



# Deep Reinforcement Learning

## Lecture 8: Learning from Demonstrations

---

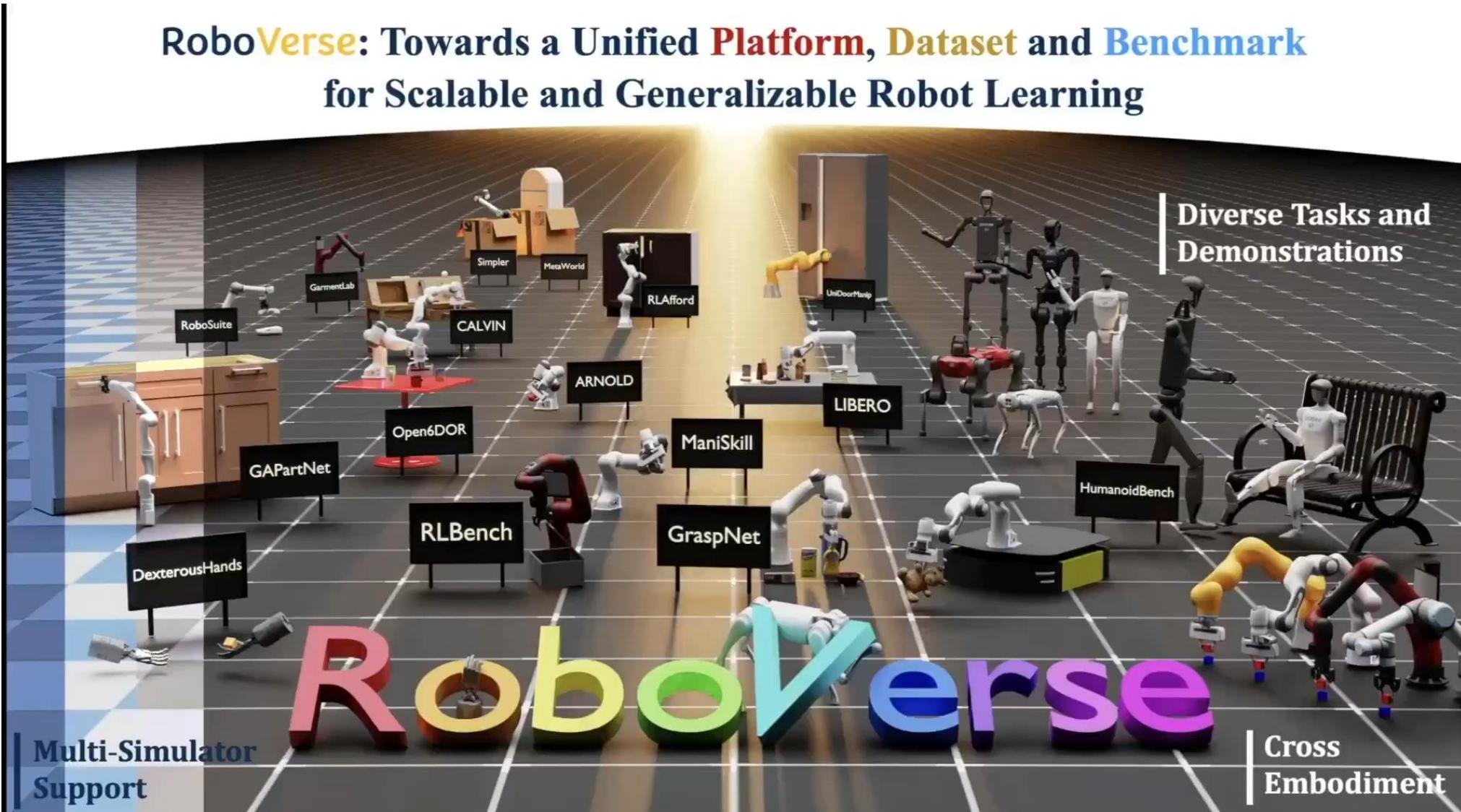
Huazhe Xu

Tsinghua University

# Project & Quiz

- This weekend!
- Midterm Quiz in week 10!! In this room! The same time as our lecture.
  - A calculator
  - A pen
  - Student ID Card
  - I'll bring you some chocolate ;p
  - Midterm quiz sample questions will be released next week!

# AI This Week



Previously, we have learned reinforcement learning…

- Data is collected from
  - the agent itself
  - some past experience of some agent
- However, in many tasks, we might have different quality of data.
  - For example, when you learn to play tennis, your coach will show you the correct motion.
  - Another example is your coach might also demonstrate how you do it wrong!
- In this lecture, we will learn how can we learn from some good data or demonstrations! (And we will leave dealing with bad/suboptimal data to offline RL.)



# In Lec8

- 1 Imitation Learning
- 2 Improving Imitation Learning
- 3 Other settings in Learning from Demonstration



# In Lec8

- 1 Imitation Learning
- 2 Improving Imitation Learning
- 3 Special settings in Learning from Demonstration



# Reward? Interaction steps?



# Imitation Learning

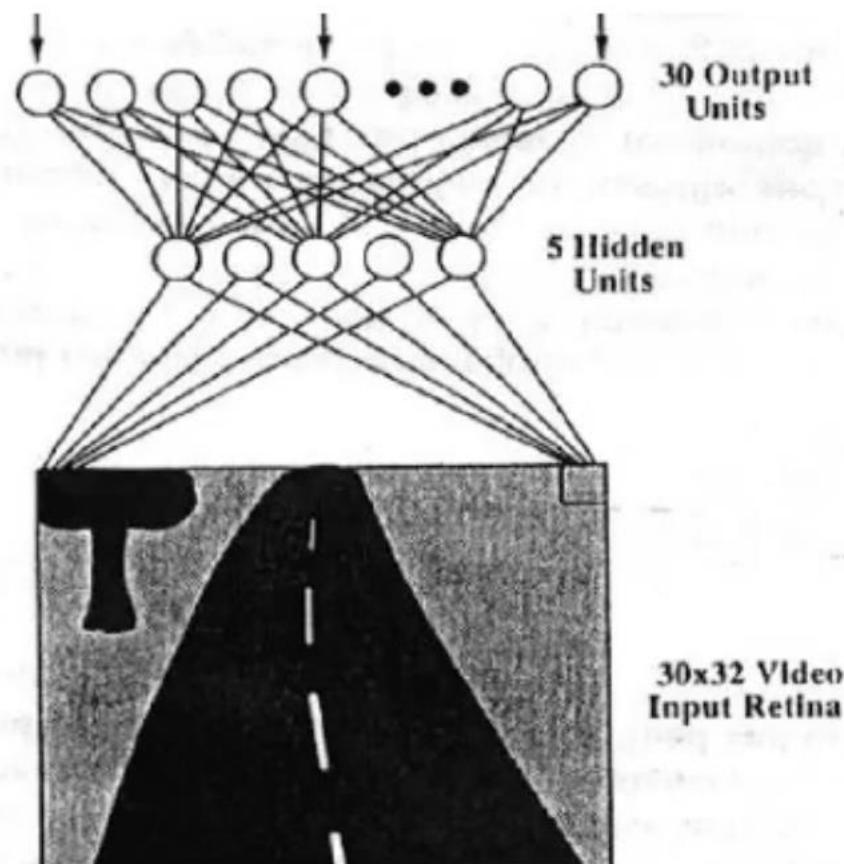
- We have many demonstrations in the format of  $(o_t, a_t)$ 
  - $o_t$  is the observation
  - $a_t$  is the action
- Our goal is to learn a policy  $\pi(a_t|o_t)$  **that maximize the likelihood or some measure of similarity** between the predicted actions and the actions from demonstrations.
- This looks familiar, right?
  - Very similar to other tasks with a supervised learning paradigm
  - E.g., image classification
  - House price regression
- If you directly minimize the loss between predicted actions and demonstration actions, this is also called **behavior cloning (BC)**.

# Imitation Learning

- Can the agent interact with the environment?
  - Yes! Later, you'll see some RL-based imitation learning.
- Can the expert (data collector) interact with environment?
  - Usually no.
  - Or we call it human(expert)-in-the-loop.

# History: Deep Imitation Learning for Driving

- ALVINN: Autonomous Land Vehicle In a Neural Network



# History: Deep Imitation Learning for Driving



ALVINN Approach

Source: Youtuber Welch Labs



ALVINN Results

# Can imitation learning actually work?

- No, in most of the cases!
  - Self-driving is not solved even today.
- Reason 1: it is similar to how we describe the accumulated error in training dynamics models. The small error will ultimately become large issues.
- Reason 2: what is even worse, once the agent goes into unknown areas, it does not know how to solve the task nor recover from the unknown area to the known area. (covariate shift)
  - More formally,  $p_{data}(\mathbf{o}_t) \neq p_{\pi_\theta}(\mathbf{o}_t)$  ! Then your policy lose control.

# General Imitation Learning

- Objective

$$\pi^* = \arg \min_{\pi} \mathbb{E}_{s \sim \rho_{\pi}^s} [\ell(\pi(\cdot \mid s), \pi_E(\cdot \mid s))]$$

- $\ell$  denotes some loss function or some distance metric.
- Distribution of  $s$  depends on rollout from  $\pi$ .
- $P(s \mid \pi) \rightarrow \rho_{\pi}(s)$  : distribution of states sampled by a policy

$$\rho_{\pi}(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s \mid \pi)$$

- Problem
- Cannot get access to the expert during sampling!

# As always, can we improve imitation learning?

- Since it is in a supervised learning scheme for now, we might want to go through the whole pipeline and improve IL by
  - Collecting **better dataset** in some smart way
  - Adjust your **data input** so that your model can make better decision
  - Leverage **better models** to produce **useful representations** for downstream tasks
  - Design **better loss function/training methods**
  - Smartly design **downstream tasks** to make the target task look easier
- We will go through this one by one.



# In Lec8

- 1 Imitation Learning
- 2 Improving Imitation Learning
- 3 Special settings in Learning from Demonstration

# IL improvements

- **Collecting better dataset in some smart way (data)**
- Adjust your data input so that your model can make better decision (inputs)
- Leverage better models to produce useful representations for downstream tasks (model)
- Design better loss function/training methods (training & loss)

# IL improvements techniques: DAgger

- Can we make the demonstration data distribution match the policy-induced data distribution?
- DAgger: Dataset Aggregation
  1. train  $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$  from human data  $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
  2. run  $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$  to get dataset  $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
  3. Ask human to label  $\mathcal{D}_\pi$  with actions  $\mathbf{a}_t$
  4. Aggregate:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

# DAgger original paper and its follow-up SMILe

---

## **A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning**

---

**Stéphane Ross**  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213, USA  
stephaneross@cmu.edu

**Geoffrey J. Gordon**  
Machine Learning Department  
Carnegie Mellon University  
Pittsburgh, PA 15213, USA  
ggordon@cs.cmu.edu

**J. Andrew Bagnell**  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213, USA  
dbagnell@ri.cmu.edu

---

## **Efficient Reductions for Imitation Learning**

---

**Stéphane Ross**  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA, USA

**J. Andrew Bagnell**  
Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA, USA

# A comparison



Expert



w/o DAgger



w/ DAgger

# What is not practical in DAgger?

- It might be hard for human to get the actions by looking at the observations. Human might need feedback after executing an tentative action.
- To collect these new data, we need to execute a suboptimal policy first. That is dangerous.
- In the Dagger in Carla example, everything is on-policy. That means your oracle/autopilot has to be ready to takeover anytime.

# IL improvements techniques: Dataset Resampling/Reweighting

- Your dataset can be imbalanced.
  - For example, in a driving scene, 95% of the data are just driving straight down the road.
- Sample more of the rare data can be helpful.
- Another way is to put more weight on it.
- How do you determine what are the rare samples?
  - A naïve way: use prediction error!

# IL improvements techniques: Data Augmentation

- Augment your data to scale up the dataset! But note there is a difference between IL and computer vision tasks.
- Commonly used data augmentation techniques include:
  - Translation
  - Color jittering
  - Image cropping
  - Black/White pepper noise
- Some data augmentation you should be careful:
  - Image left-right flip: sometimes this is useful if your action can also be flipped.
  - Image masking: you have to mask out pixels smartly.
- Rarely used data augmentation
  - Image rotation

# Better Data Collection Techniques

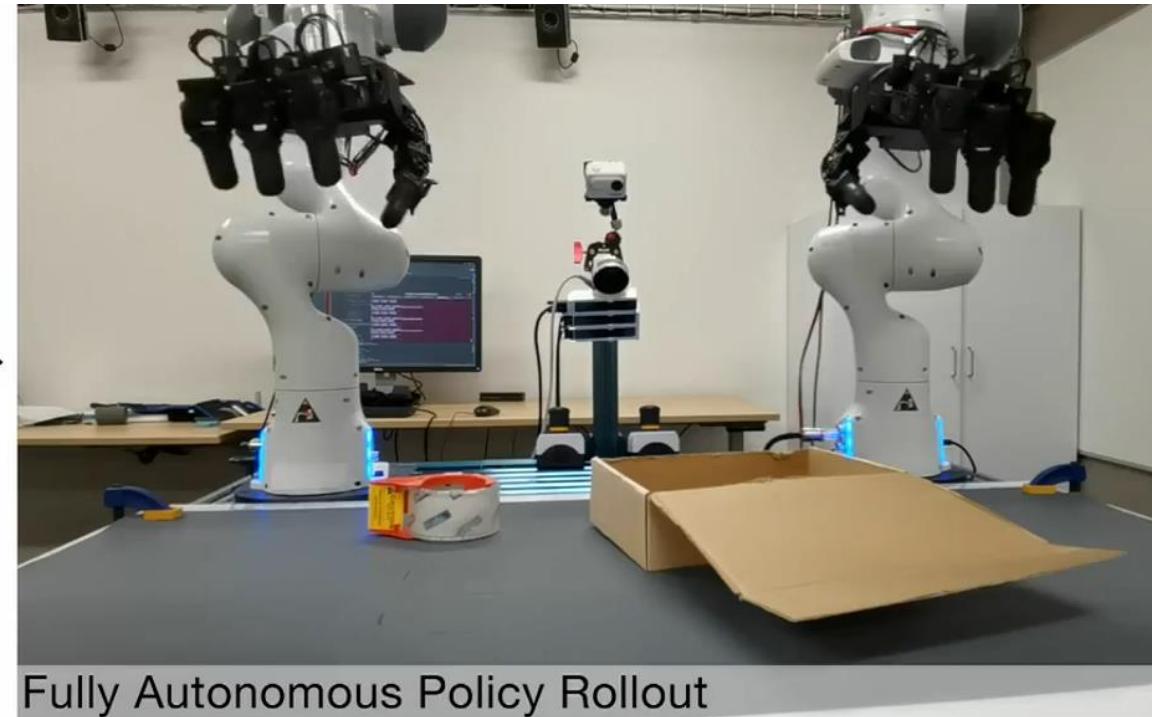


Mobile Aloha from Chelsea Lab

# Better Data Collection Techniques



**Hand motion capture + Point cloud**

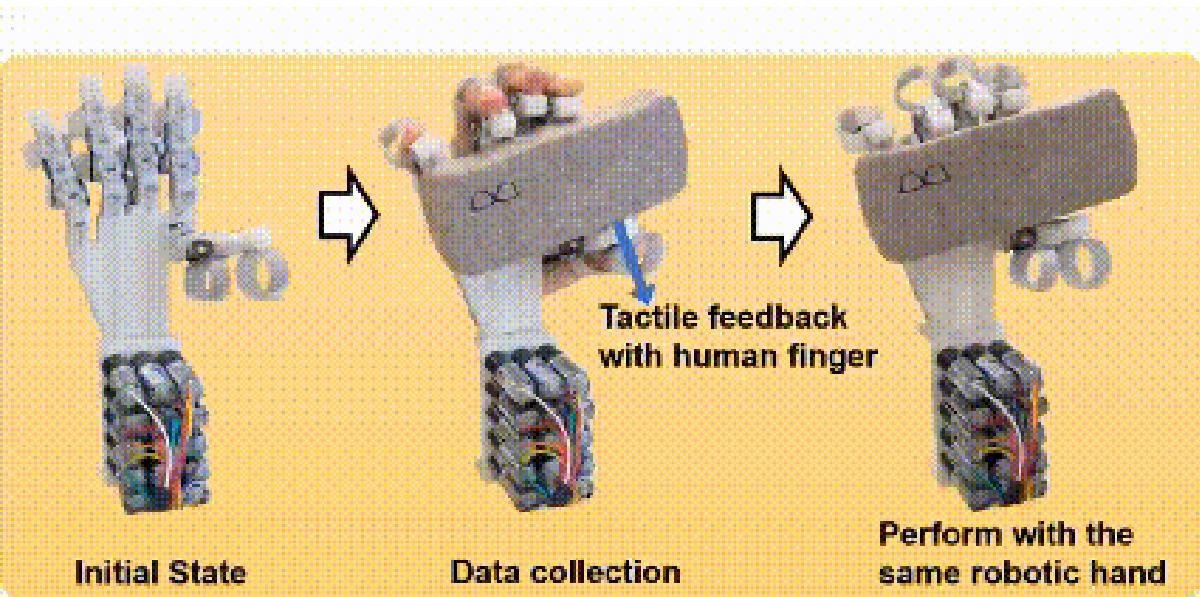


**Fully Autonomous Policy Rollout**

DexCap from Stanford Fei-Fei lab



AirExo from SJTU Cewu's Lab



HiroHand from TEA lab



UMI from Stanford Shuran Lab

# IL improvements

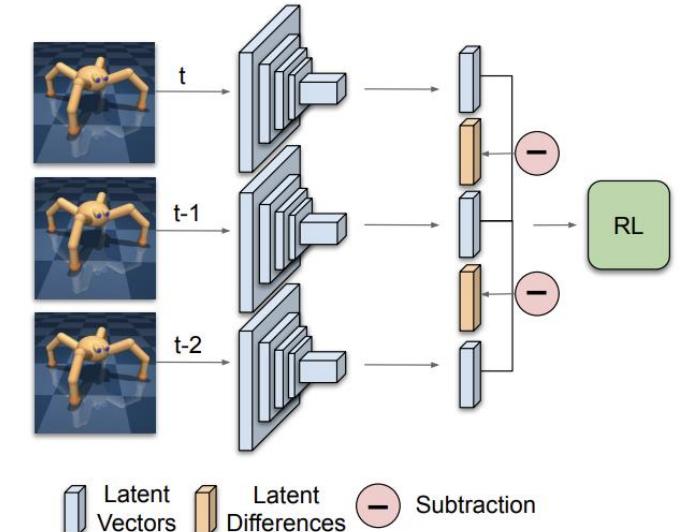
- Collecting better dataset in some smart way (data)
- **Adjust your data input so that your model can make better decision (inputs)**
- Leverage better models to produce useful representations for downstream tasks (model)
- Design better loss function/training methods (training & loss)

# IL improvement techniques: Putting histories into your observation.

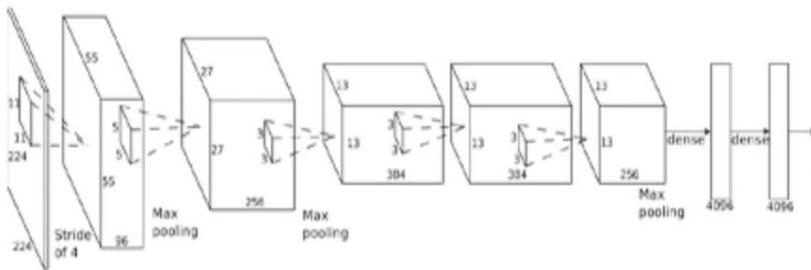
- Instead of using a single image, you stack multiple images into your observation. (Familiar, right? Same thing as in DQN!)
- The task might be non-markovian, we want the history to help.
- Other ways:
  - Use sequence models: RNN or Transformer.
  - Flare: use concatenated difference between frames (works in RL! But can also be applied to IL.)

## REINFORCEMENT LEARNING WITH LATENT FLOW

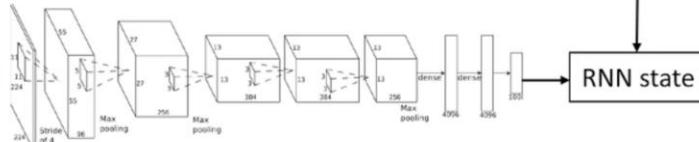
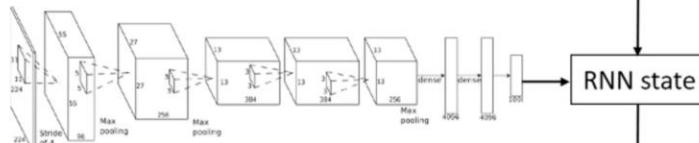
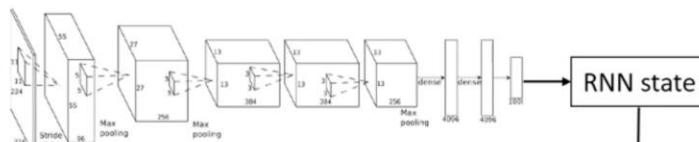
Wenling Shang<sup>2,†</sup>, Xiaofei Wang<sup>1,†</sup>,  
Aravind Srinivas<sup>1</sup>, Aravind Rajeswaran<sup>3</sup>, Yang Gao<sup>1</sup>,  
Pieter Abbeel<sup>1</sup> & Michael Laskin<sup>1</sup>  
University of California Berkeley<sup>1</sup>, Deepmind<sup>2</sup>, University of Washington<sup>3</sup>



# Illustration of Stacked Frames and RNN-Based Approach



Stacked Frames

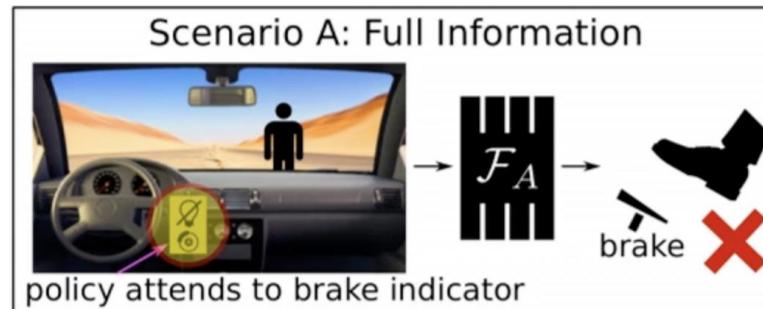


Using RNN to deal with history

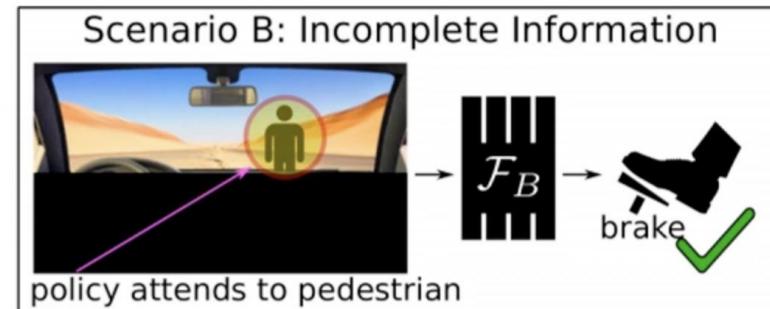
Adapted from CS 285 Levine's lecture

# Using histories can have causal confusions.

- The driving example
- Scenario A: whenever a driver steps on the brake, the light will be lit up. If use history observations, the policy would just check the light instead of the pedestrian in front of the car.
- Scenario B: The light is broken. The policy would attend to the pedestrian.



“causal confusion”



see: de Haan et al., “Causal Confusion in Imitation Learning”

# Papers addressing causal confusion

---

## **Fighting Copycat Agents in Behavioral Cloning from Observation Histories**

---

**Chuan Wen<sup>\*1</sup>, Jierui Lin<sup>\*2</sup>, Trevor Darrell<sup>2</sup>, Dinesh Jayaraman<sup>3</sup>, Yang Gao<sup>†124</sup>**

<sup>1</sup>Institute for Interdisciplinary Information Sciences, Tsinghua University

<sup>2</sup>UC Berkeley, <sup>3</sup>University of Pennsylvania, <sup>4</sup>Shanghai Qi Zhi Institute

# IL improvement techniques: Multi-view Inputs

- If you are using images as input, it is partially observable.
- Consider adding more cameras or stereo cameras to provide more 3D information.
- This is very common in robotics and driving.

# IL improvements

- Collecting better dataset in some smart way (data)
- Adjust your data input so that your model can make better decision (inputs)
- **Leverage better models for downstream tasks (model)**
- Design better loss function/training methods (training & loss)

# IL improvement techniques: pre-trained models to extract representations

- This is an active research area and is usually referred as representation learning (for sensorimotor).
- It usually works for both BC and visual RL.

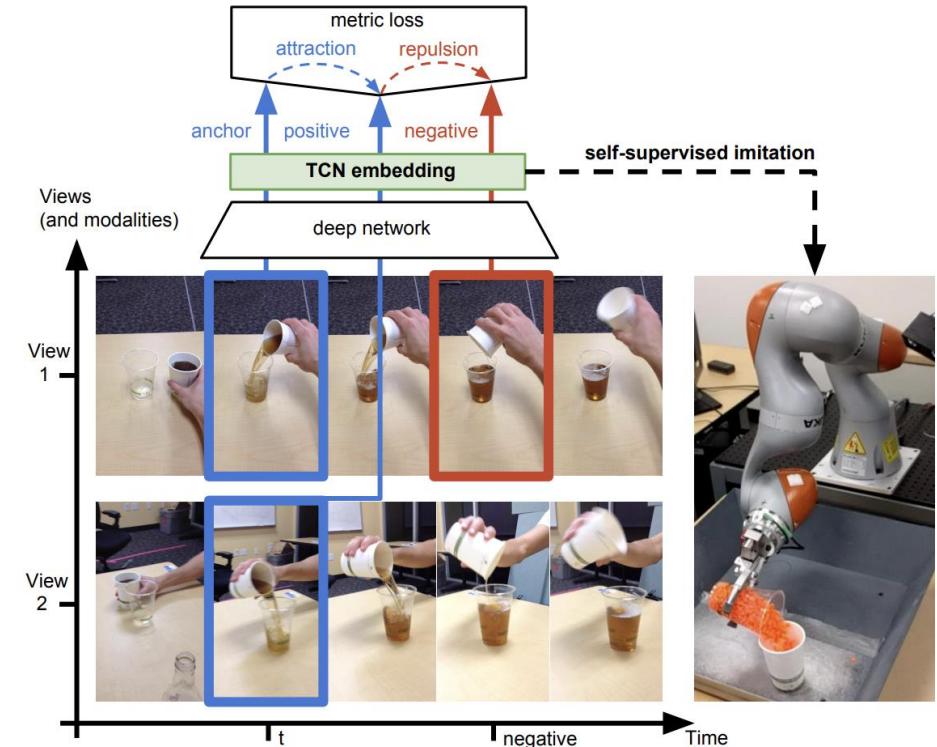
Papers:

- Temporal Contrastive Network (TCN)
- The (Un)Surprising Effectiveness of Pre-Trained Vision Models for Control (PVR)
- R3M: A Universal Visual Representation for Robot Manipulation (R3M)
- Masked Visual Pre-training for Motor Control (MVP)
- Where are we in the search for an Artificial Visual Cortex for Embodied Intelligence? (VC1)

# Time-Contrastive Networks: Self-Supervised Learning from Video (TCN)

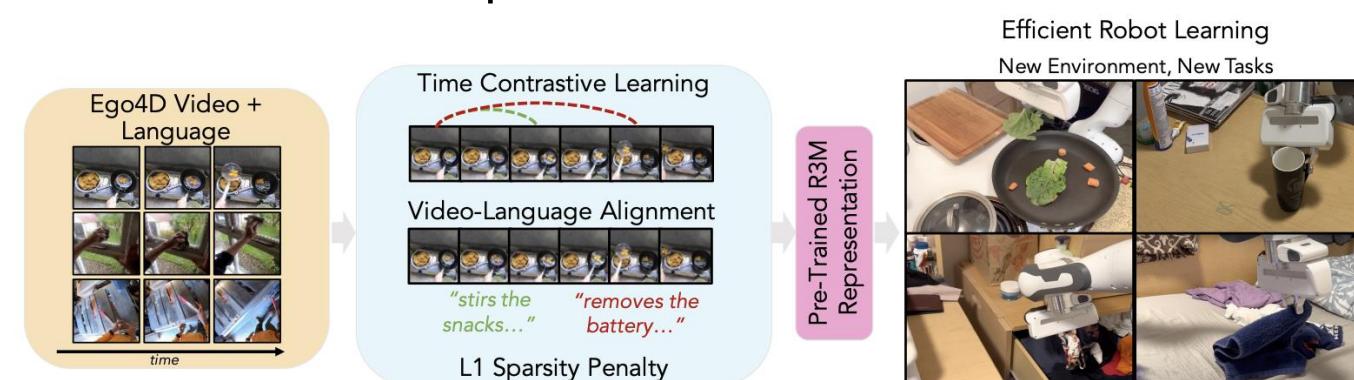
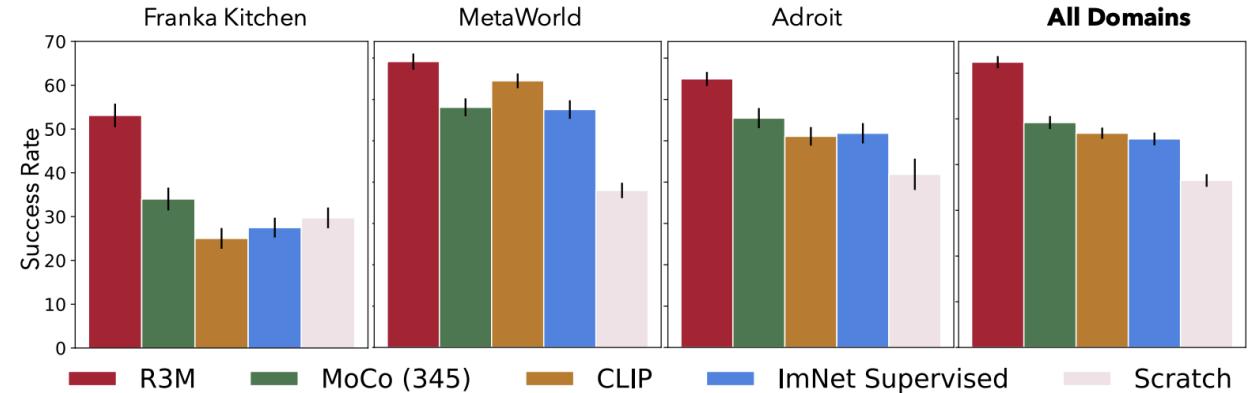
- High-level description:
  1. Record synchronized multi-view videos
  2. Use networks to extract representations
  3. Train the network such that
    1. Temporally close frames' representation are also close to each other.
    2. Temporally distant frames' representation are also far away!

This technique would provide **better representation for robotic tasks!**



# R3M

- High-level description:
  1. Prepare diverse human videos such as Ego4D
  2. Train with 3 losses:
    1. Video-language joint loss for capturing the semantics
    2. TCN loss for capturing temporal information
    3. L1 sparsity loss
  3. When tested on a robot, it works better in performance.



# How to better leverage pre-trained model?

- Since these are recent papers, the conclusion hasn't been tested by time.
- A recent ICML paper from UCSD and THU:

## On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline

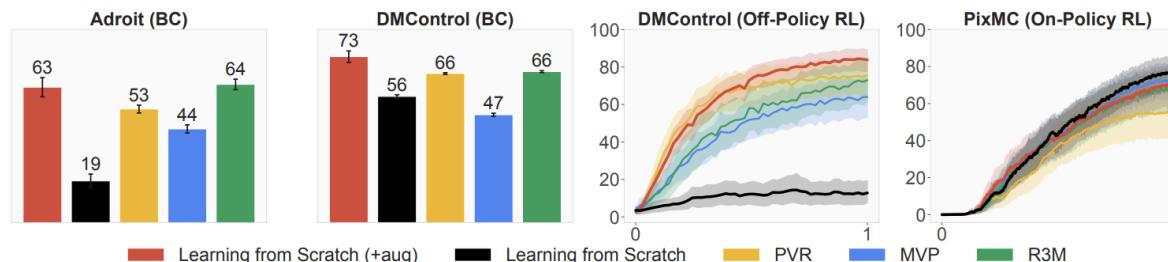
Nicklas Hansen<sup>12\*</sup>, Zhecheng Yuan<sup>3\*</sup>, Yanjie Ze<sup>14\*</sup>, Tongzhou Mu<sup>1\*</sup>,

Aravind Rajeswaran<sup>2†</sup>, Hao Su<sup>1†</sup>, Huazhe Xu<sup>35†</sup>, Xiaolong Wang<sup>1†</sup>

<sup>1</sup>University of California San Diego, <sup>2</sup>Meta AI,

<sup>3</sup>Tsinghua University, <sup>4</sup>Shanghai Jiao Tong University, <sup>5</sup>Shanghai Qi Zhi Institute

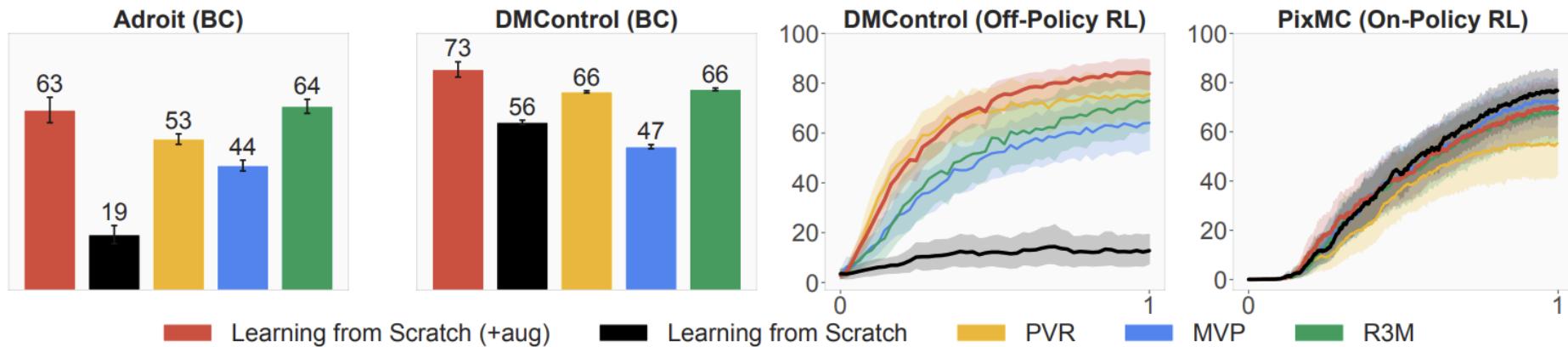
\*Equal contributions, †Equal advising



# On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline

Learning-from-scratch baseline has competitive performance with recent methods that leverage frozen visual representations trained on large-scale vision datasets.

**Which one is correct? That can be a good topic for you to figure out!**



# IL improvement techniques: Pre-training Inverse Dynamics Models

- What is an inverse dynamics model?
  - Given state  $s_t$  and next state  $s_{t+1}$ , predict action  $a_t = f(s_t, s_{t+1})$
  - Imagine what action would help you to jump from current state to the next state.
- What can we do if we have an inverse dynamics model?
  - We can use it to label actions if we only have observations, e.g., in video games or youtube videos!
- VPT:
  - You have a small set of demonstration data in Minecraft game.
  - You have near infinite videos of Minecraft on youtube/bilibili.
  - You train an inverse model on small set and use this model to label large sets.
  - Do IL and the execution of your learned policy can help you to enrich the small dataset.
- I put it into the pre-training schemes. But you can also regard it as a way to collect better data.

# Video PreTraining (VPT)

## Video PreTraining (VPT): Learning to Act by Watching Unlabeled Online Videos

Bowen Baker\*†  
bowen@openai.com

Ilge Akkaya\*†  
ilge@openai.com

Peter Zhokhov\*†  
peterz@openai.com

Joost Huizinga\*†  
joost@openai.com

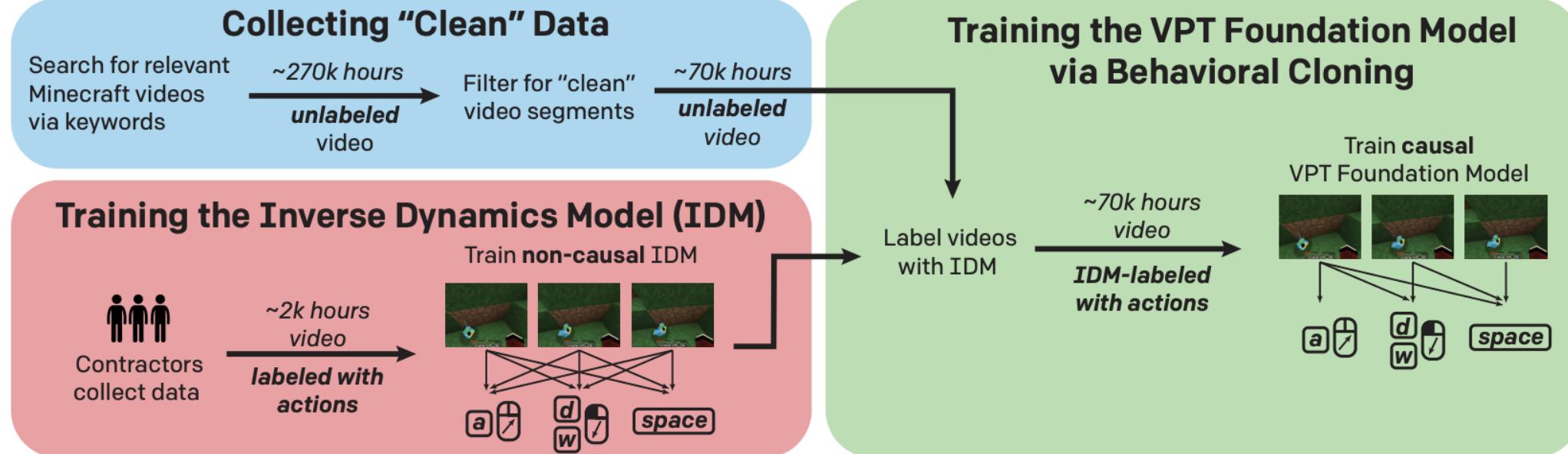
Jie Tang\*†  
jietang@openai.com

Adrien Ecoffet\*†  
adrien@openai.com

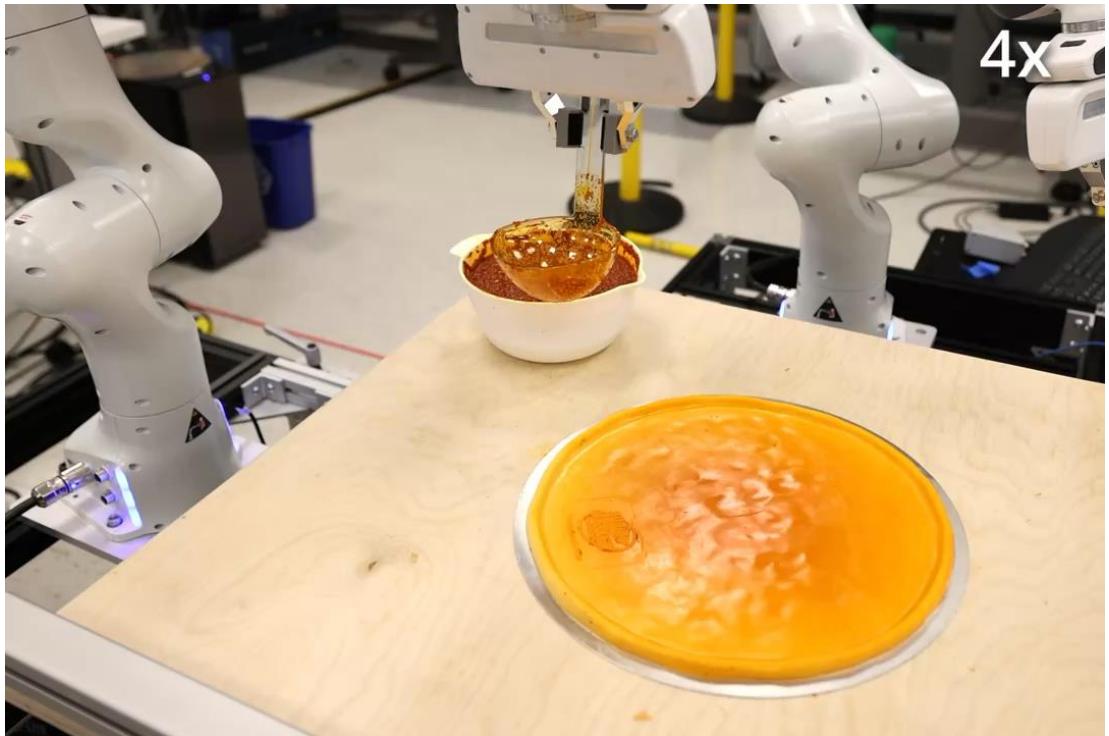
Brandon Houghton\*†  
brandon@openai.com

Raul Sampedro\*†  
raulsamg@gmail.com

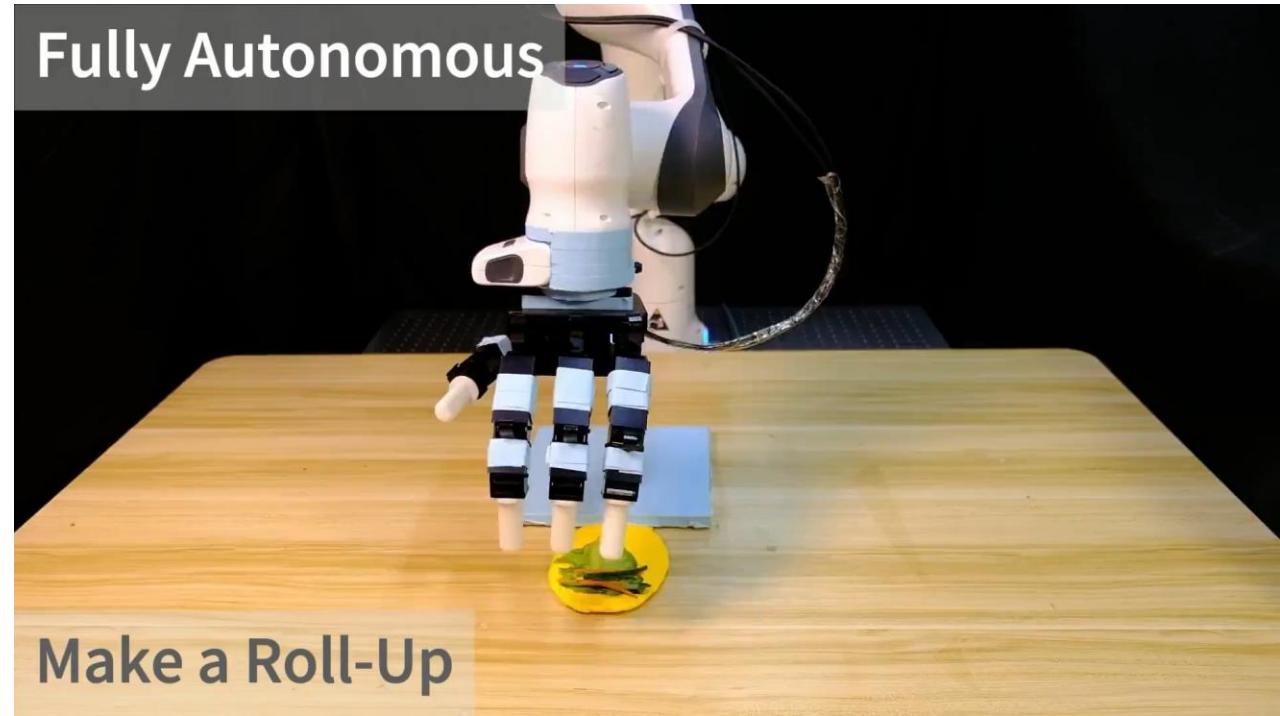
Jeff Clune\*††  
jclune@gmail.com



# Diffusion Policy



DP from Columbia U



DP3 from THU TEA lab

# IL improvements

- Collecting better dataset in some smart way (data)
- Adjust your data input so that your model can make better decision (inputs)
- Leverage better models to produce useful representations for downstream tasks (model)
- **Design better loss function/training methods (training & loss)**

# IL improvement techniques: Multi-modality-Aware Training

- While the dataset can be good (i.e., from expert and accomplish tasks near optimally), the innate nature of human behaviors are multimodal.
- An example:

Xu et al, End-to-end Learning of Driving Models from Large-scale Video Datasets



(c) multiple possible actions:  
turn left or go straight

- Shall we use classification loss like cross-entropy or regression loss like MSE?
  - Cross-entropy! MSE would make your agent going directly toward the obstacle.

# IL improvement techniques: Multi-modality-Aware Training

- What other methods can you use to capture multi-modal behavior?
  - Use Mixture-of-Gaussian instead of gaussian distribution.
  - Learn a policy that is conditioned on a latent variable.
    - An intuitive explanation: the model takes in both an image & some noise, the noise will determine the modality.
    - Not going into details: Variation Autoencoder, Normalizing Flow, Stein Variational Gradient Descent
  - Autoregressive discretization
    - If we have discrete actions, everything is solved!
    - But discretize the actions can be intractable, e.g., your action have high dimensionality.
    - Then you discretize the actions dimension one by one: given an image we output a discretized action dim 1, then input action dim 1 to another neural network to output dim 2.

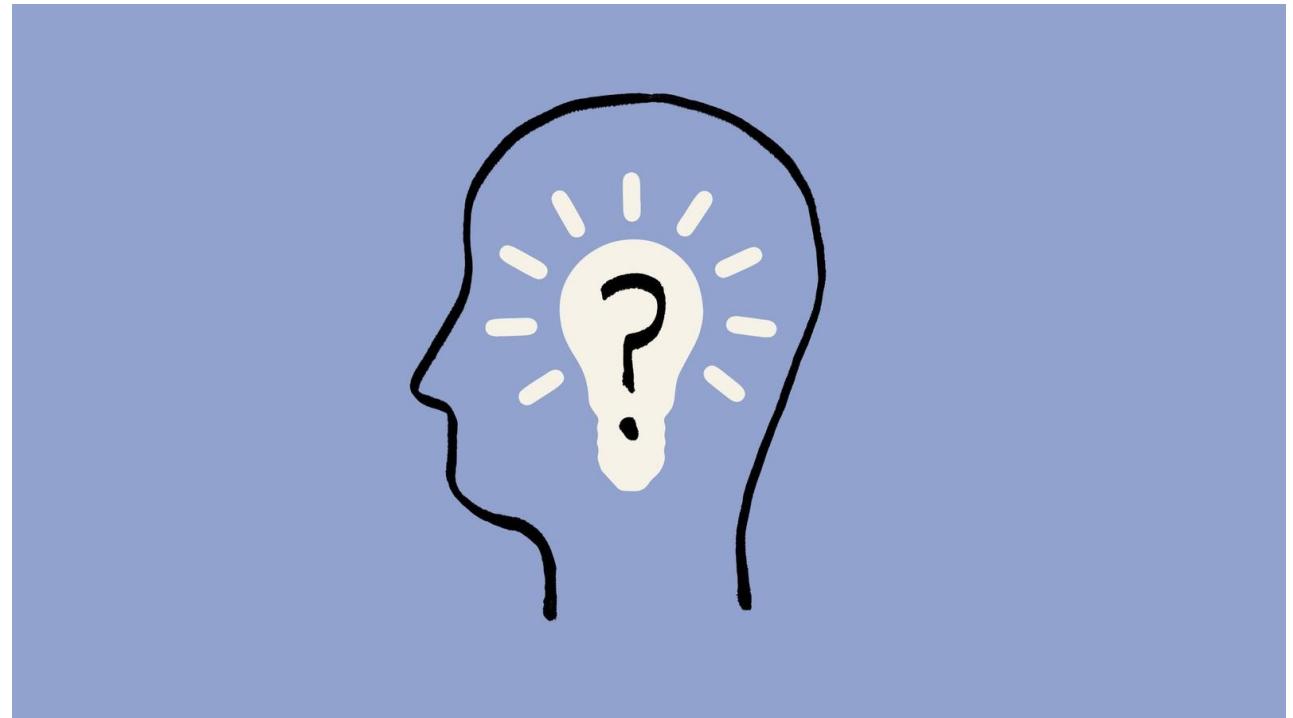


# In Lec10

- 1 Imitation Learning
- 2 Improving Imitation Learning
- 3 Special settings in Learning from Demonstration

Previously, most of the settings are offline.

- What if online interactions are allowed? What can you think of?



# Reinforcement Learning from Demonstrations

- There are multiple ways to use demonstrations.
- Today we introduce Deep Q-learning from Demonstrations (DQfD).

## Deep Q-learning from Demonstrations

### **Todd Hester**

Google DeepMind

[toddhester@google.com](mailto:toddhester@google.com)

### **Matej Vecerik**

Google DeepMind

[matejvecerik@google.com](mailto:matejvecerik@google.com)

### **Olivier Pietquin**

Google DeepMind

[pietquin@google.com](mailto:pietquin@google.com)

### **Marc Lanctot**

Google DeepMind

[lanctot@google.com](mailto:lanctot@google.com)

### **Tom Schaul**

Google DeepMind

[schaul@google.com](mailto:schaul@google.com)

### **Bilal Piot**

Google DeepMind

[piot@google.com](mailto:piot@google.com)

### **Dan Horgan**

Google DeepMind

[horgan@google.com](mailto:horgan@google.com)

### **John Quan**

Google DeepMind

[johnquan@google.com](mailto:johnquan@google.com)

### **Andrew Sendonaris**

Google DeepMind

[sendos@yahoo.com](mailto:sendos@yahoo.com)

### **Ian Osband**

Google DeepMind

[iosband@google.com](mailto:iosband@google.com)

### **Gabriel Dulac-Arnold**

Google DeepMind

[gabe@squirrelsoup.net](mailto:gabe@squirrelsoup.net)

### **John Agapiou**

Google DeepMind

[jagapiou@google.com](mailto:jagapiou@google.com)

### **Joel Z. Leibo**

Google DeepMind

[jzl@google.com](mailto:jzl@google.com)

### **Audrunas Gruslys**

Google DeepMind

[audrunas@google.com](mailto:audrunas@google.com)

# DQfD

- DQfD uses not only the collected experience but also expert demonstrations, i.e., expert demos in the replay buffer.
- At the heart of this approach is this loss:

$$J(Q) = J_{DQ}(Q) + \lambda_1 J_n(Q) + \lambda_2 J_E(Q) + \lambda_3 J_{L2}(Q)$$

- $J_{DQ}$  is the bellman loss,  $J_n$  is the n-step loss, and  $J_{L2}$  is L2 regularization.
- The major change is  $J_E$ . In this loss, with any state, the Q-value of some action  $a$  is at least a margin smaller than Q-value of an expert action.  $a_E$  is the expert action.  $l$  is a hinge loss between  $a_E$  and  $a$ .

$$J_E(Q) = \max_{a \in A} [Q(s, a) + l(a_E, a)] - Q(s, a_E)$$

# Exemplar papers with demonstrations

---

## VRL3: A Data-Driven Framework for Visual Deep Reinforcement Learning

---

Che Wang<sup>1,2\*</sup>

Xufang Luo<sup>3</sup>

Keith Ross<sup>1</sup>

Dongsheng Li<sup>3</sup>

<sup>1</sup> New York University Shanghai

<sup>2</sup> New York University

<sup>3</sup> Microsoft Research Asia, Shanghai, China

## Reinforcement learning with Demonstrations from Mismatched Task under Sparse Reward

Yanjiang Guo<sup>1</sup>, Jingyue Gao<sup>1</sup>, Zheng Wu<sup>2</sup>, Chengming Shi<sup>1</sup>, Jianyu Chen<sup>1,3,\*</sup>

<sup>1</sup>Institute for Interdisciplinary Information Sciences, Tsinghua University

<sup>2</sup>University of California, Berkeley

<sup>3</sup>Shanghai Qizhi Institute

{guoyj22, gaojy19, shicm19}@mails.tsinghua.edu.cn, zheng\_wu@berkeley.edu

\* Correspondence to: Jianyu Chen (jianyuchen@tsinghua.edu.cn)

## Learning Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations

Aravind Rajeswaran<sup>1\*</sup>, Vikash Kumar<sup>1,2\*</sup>, Abhishek Gupta<sup>3</sup>, Giulia Vezzani<sup>4</sup>,  
John Schulman<sup>2</sup>, Emanuel Todorov<sup>1</sup>, Sergey Levine<sup>3</sup>

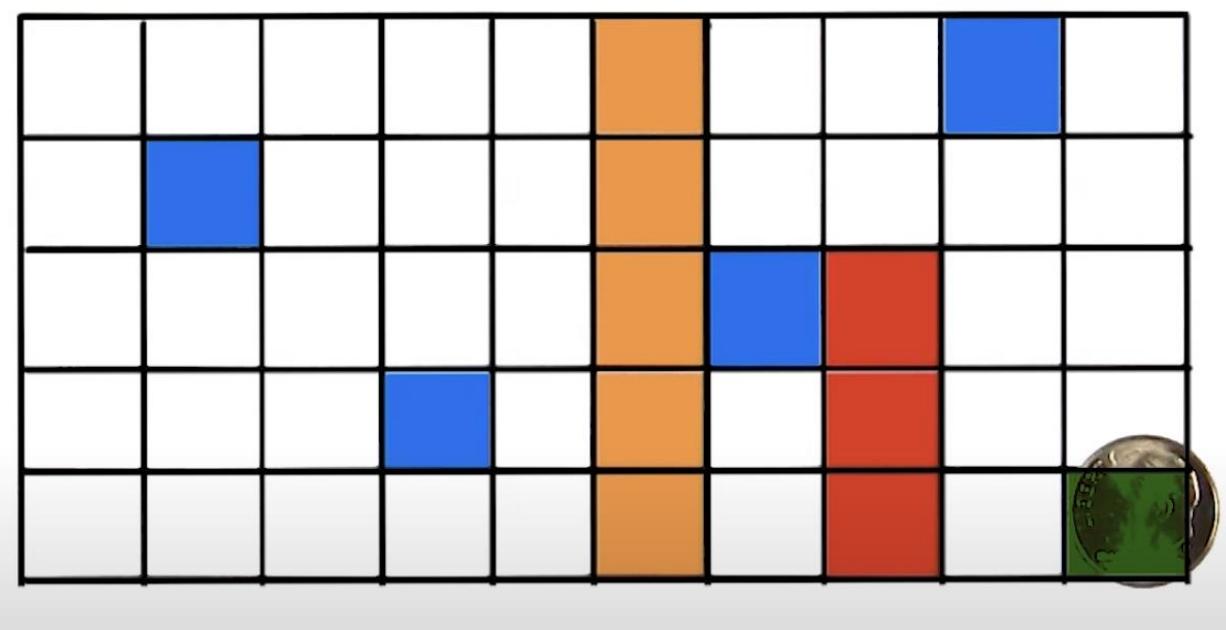
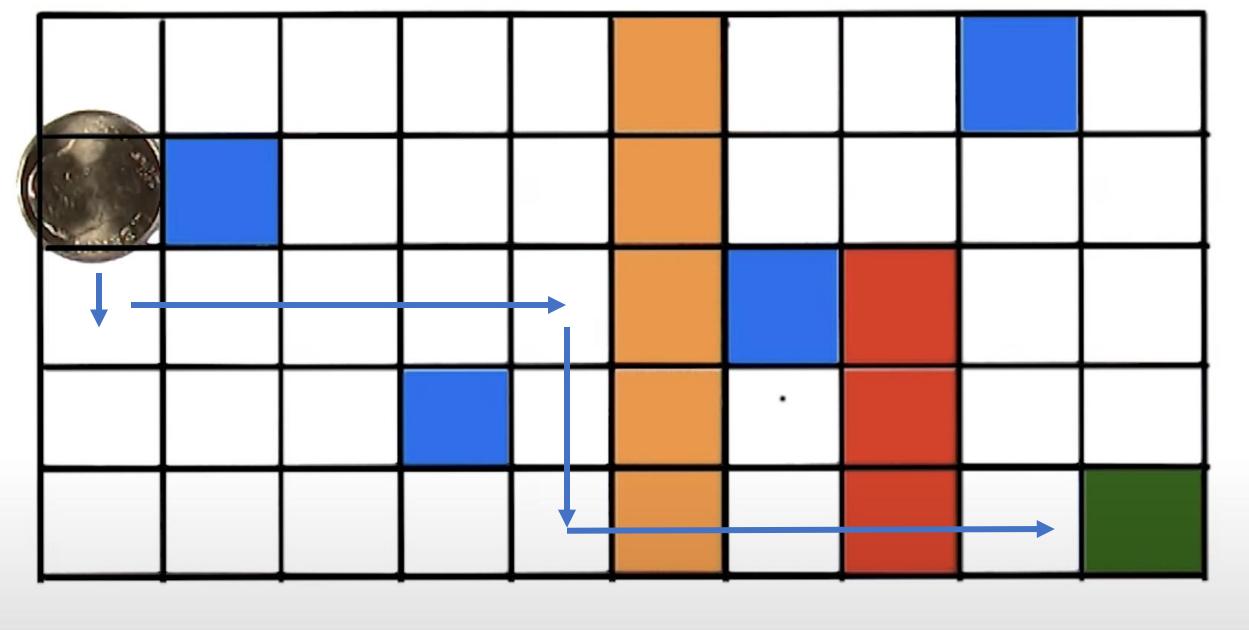
# Other Settings in Imitation Learning

- Learning from observations where actions are missing
- Model-based imitation learning where dynamics model is provided or learned
- Third-person imitation learning
  - Just like how humans learn from their parents/teachers, we don't always learn things from a first person view.
  - Moreover, your teachers appearance and shape might be very different from you. But you can still learn from him/her.
- Learning from suboptimal demonstrations
- Learning from instructions rather than demonstrations

# Another way to learn from demonstrations: Inverse RL

- Instead of learning policies from demonstrations
- Can we learn about the reward function?
- Can we then use these reward function to train an RL agent?
- We will learn these in future lectures if time permits ☺

# Understand IRL through an Example



- What can you infer?

# Why do we need inverse RL?

---

- For science!
  - Model animal and human behaviors
- Imitation
  - Presupposition: reward function provides the most succinct and transferable definition of the task
  - Modeling of other agents, both adversarial and cooperative

---

## Dynamic Inverse Reinforcement Learning for Characterizing Animal Behavior

---

**Zoe C. Ashwood**<sup>1,2,\*</sup>   **Aditi Jha**<sup>1,3,\*</sup>   **Jonathan W. Pillow**<sup>1</sup>

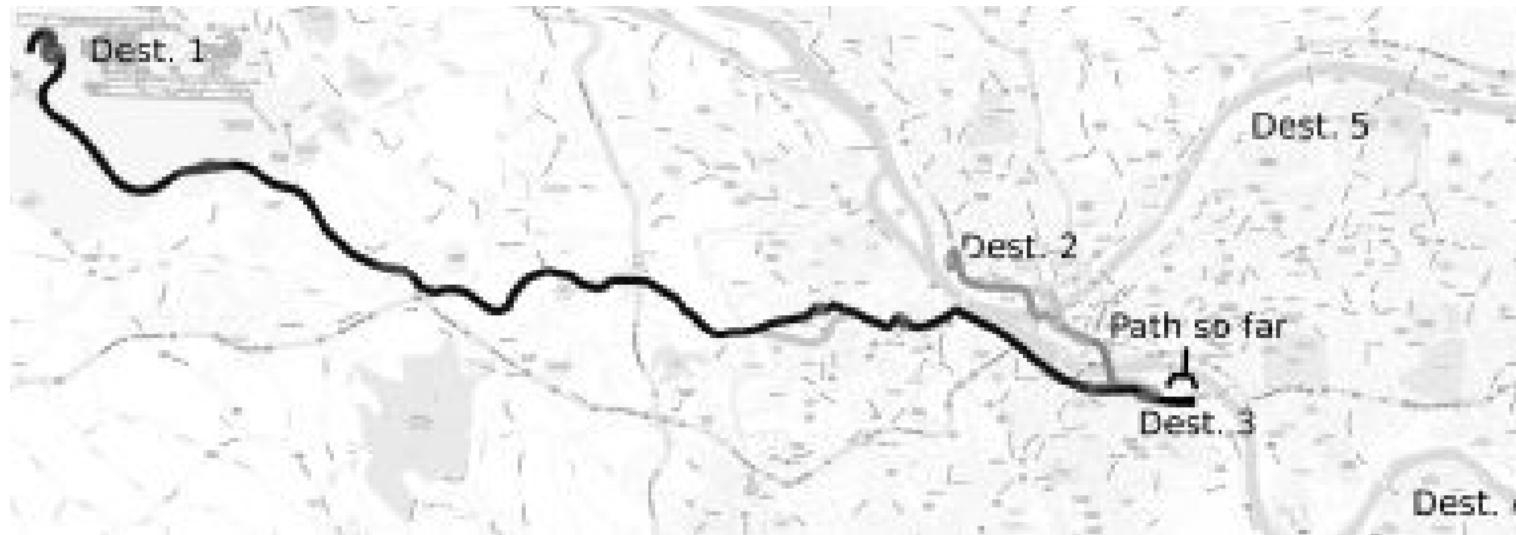
<sup>1</sup>Princeton Neuroscience Institute, Princeton University

<sup>2</sup>Dept. of Computer Science, Princeton University

<sup>3</sup>Dept. of Electrical and Computer Engineering, Princeton University  
`{zashwood, aditijha, pillow}@princeton.edu`

# IRL in Urban Navigation

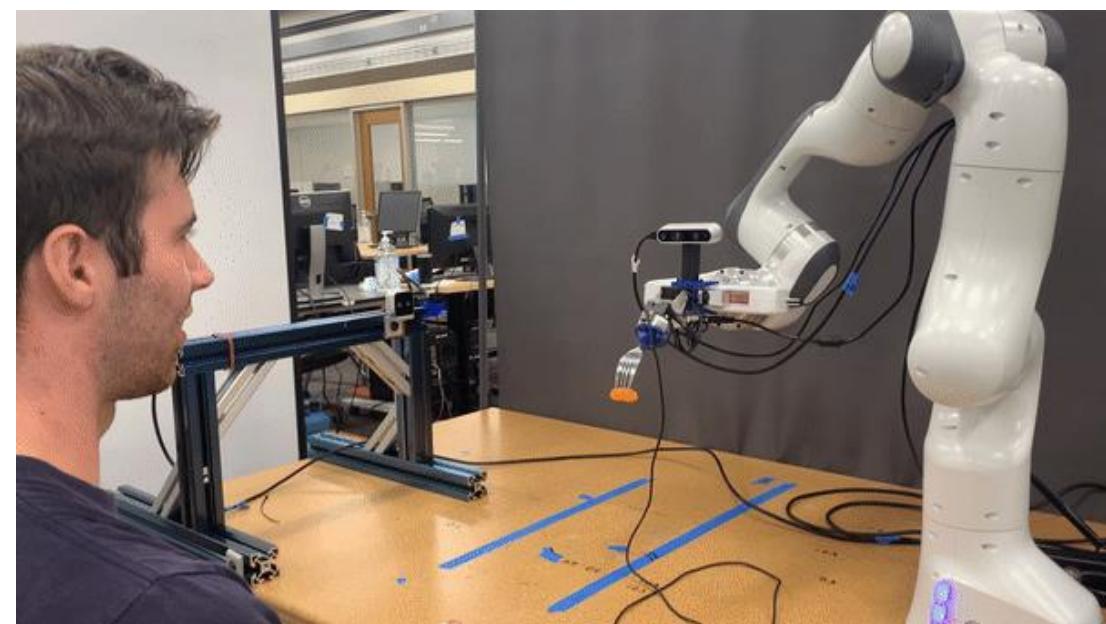
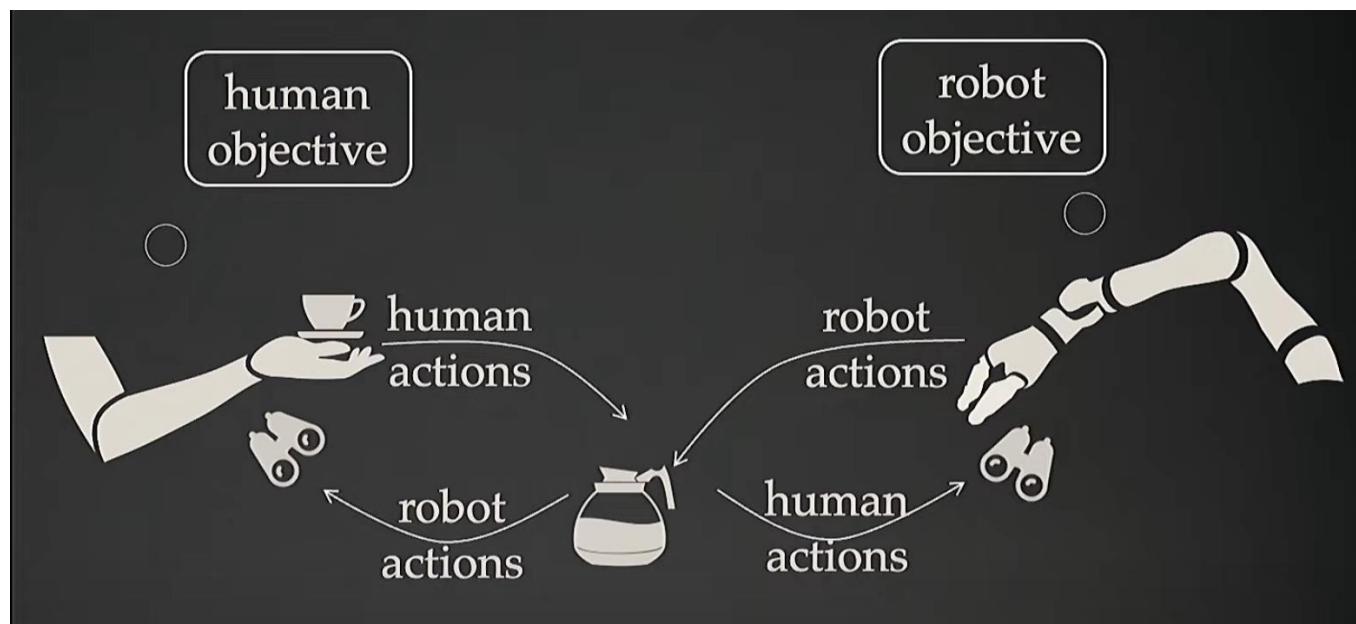
- Infer destination based of partial navigation trajectory



## Maximum Entropy Inverse Reinforcement Learning

**Brian D. Ziebart, Andrew Maas, J. Andrew Bagnell, and Anind K. Dey**  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213  
[bziebart@cs.cmu.edu](mailto:bziebart@cs.cmu.edu), [amaas@andrew.cmu.edu](mailto:amaas@andrew.cmu.edu), [dbagnell@ri.cmu.edu](mailto:dbagnell@ri.cmu.edu), [anind@cs.cmu.edu](mailto:anind@cs.cmu.edu)

# IRL in Human-Robot Interaction (HRI)



# IRL v.s. BC

Which has the most succinct description:  $\pi^*$  v.s.  $R^*$ ?

- It depends.
  - But for RL and planning, reward can be better.
  - From a human video, it is hard to tell what is the optimal policy but easier to tell about the reward.

# Basic Principle

- Find a reward function  $R^*$  which explains the expert behavior.
- Find  $R^*$  such that

$$\mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid \pi^* \right] \geq \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid \pi \right] \quad \forall \pi$$

- Challenges:
  - $R=0$  is always a solution. Reward function ambiguity!
  - We observe trajectories rather than policies. How to compute the left-hand side?
  - Expert has to be optimal.
  - It is not possible to enumerate policies.

# Feature-Based Reward Function

Let  $R(s) = w^\top \phi(s)$ , where  $w \in \mathfrak{R}^n$ , and  $\phi : S \rightarrow \mathfrak{R}^n$ .

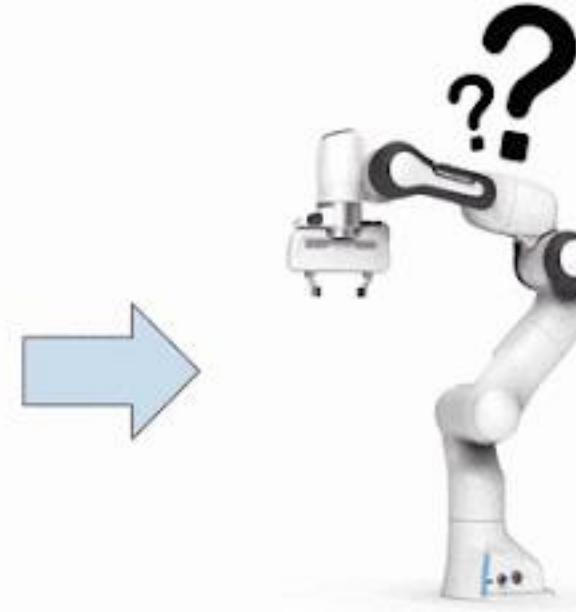
$$\begin{aligned} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi \right] &= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t w^\top \phi(s_t) \mid \pi \right] \\ &= w^\top \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \phi(s_t) \mid \pi \right] \\ &= w^\top \mu(\pi) \end{aligned}$$

- $\mu$  is the feature expectation.
- Recall that  $\mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid \pi^* \right] \geq \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid \pi \right] \quad \forall \pi$

Find  $w^*$  such that  $w^{*\top} \mu(\pi^*) \geq w^{*\top} \mu(\pi) \quad \forall \pi$

# What can you do if demos are videos?

- Follow the idea of IRL
  - Reward learning from videos!
  - How can you achieve that?



# What are the challenges?

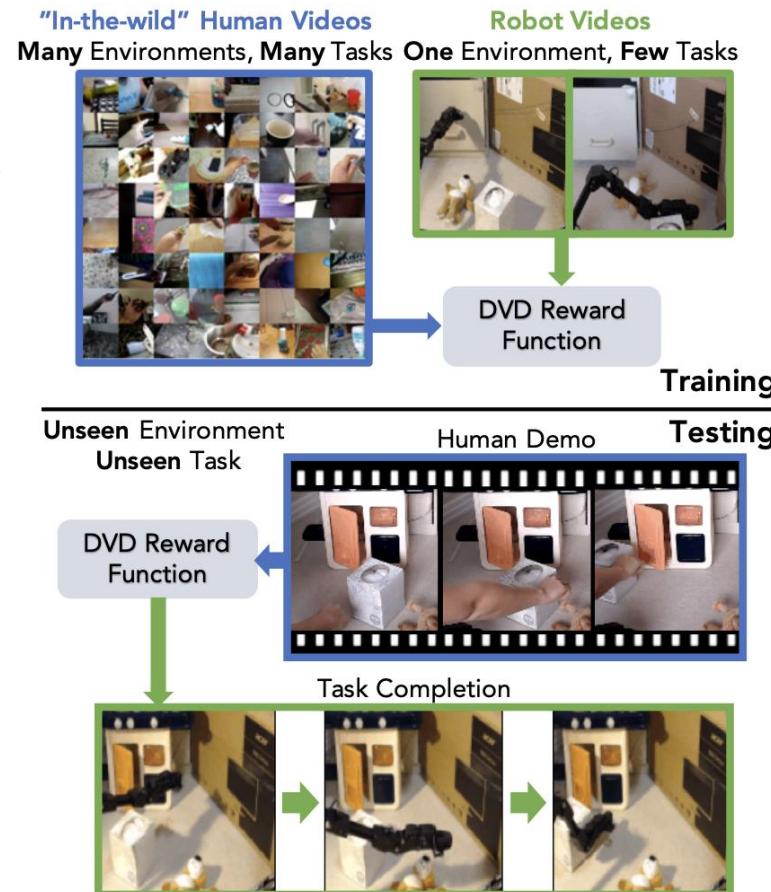
- Different domain!



# Learning Generalizable Robotic Reward Functions from “In-The-Wild” Human Videos

Annie S. Chen, Suraj Nair, Chelsea Finn  
Stanford University

- Domain-agnostic Discriminator
  - A neural network that learns similarity Among tasks.
    - + Positive samples, same task
    - Negative samples, different taskTasks are from human or robot.
- Video Predictor
  - Serves as a dynamics model



# Learning Generalizable Robotic Reward Functions from “In-The-Wild” Human Videos

Annie S. Chen, Suraj Nair, Chelsea Finn  
Stanford University

- Algorithm

---

## Algorithm 1 DOMAIN-AGNOSTIC VIDEO DISCRIMINATOR (DVD)

---

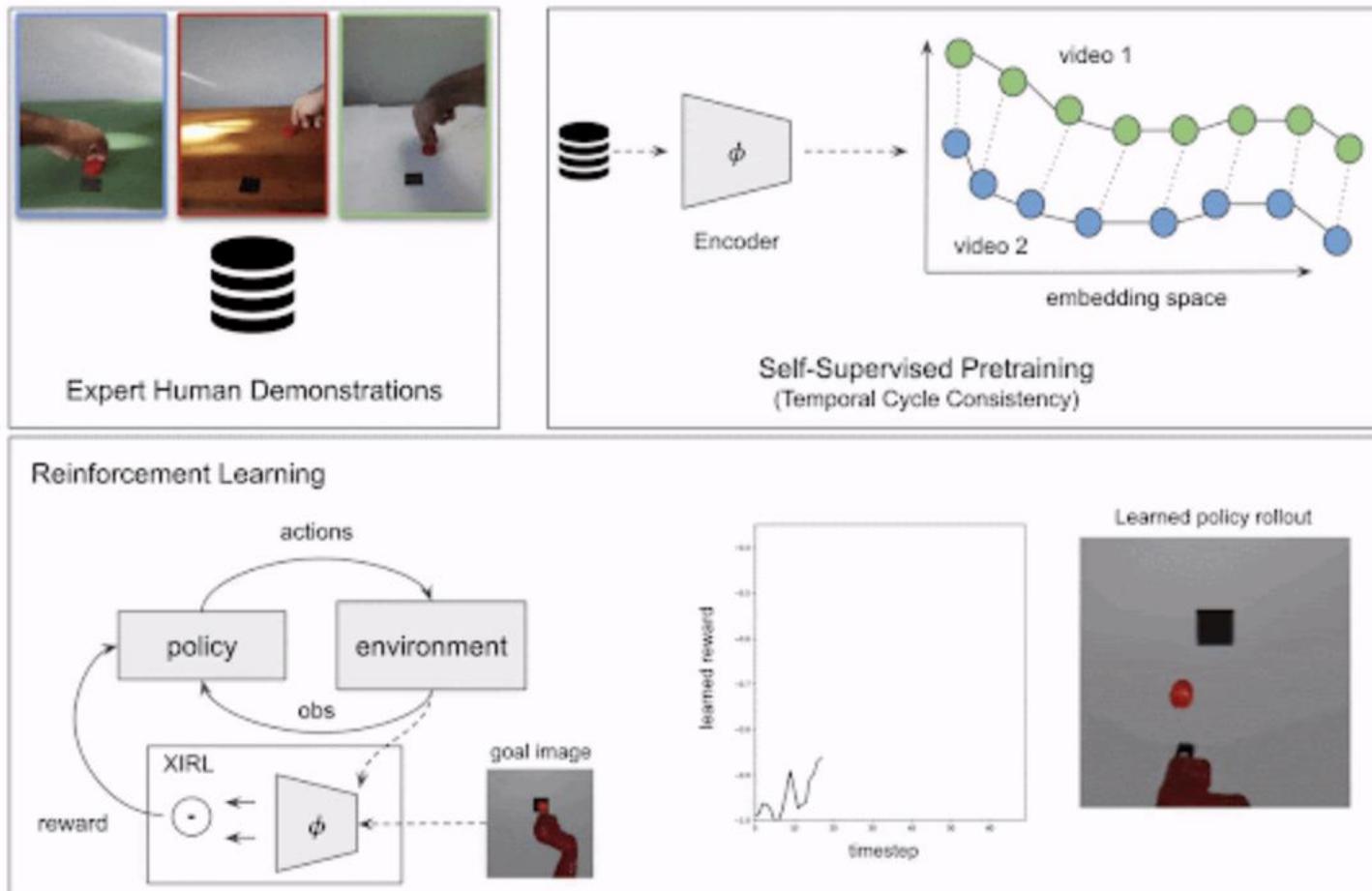
```
1: // Training DVD
2: Require:  $\mathcal{D}^h$  human demonstration data for  $N$  tasks  $\{\mathcal{T}_n\}$ 
3: Require:  $\mathcal{D}^r$  robot demonstration data for  $M$  tasks  $\{\mathcal{T}_m\} \subseteq \{\mathcal{T}_n\}$ 
4: Require: Pre-trained video encoder  $f_{enc}$ 
5: Randomly initialize  $\theta$ 
6: while training do
7:   Sample anchor video  $d_i \in \mathcal{D}^h \cup \mathcal{D}^r$ 
8:   Sample positive video  $d'_i \in \{\mathcal{D}_{\mathcal{T}_i}^h\} \cup \{\mathcal{D}_{\mathcal{T}_i}^r\} \setminus d_i$ 
9:   Sample negative video  $d_j \in \{\mathcal{D}_{\mathcal{T}_j}^h\} \cup \{\mathcal{D}_{\mathcal{T}_j}^r\} \forall j \neq i$ 
10:  Update  $\mathcal{R}_\theta$  with  $d_i, d'_i, d_j$  according to Eq. 1
11: // Planning Conditioned on Video Demo
12: Require: Trained reward function  $\mathcal{R}_\theta$  & video prediction model  $p_\phi$ 
13: Require: Human video demo  $d_i$  for task  $\mathcal{T}_i$ 
14: for trials  $1, \dots, n$  do
15:   Sample  $\{a_{1:H}^{1:G}\}$  & get predictions  $\{\tilde{s}_{1:H}^g\} \sim \{p_\phi(s_0, a_{1:H}^g)\}$ 
16:   Step  $a_{1:H}^*$  which maximizes  $\mathcal{R}_\theta(\tilde{s}_{1:H}^g, d_i)$ 
```

---

# XIRL: Cross-embodiment Inverse Reinforcement Learning

Kevin Zakka<sup>1,3\*</sup>, Andy Zeng<sup>2</sup>, Pete Florence<sup>2</sup>, Jonathan Tompson<sup>2</sup>,  
Jeannette Bohg<sup>1</sup>, and Debidatta Dwibedi<sup>2</sup>

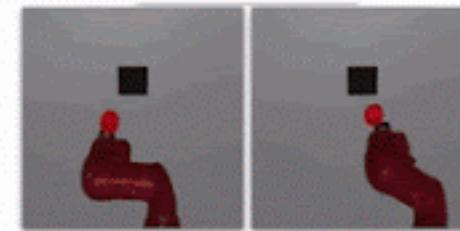
<sup>1</sup>Stanford University, <sup>2</sup>Robotics at Google, <sup>3</sup>UC Berkeley



# From IRL to RL

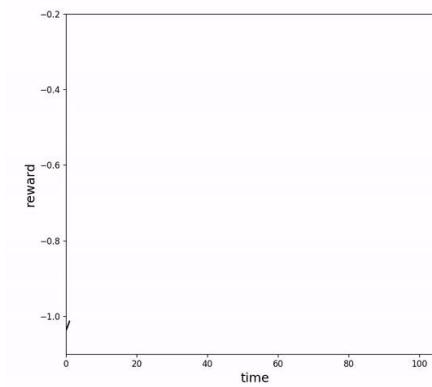
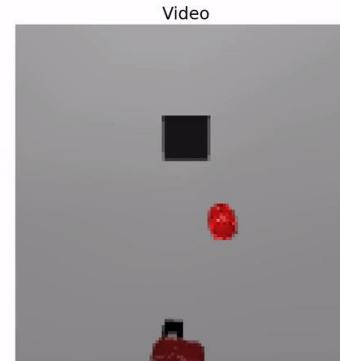
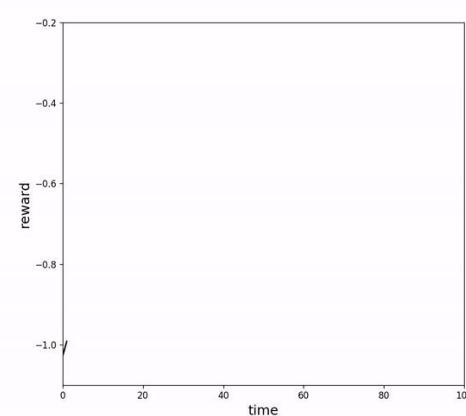
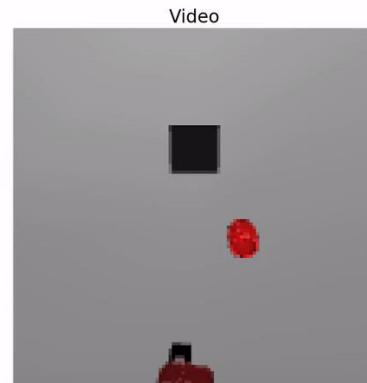
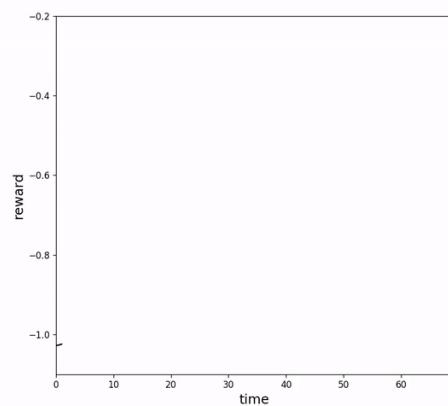
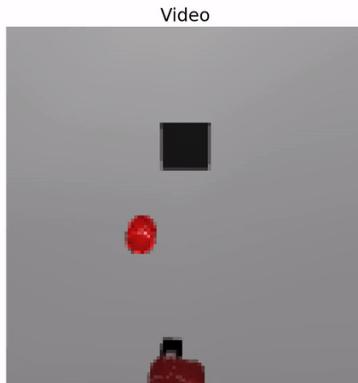


**Reward Learning**



**Reinforcement Learning**

# The XIRL-learned reward



---

# MINE DOJO: Building Open-Ended Embodied Agents with Internet-Scale Knowledge

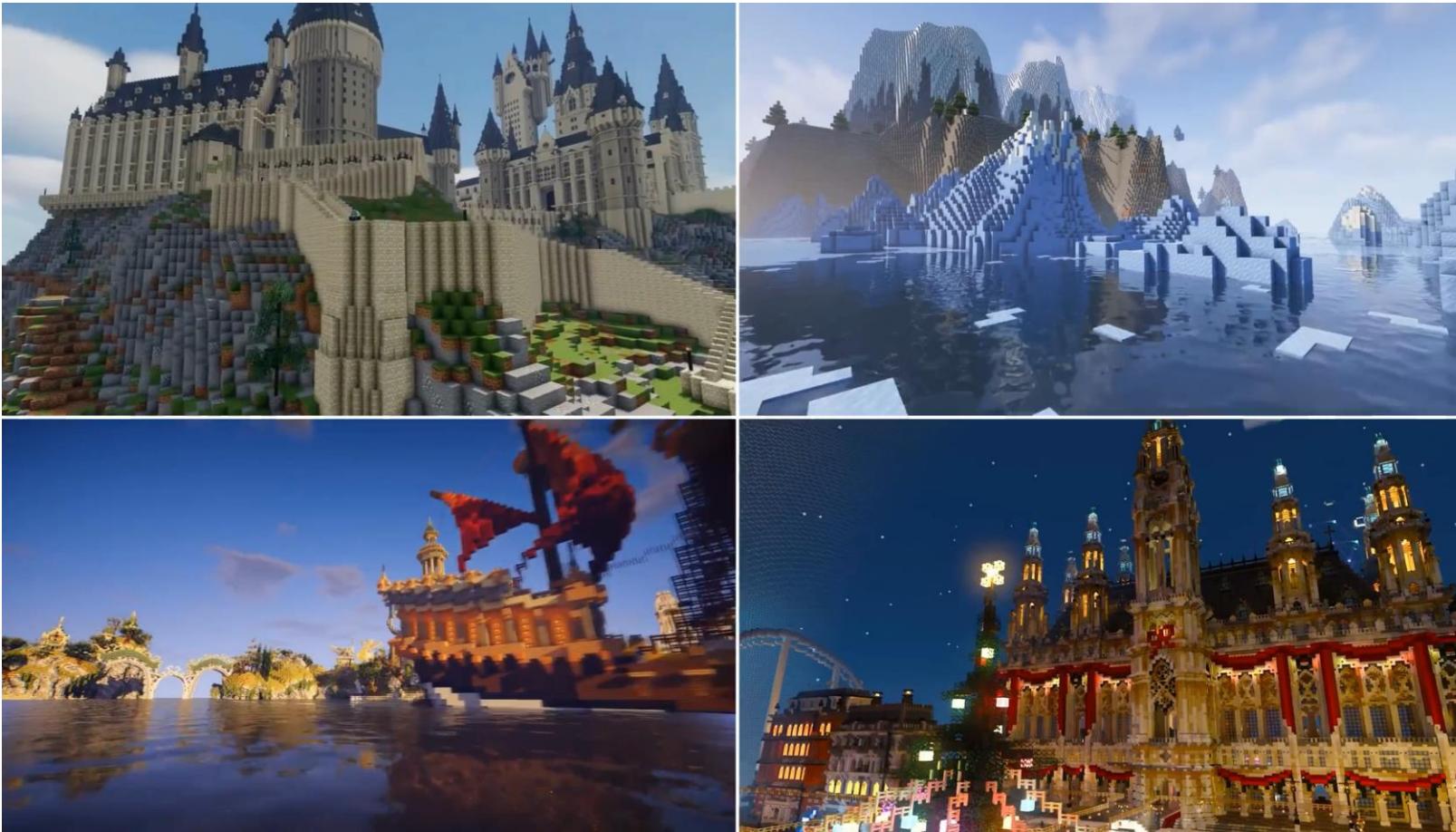
---

Linxi Fan<sup>1</sup>, Guanzhi Wang<sup>2\*</sup>, Yunfan Jiang<sup>3\*</sup>, Ajay Mandlekar<sup>1</sup>, Yuncong Yang<sup>4</sup>,  
Haoyi Zhu<sup>5</sup>, Andrew Tang<sup>4</sup>, De-An Huang<sup>1</sup>, Yuke Zhu<sup>1,6†</sup>, Anima Anandkumar<sup>1,2†</sup>

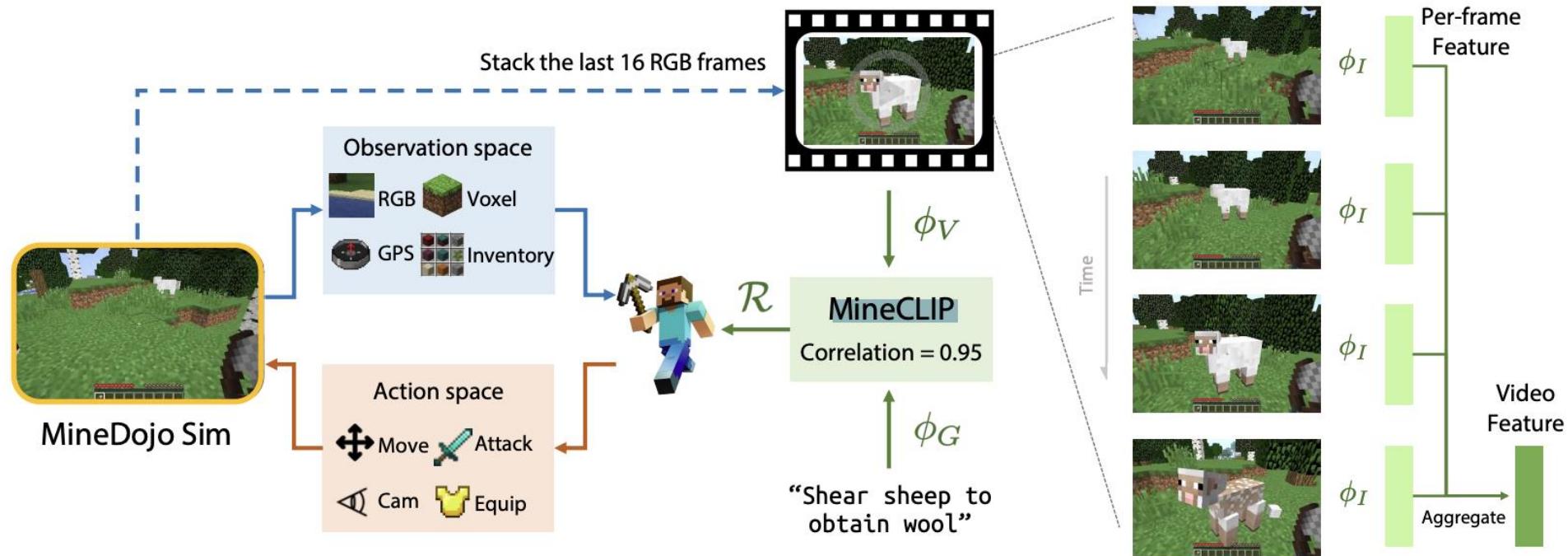
<sup>1</sup>NVIDIA, <sup>2</sup>Caltech, <sup>3</sup>Stanford, <sup>4</sup>Columbia, <sup>5</sup>SJTU, <sup>6</sup>UT Austin

\*Equal contribution †Equal advising

<https://minedojo.org>



# Learning Rewards from Language



Thank you!