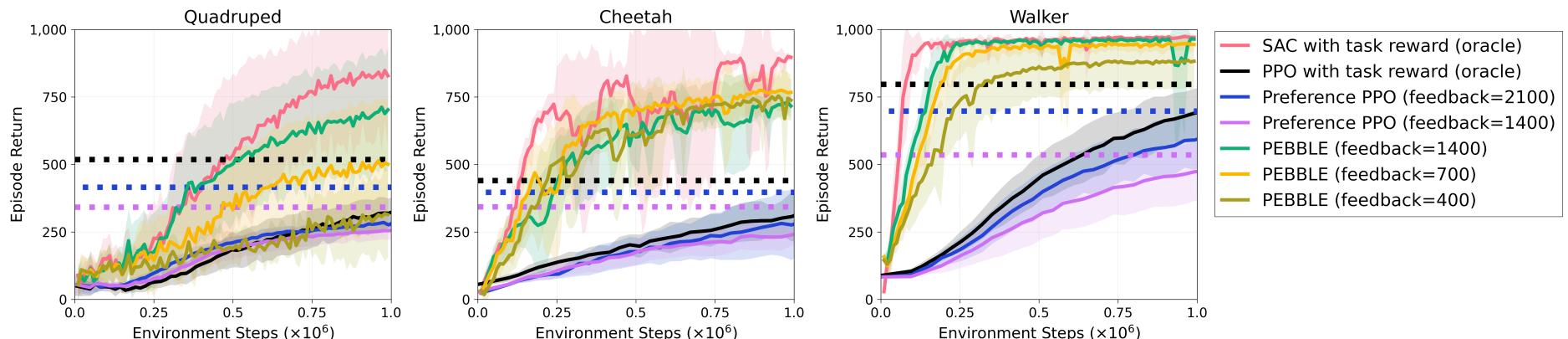
- No unsupervised pre-training → random exploration (green)
- Other parts (red & blue) are same except uses on-policy RL (PPO [2]) for updating policy
- **RL agents trained with true rewards (upper bound)**

[1] Christiano, P., Leike, J., Brown, T.B., Martic, M., Legg, S. and Amodei, D., Deep reinforcement learning from human preferences. NeurIPS, 2017.

[2] Schulman, J., Wolski, F., Dhariwal, P., Radford, A. and Klimov, O., Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.

# Locomotion Tasks

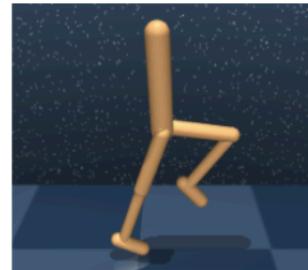
- Learning curves (10 random seeds)



Quadruped



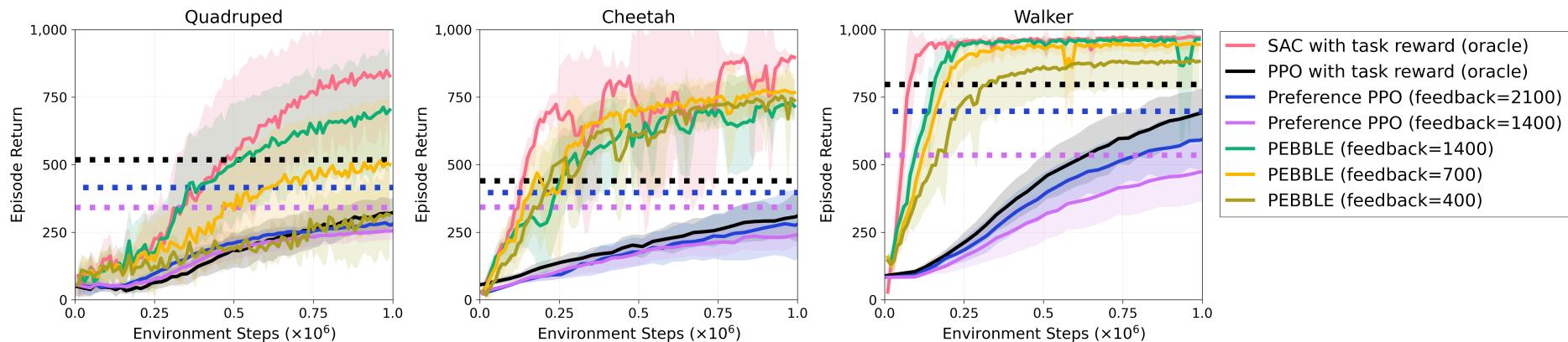
Cheetah



Walker

# Locomotion Tasks

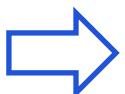
- Learning curves (10 random seeds)



\* Asymptotic performance of PPO and Preference PPO is indicated by dotted lines of the corresponding color

Given a budget of 1400 queries,

- PEBBLE (green) reaches the same performance as SAC (pink)
- Preference PPO (purple) is unable to match PPO (black)



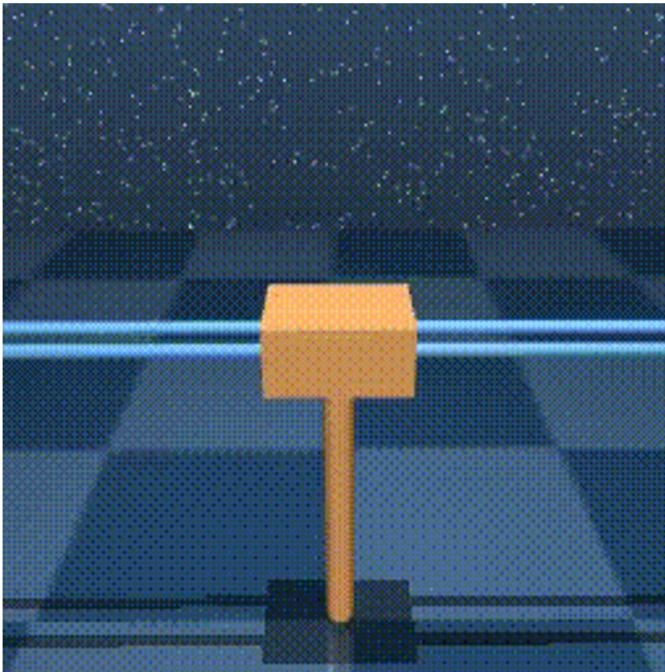
PEBBLE is indeed more feedback-efficient

# Can Human Teach Novel Behaviors?

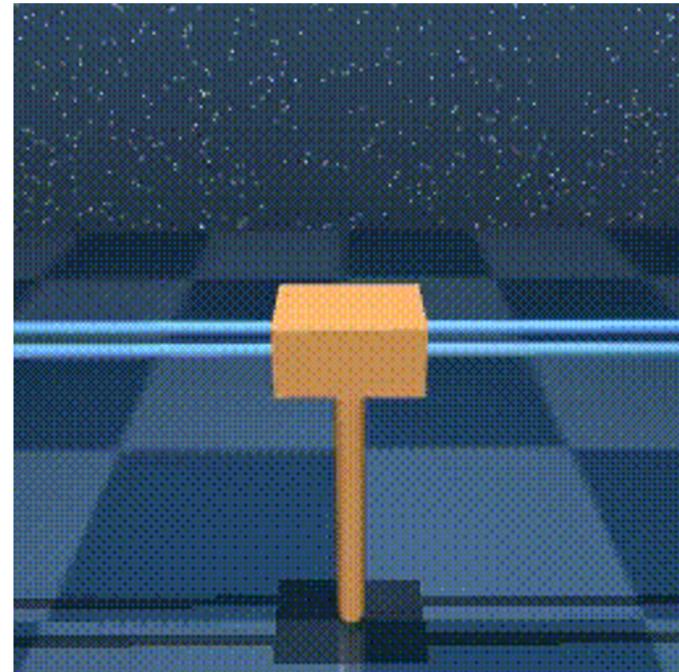
---

# Can Human Teach Novel Behaviors?

- 40 queries in less than 5 mins



Counter clockwise



Clockwise

# Can Human Teach Novel Behaviors?

- 200 queries in less than 30 mins



Waving left front leg

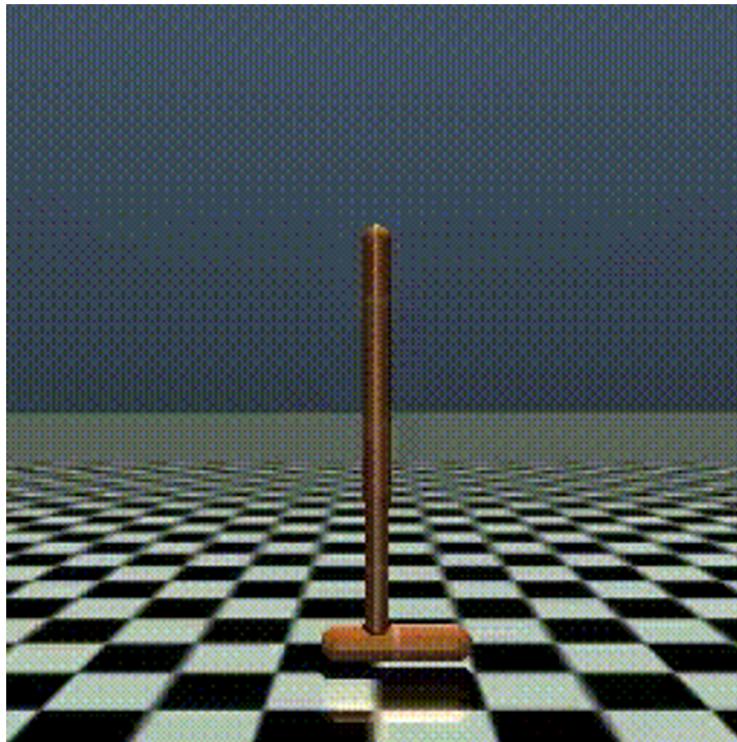


Waving right front leg

# Can Human Teach Novel Behaviors?

---

- 400 queries in less than one hour

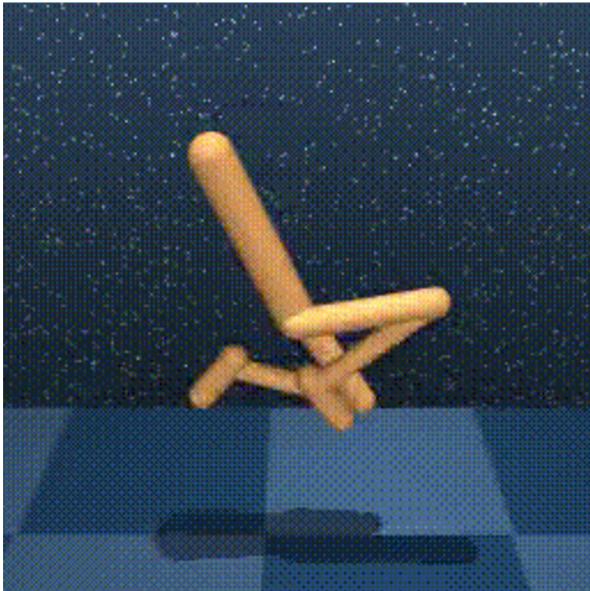


# Can We Avoid Reward Exploitation?

---

# Can We Avoid Reward Exploitation?

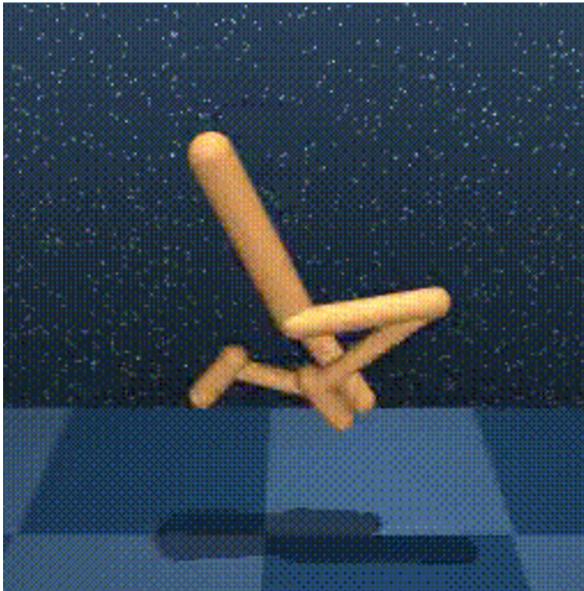
---



SAC with task reward on walker, walk  
**(use one leg even if score ~=1000)**

# Can We Avoid Reward Exploitation?

- 150 queries in less than 20 mins



SAC with task reward on walker, walk  
**(use one leg even if score  $\approx 1000$ )**



SAC trained with human feedback  
**(use both legs)**

# Summary

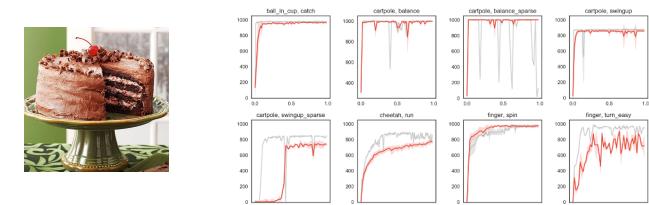
## ■ NN Architectures

Pretrained Transformers as Universal Computation Engines  
Kevin Lu, Aditya Grover, Pieter Abbeel, Igor Mordatch



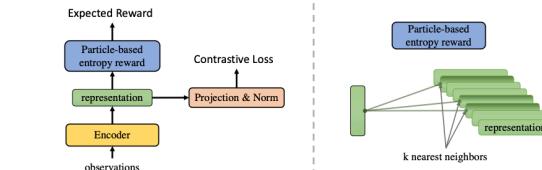
## ■ Representation Learning for Reinforcement Learning

CURL: Contrastive Unsupervised Representation Learning for RL  
Misha Laskin\*, Aravind Srinivas\*, Pieter Abbeel



## ■ Unsupervised Skill Discovery

Behavior from the Void: Unsupervised Active Pre-Training  
Hao Liu, Pieter Abbeel



## ■ Human-in-the-loop Reinforcement Learning

PEBBLE: Feedback-Efficient Interactive RL via Relabeling Experience and Unsupervised Pre-training  
Kimin Lee\*, Laura Smith\*, Pieter Abbeel

