



# Deep Reinforcement Learning

## Lecture 9: Diffusion Models for RL/IL

---

Huazhe Xu  
Tsinghua University

# AI This Week



# In Lec9

- 1 Diffusion Models
- 2 Diffusion in Policy Learning
- 3 Flow Matching





# In Lec9

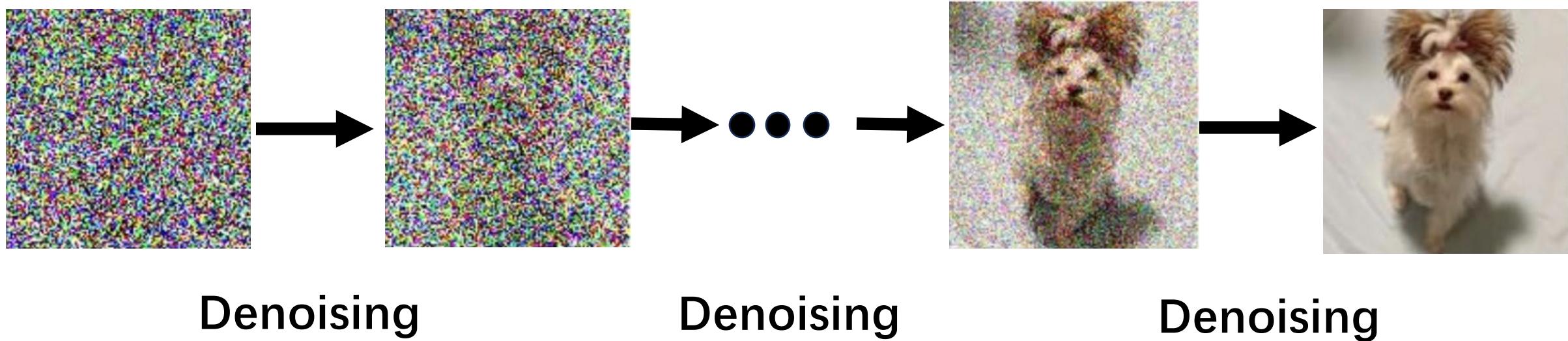
1 Diffusion Models

2 Diffusion Policy

3 Flow Matching

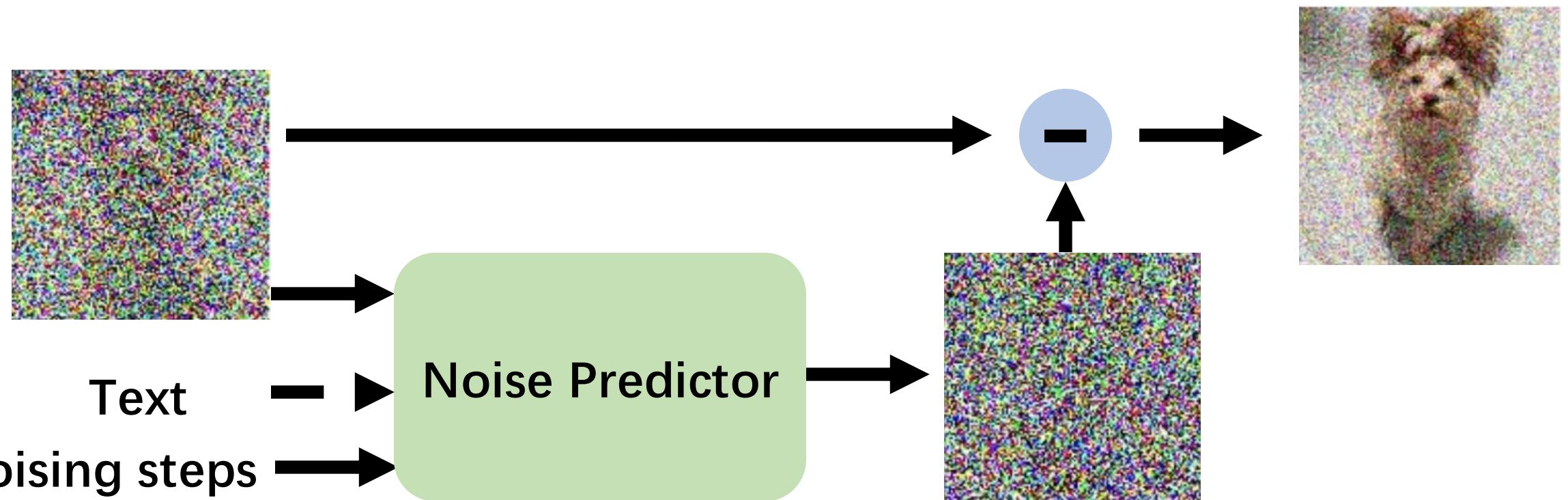


# What is diffusion?



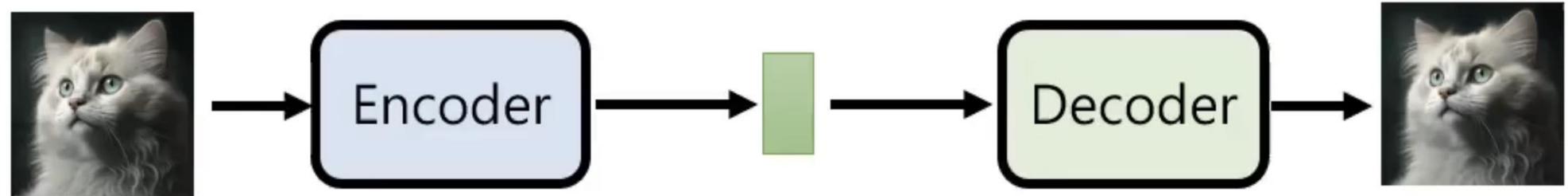
**Reverse process of adding noise to the original signals.**

# Denoising Module

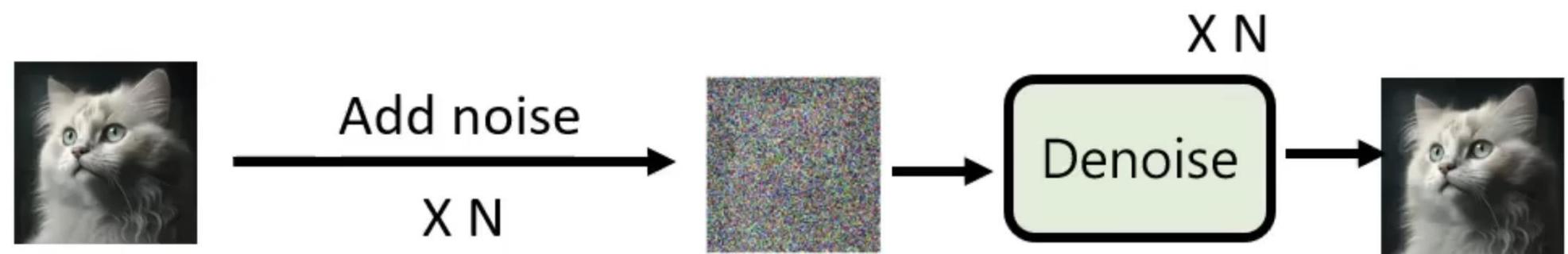


Predict the noise you added!

## VAE



## Diffusion



---

## Algorithm 1 Training

---

```
1: repeat  
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$            Sample a clean image  
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$   
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$        Sample a noise  
5:   Take gradient descent step on  
      $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$   
6: until converged
```

---

Combine image with noise

$\bar{\alpha}_1, \bar{\alpha}_2, \dots, \bar{\alpha}_T$   
→ smaller

---

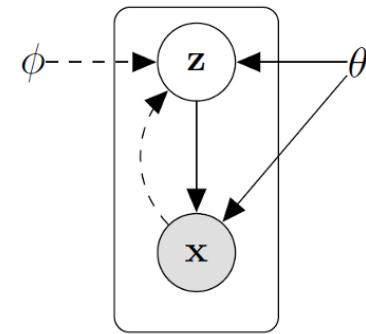
## Algorithm 2 Sampling

---

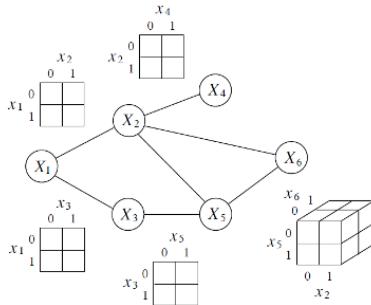
```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$            Sample a noise  
2: for  $t = T, \dots, 1$  do  
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$        Sample another  
     noise  
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$   
5: end for  
6: return  $\mathbf{x}_0$ 
```

---

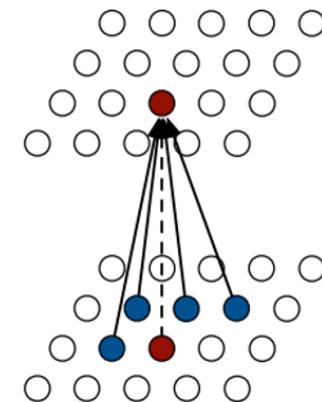
# Generative Modeling by Estimating Gradients of the Data Distribution



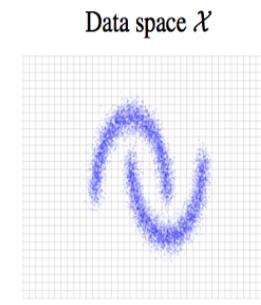
Bayesian networks  
(e.g., VAEs)



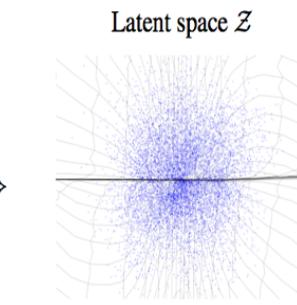
MRF



Autoregressive  
models



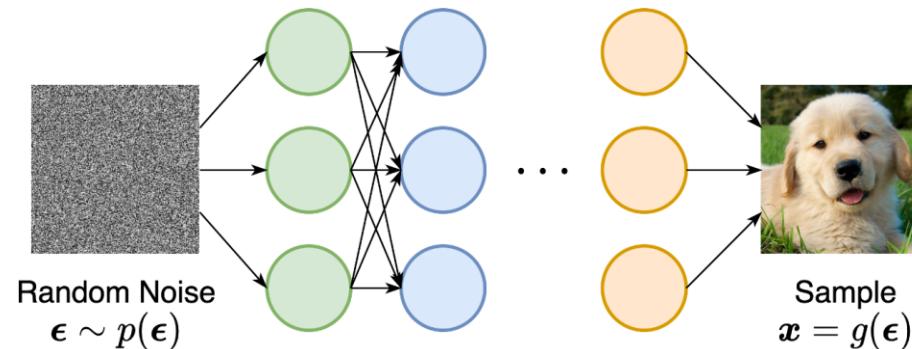
Data space  $\mathcal{X}$



Latent space  $\mathcal{Z}$

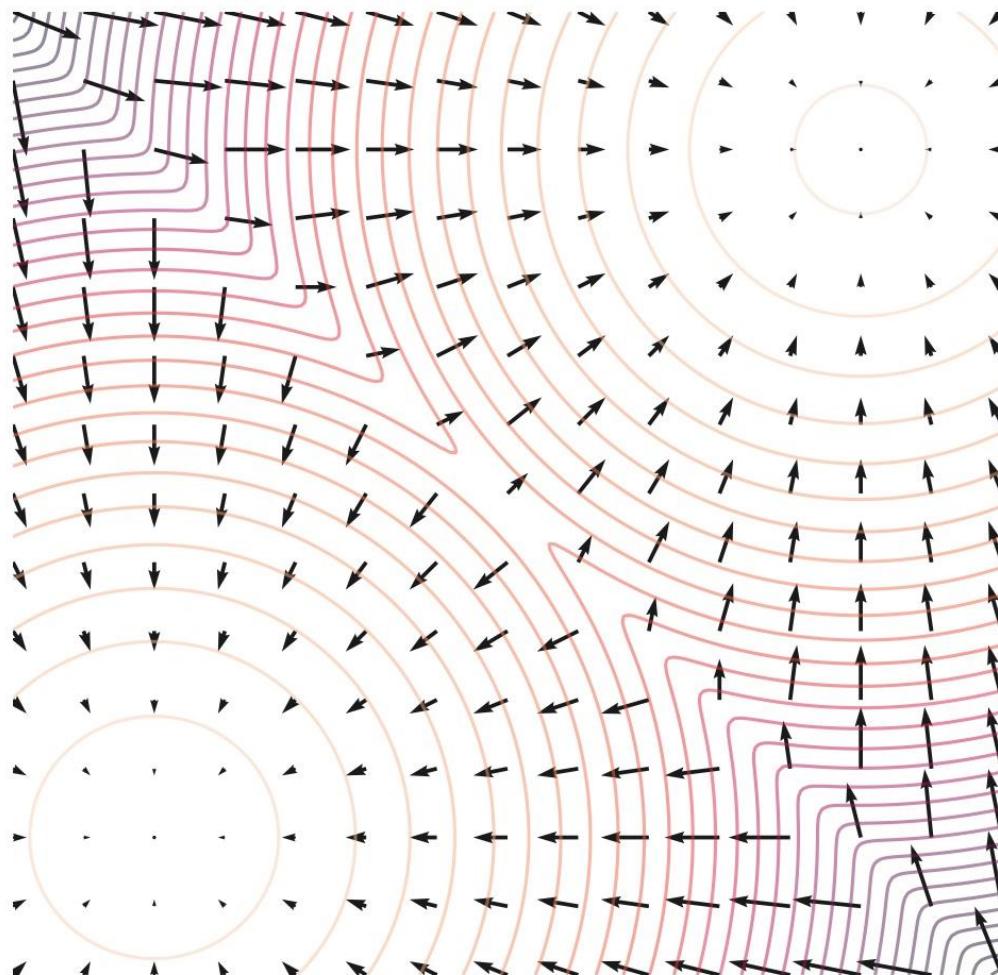
Likelihood-based  
Models

Special constraints



Implicit Generative  
Models  
Unstable

# Score-based Models



**Contours: Mixture of Gaussian**

**Vector Field: Score function**

Scores are not probabilities. No restrictions needed.

- Give dataset  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \sim p(x)$ , generate a new sample
- Max Likelihood approach may have p.d.f  $p_\theta(\mathbf{x}) = \frac{e^{-f_\theta(\mathbf{x})}}{Z_\theta}$ 
  - Normalizing, Usually intractable
- We can train  $\max_\theta \sum_{i=1}^N \log p_\theta(\mathbf{x}_i)$
- The score function  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$
- We get rid of the Z by using

$$\mathbf{s}_\theta(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}) = -\nabla_{\mathbf{x}} f_\theta(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log Z_\theta}_{=0} = -\nabla_{\mathbf{x}} f_\theta(\mathbf{x}).$$

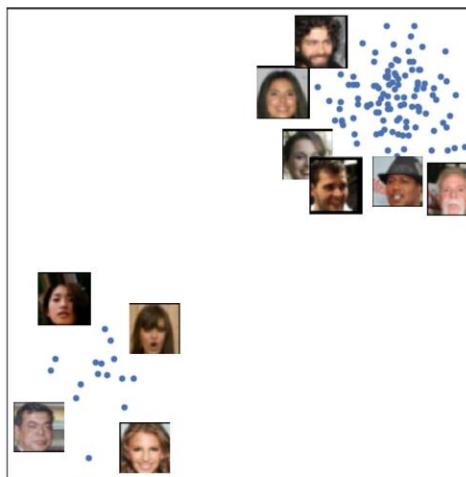
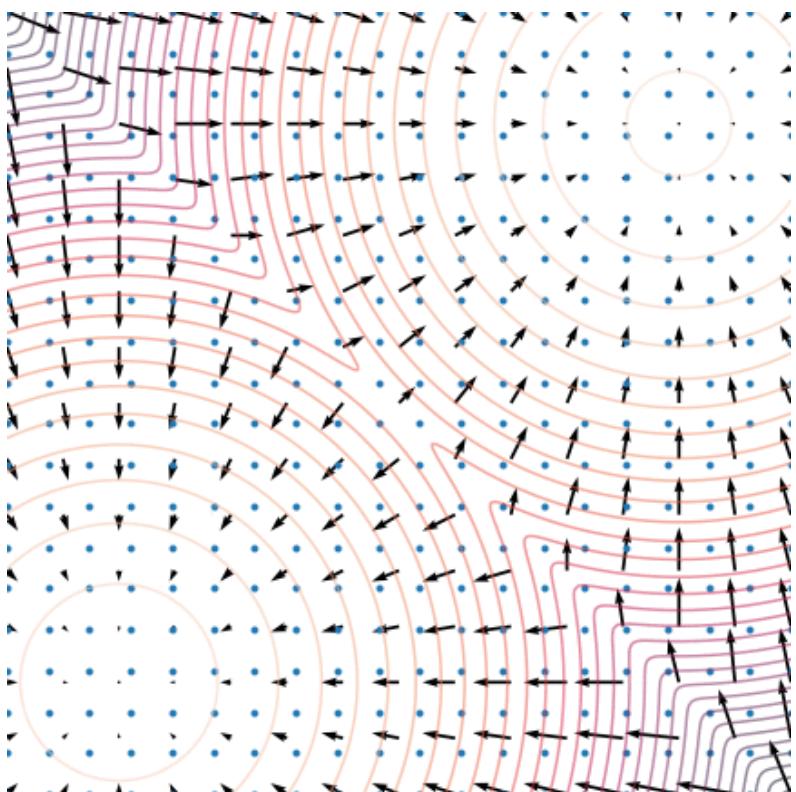
- Loss

$$\mathbb{E}_{p(\mathbf{x})} \left[ \|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x})\|_2^2 \right]$$

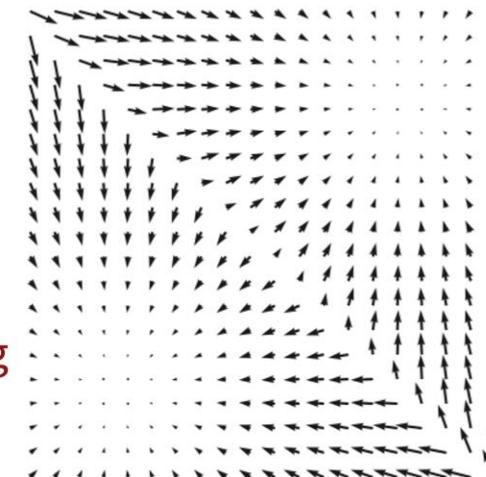
$$J(\boldsymbol{\theta}) \triangleq \mathbb{E}_{p_d} \left[ \text{tr}(\nabla_{\mathbf{x}} \mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})) + \frac{1}{2} \|\mathbf{s}_m(\mathbf{x}; \boldsymbol{\theta})\|_2^2 \right]$$

- Sampling by Langevin Dynamics

$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \epsilon \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \sqrt{2\epsilon} \mathbf{z}_i, \quad i = 0, 1, \dots, K,$$



score  
matching



Scores

$$\mathbf{s}_{\theta}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

Langevin  
dynamics



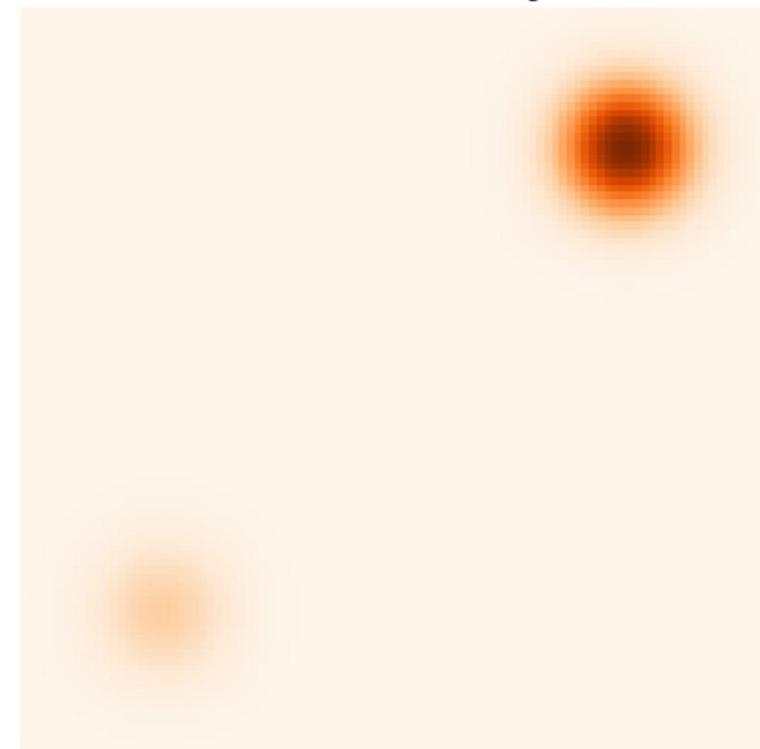
New samples

- Inaccurate Estimation

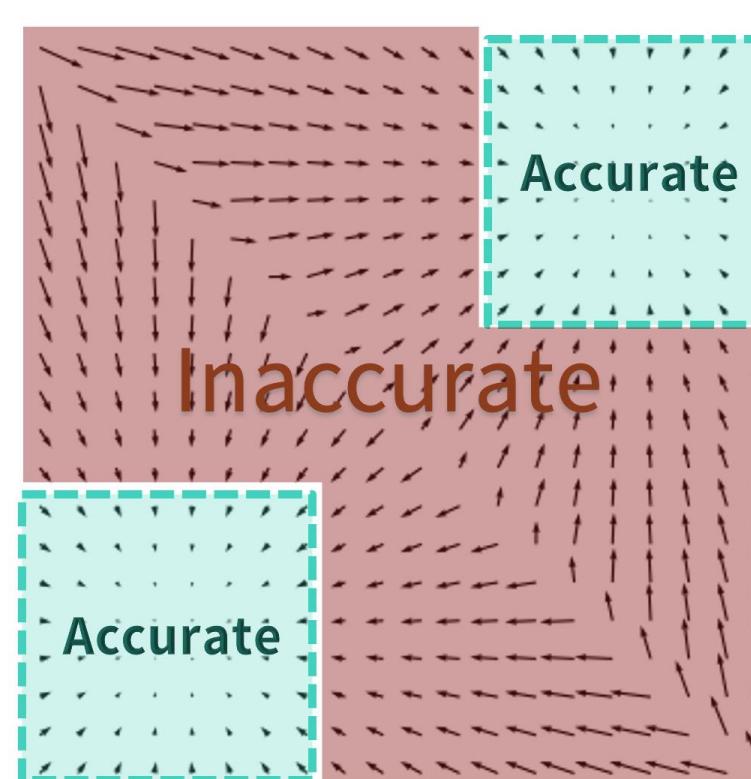
$$\mathbb{E}_{p(\mathbf{x})} \left[ \|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x})\|_2^2 \right] = \int p(\mathbf{x}) \|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x})\|_2^2 \, d\mathbf{x}$$

Weight

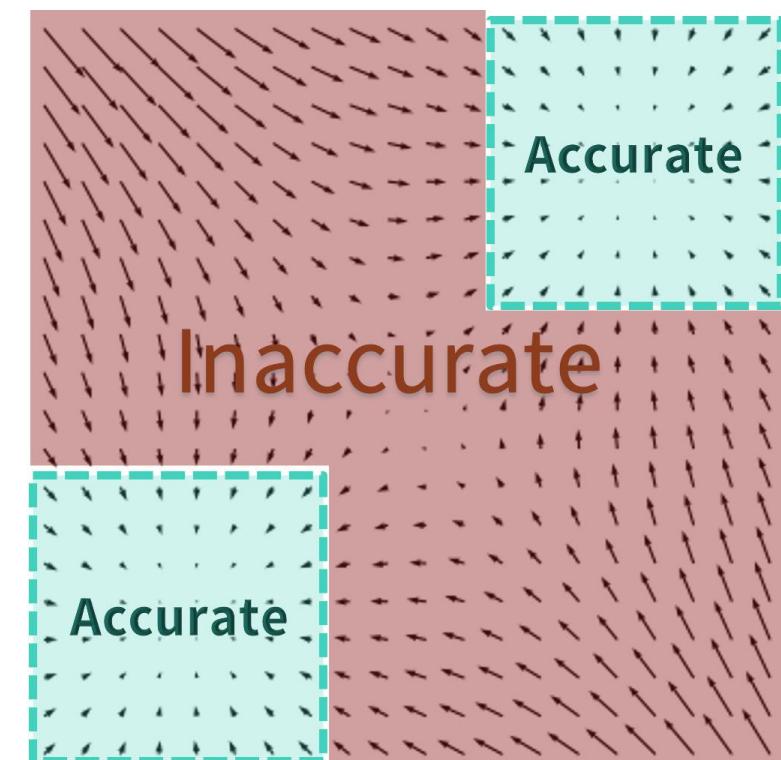
Data density



Data scores

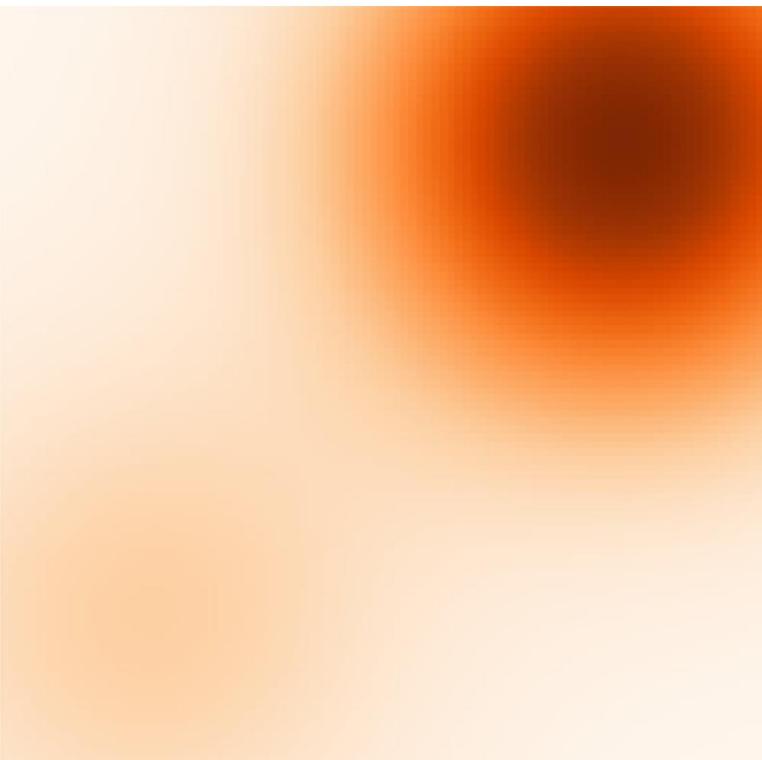


Estimated scores

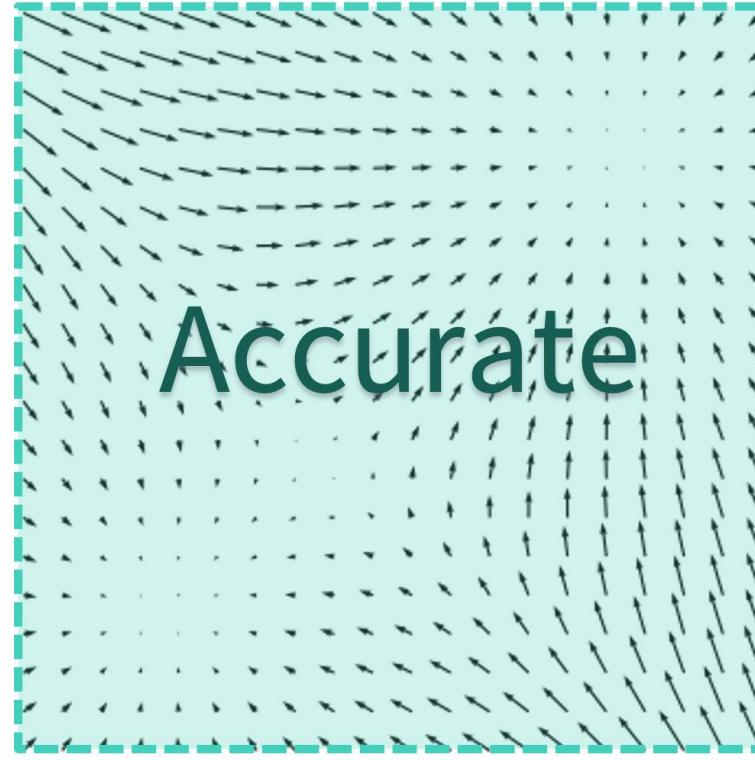


- Perturb the samples

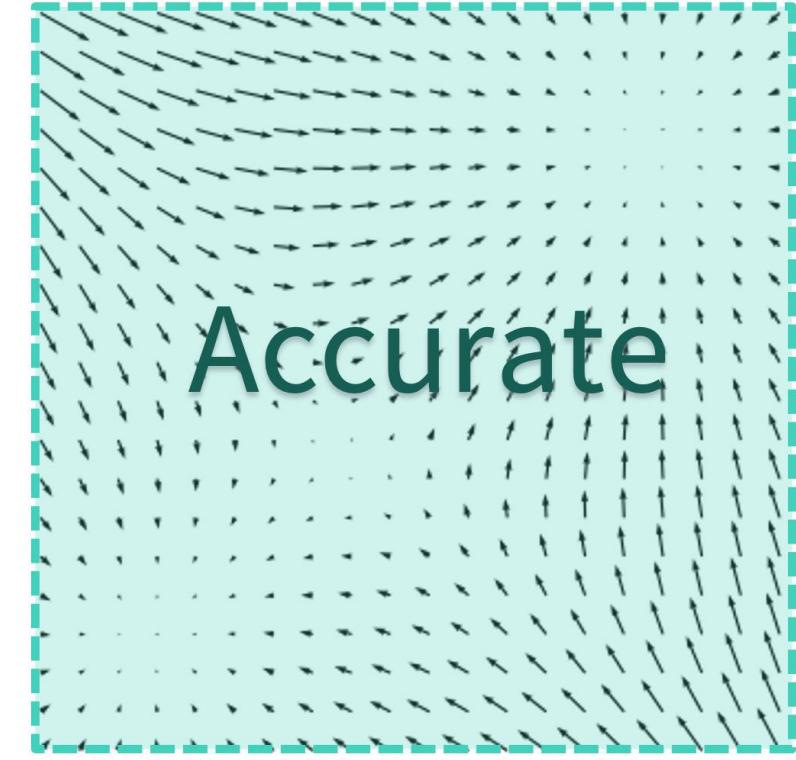
Perturbed density



Perturbed scores



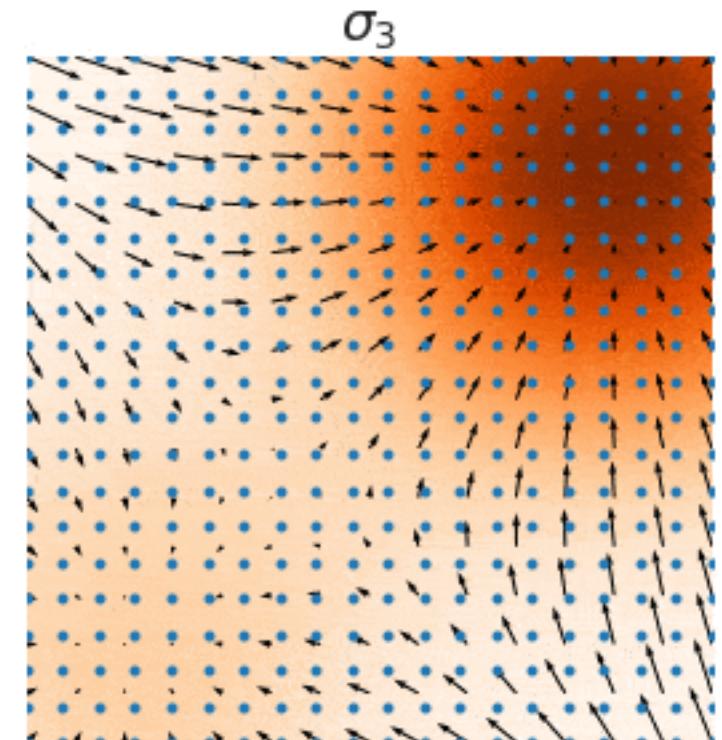
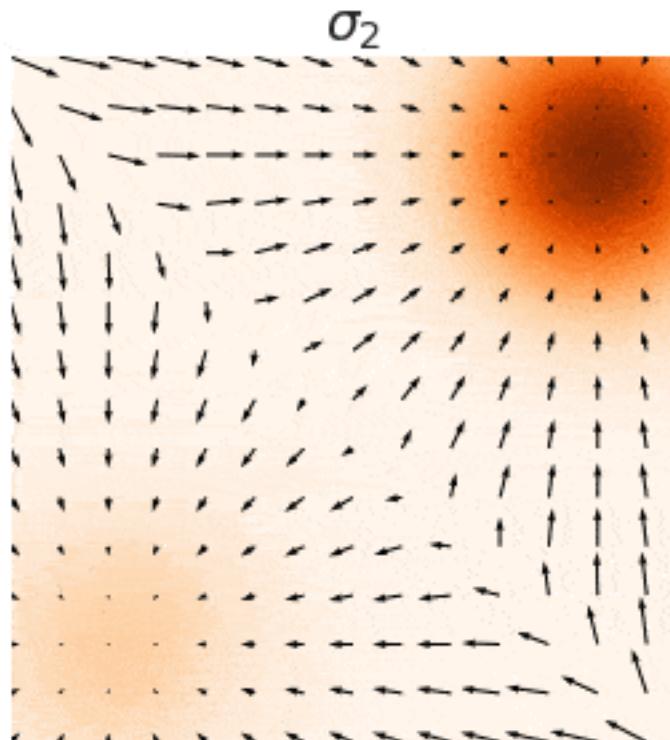
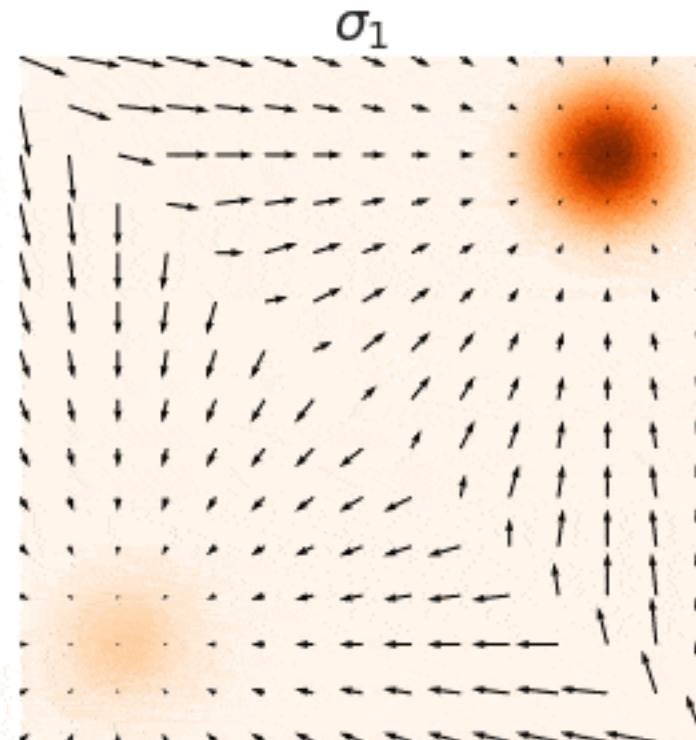
Estimated scores

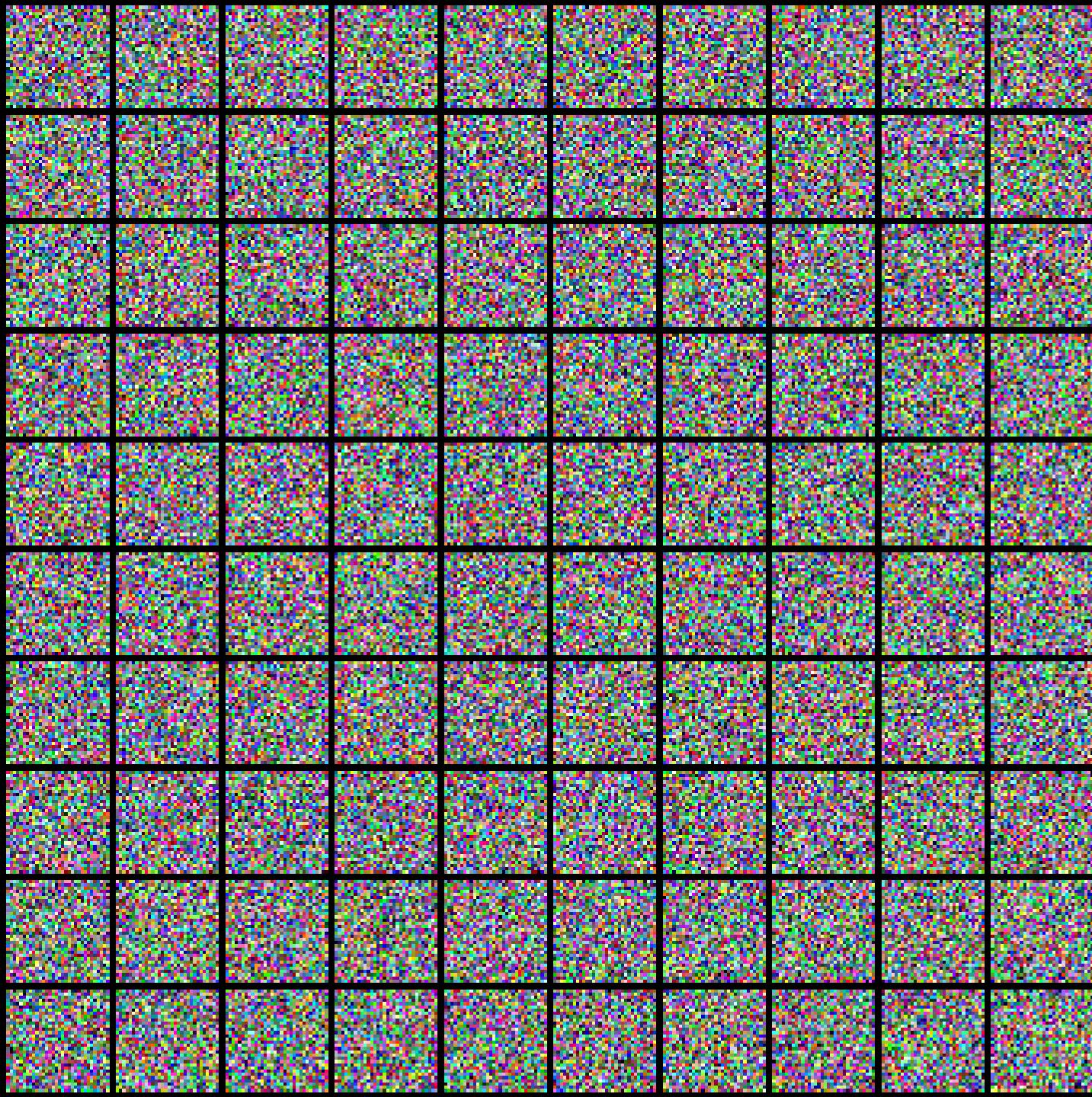


Bigger noise: better scores, worse data dist

$$\sigma_1 < \sigma_2 < \dots < \sigma_L$$

$$\sigma_1 < \sigma_2 < \cdots < \sigma_L$$







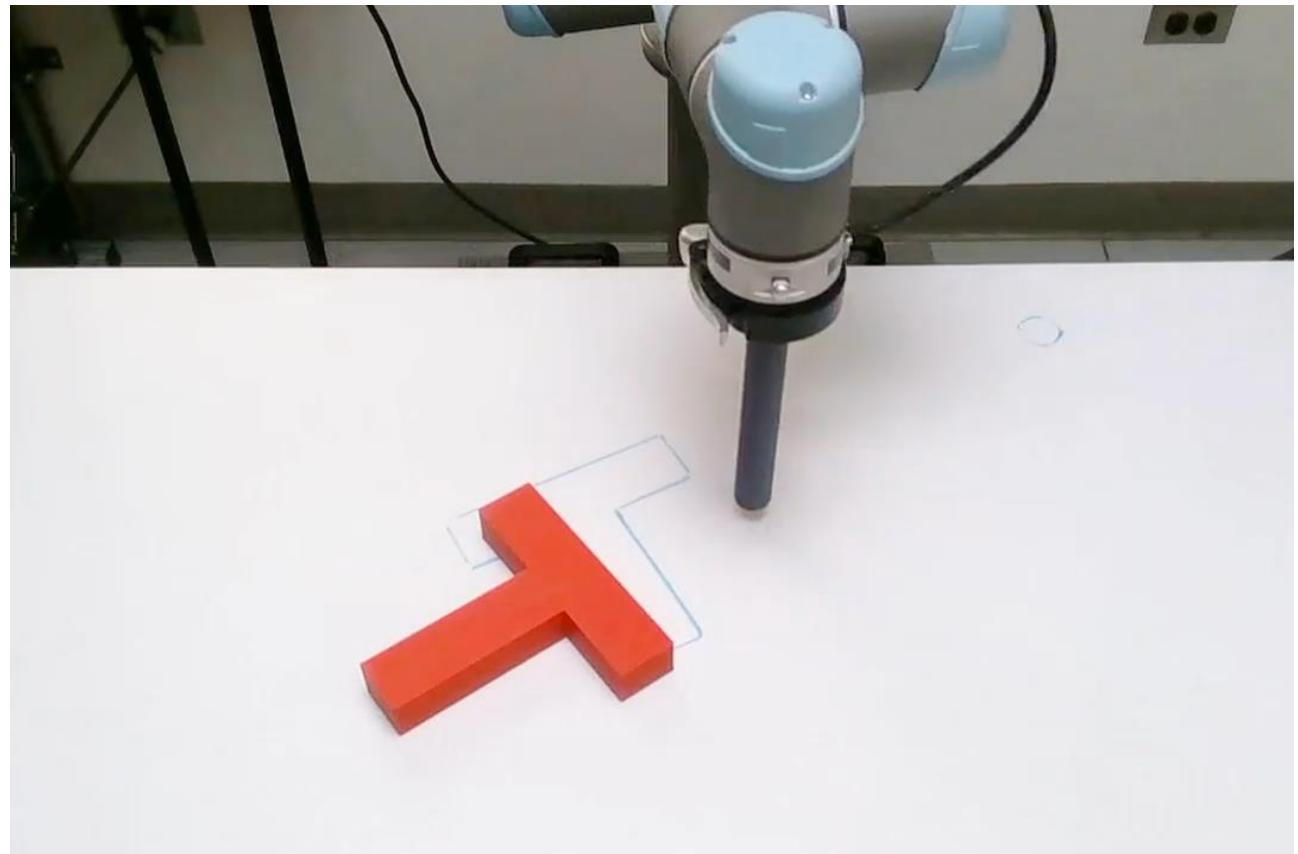
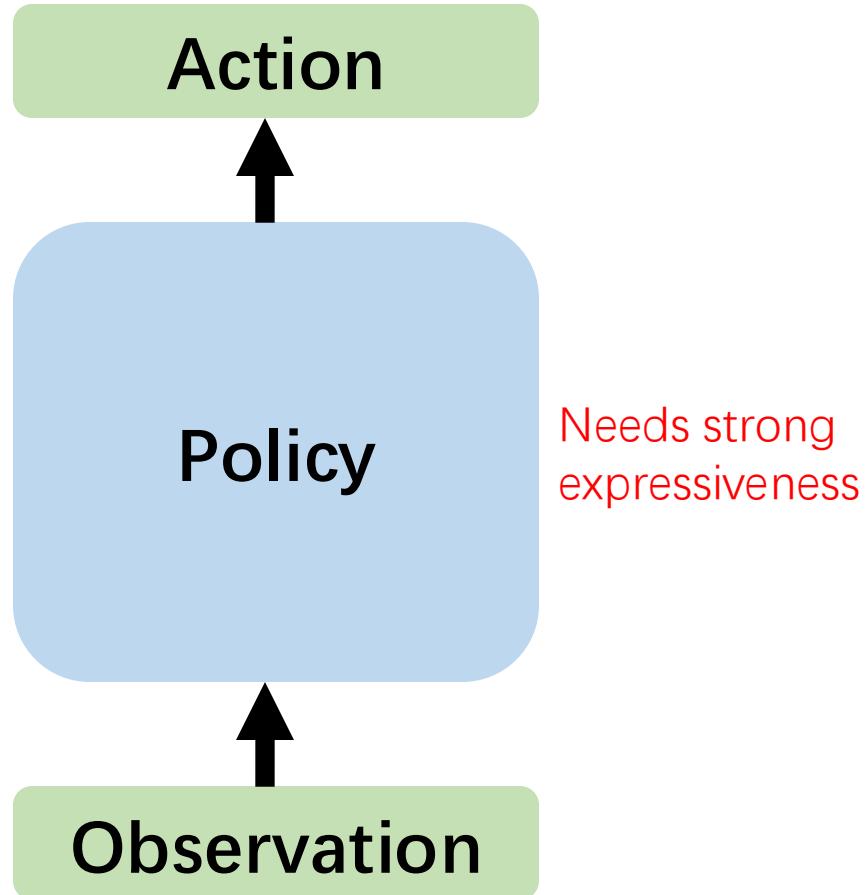
# In Lec9

1 Diffusion Models

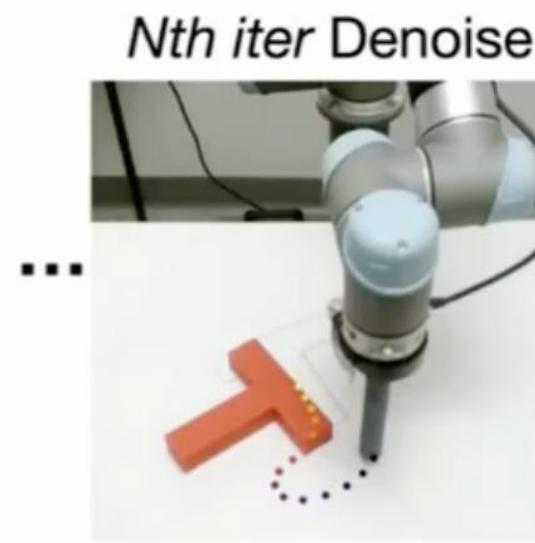
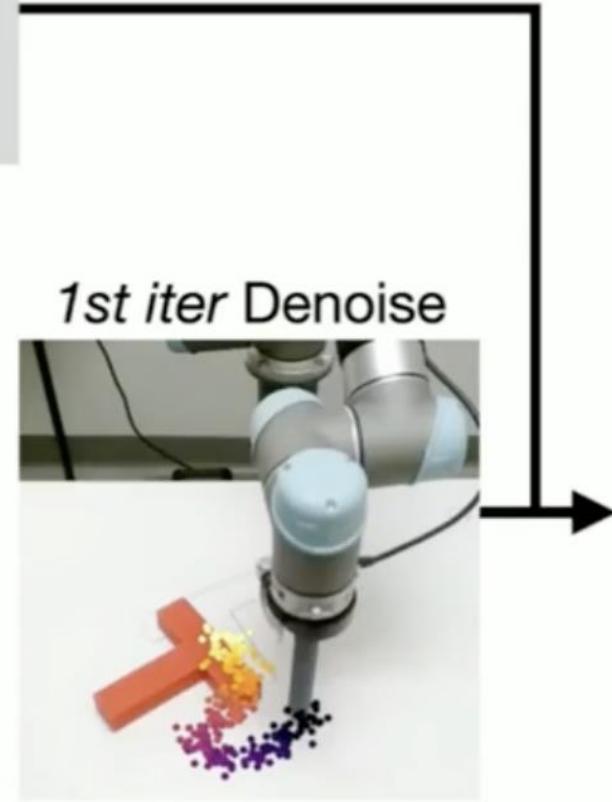
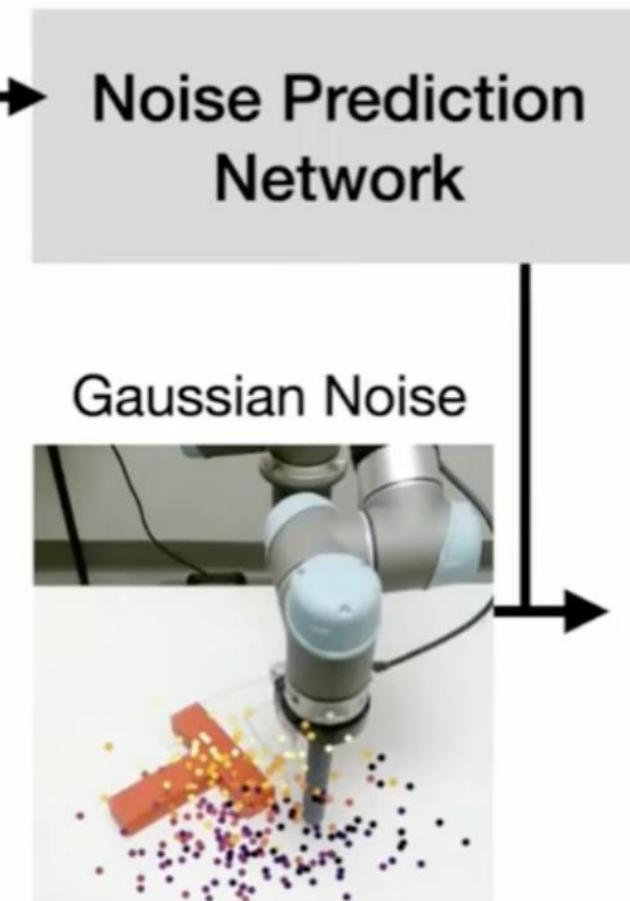
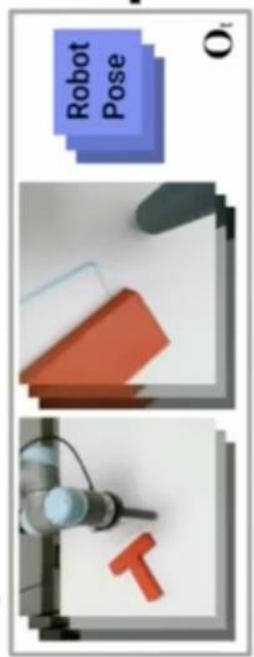
2 Diffusion in Policy Learning

3 Flow Matching

# Diffusion Policy



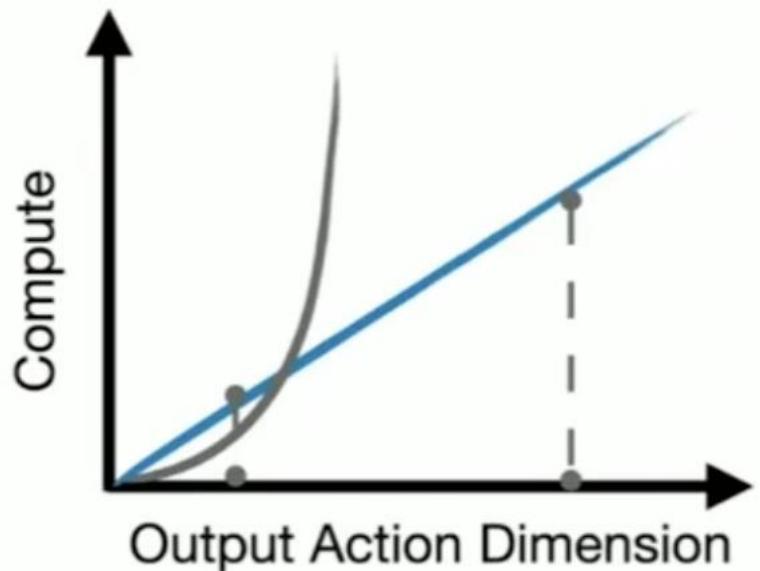
**Input: Observation**



**Output: Action**



**Action  
Multimodality**

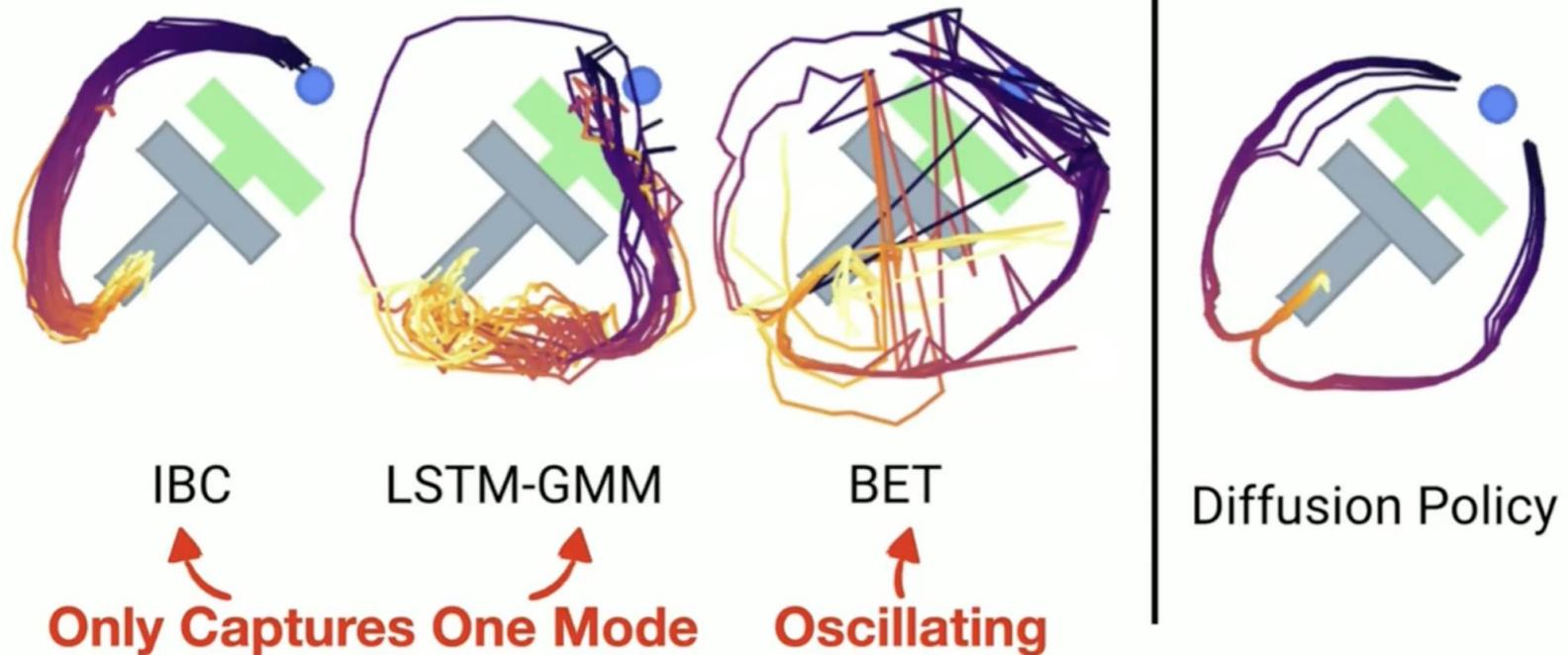
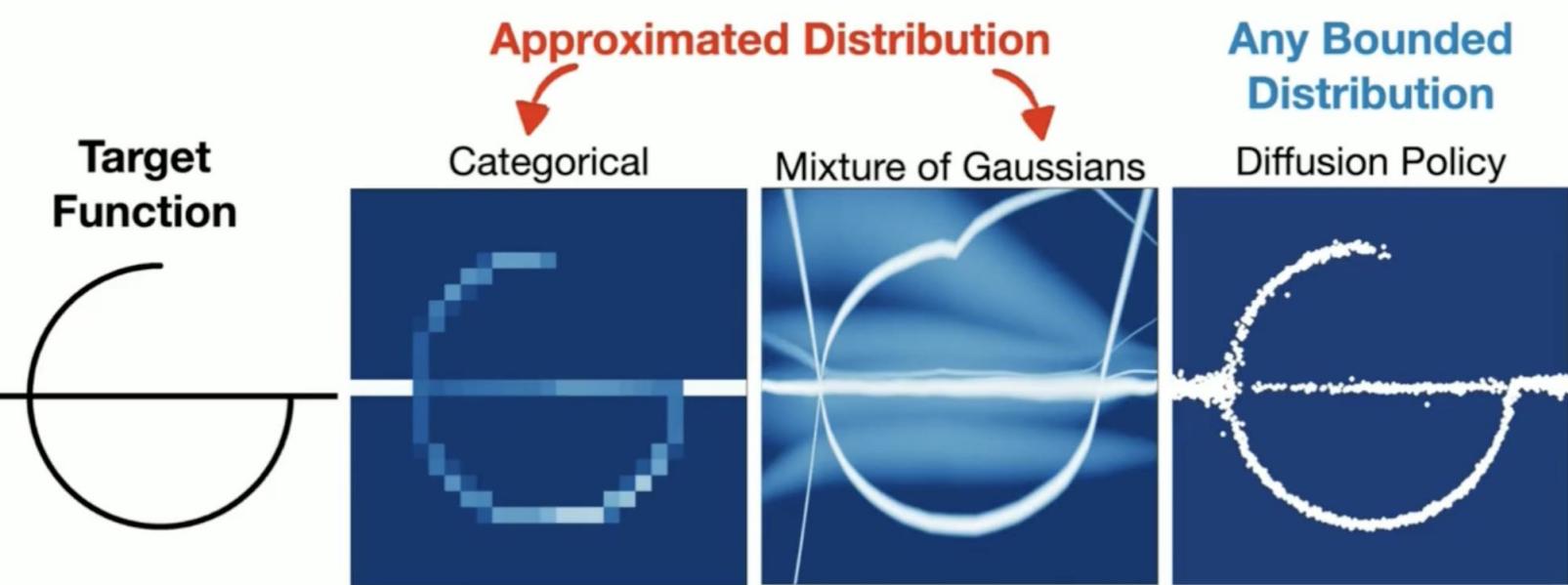


**Action Space  
Scalability**

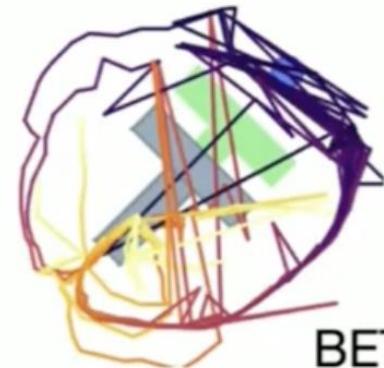
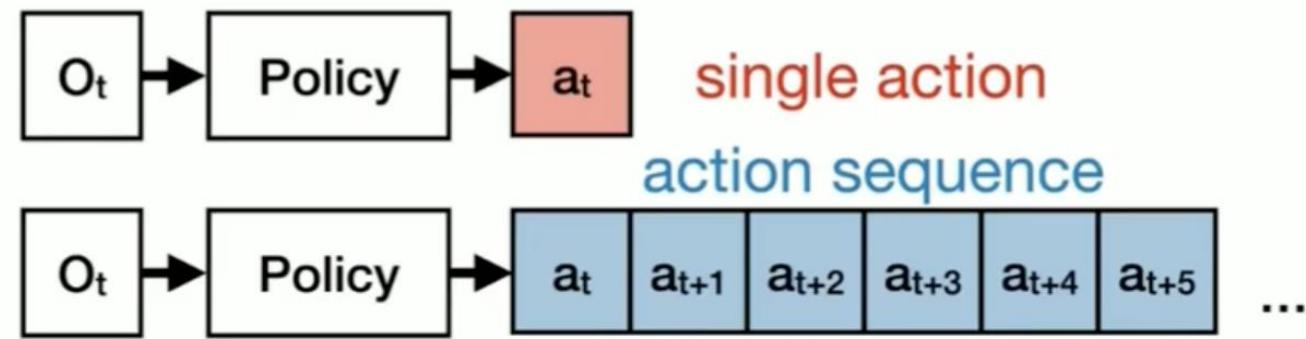
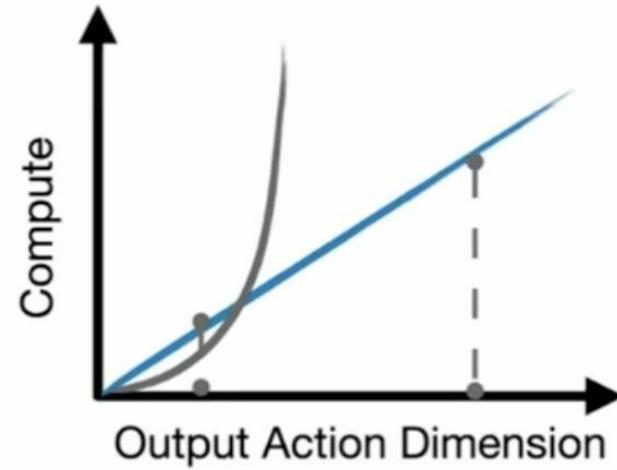


**Training  
Stability**

# Action Modality



# Action Space Scalability

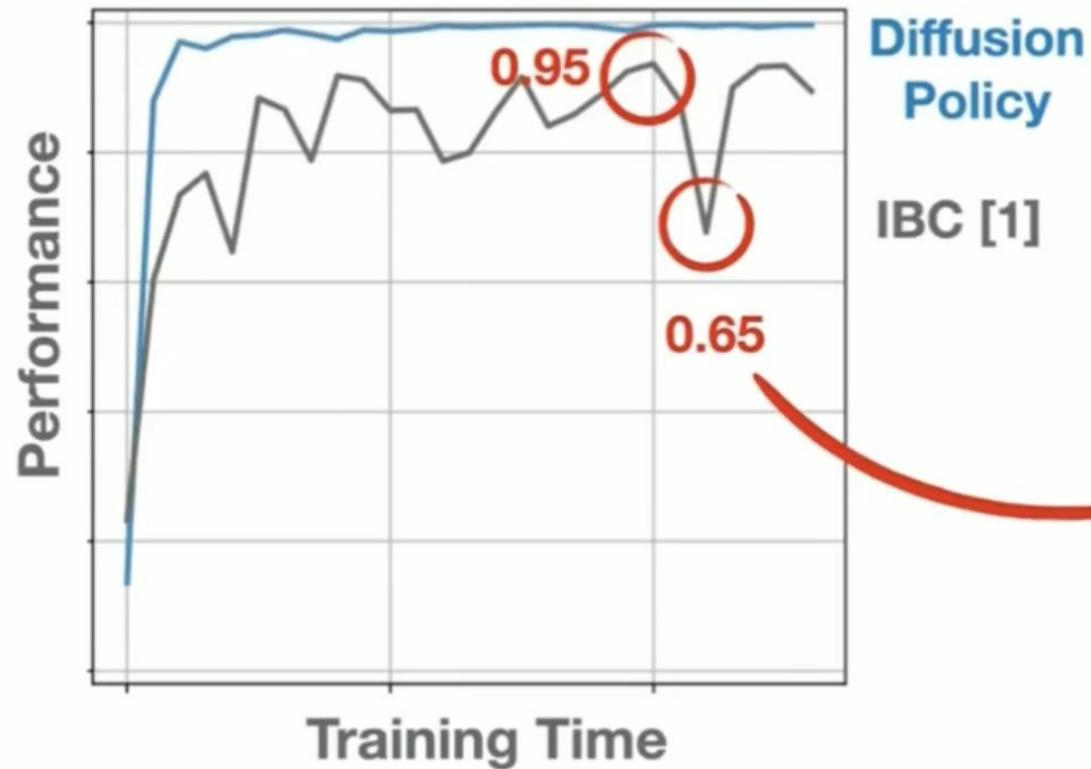


Oscillation between Modes

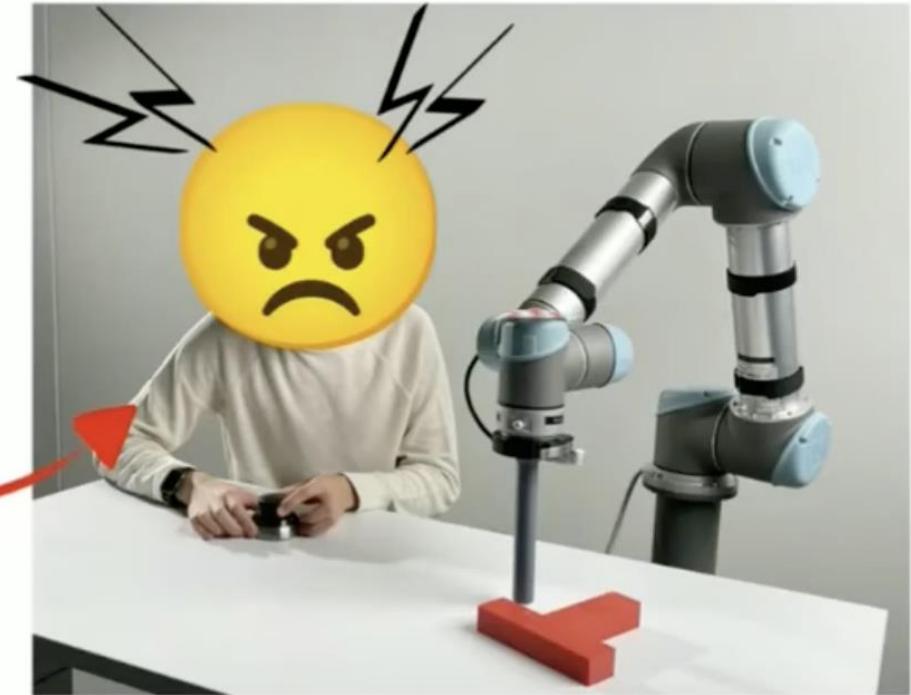


Temporal Mode Consistency

# Training Stability

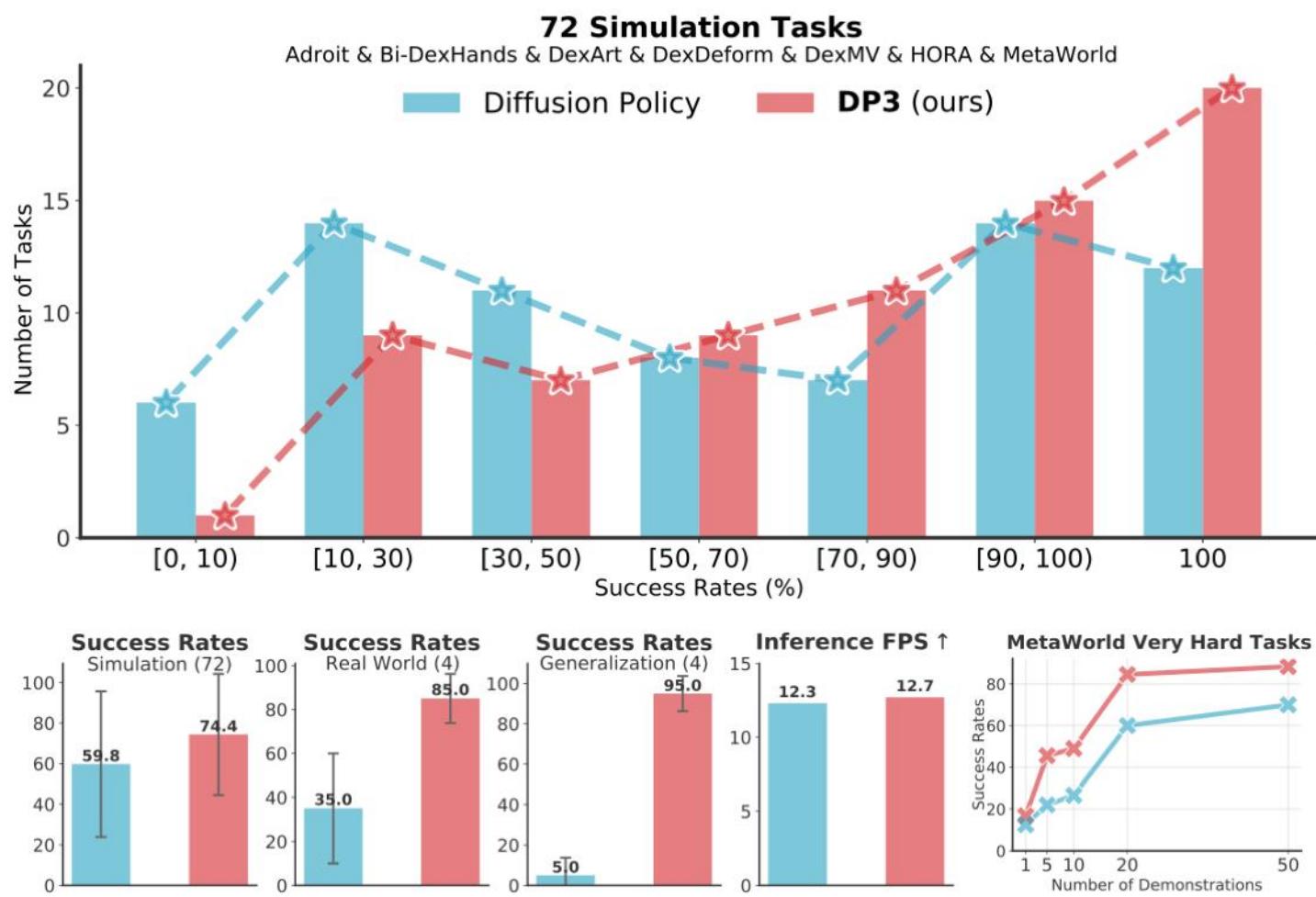


Diffusion Policy  
IBC [1]

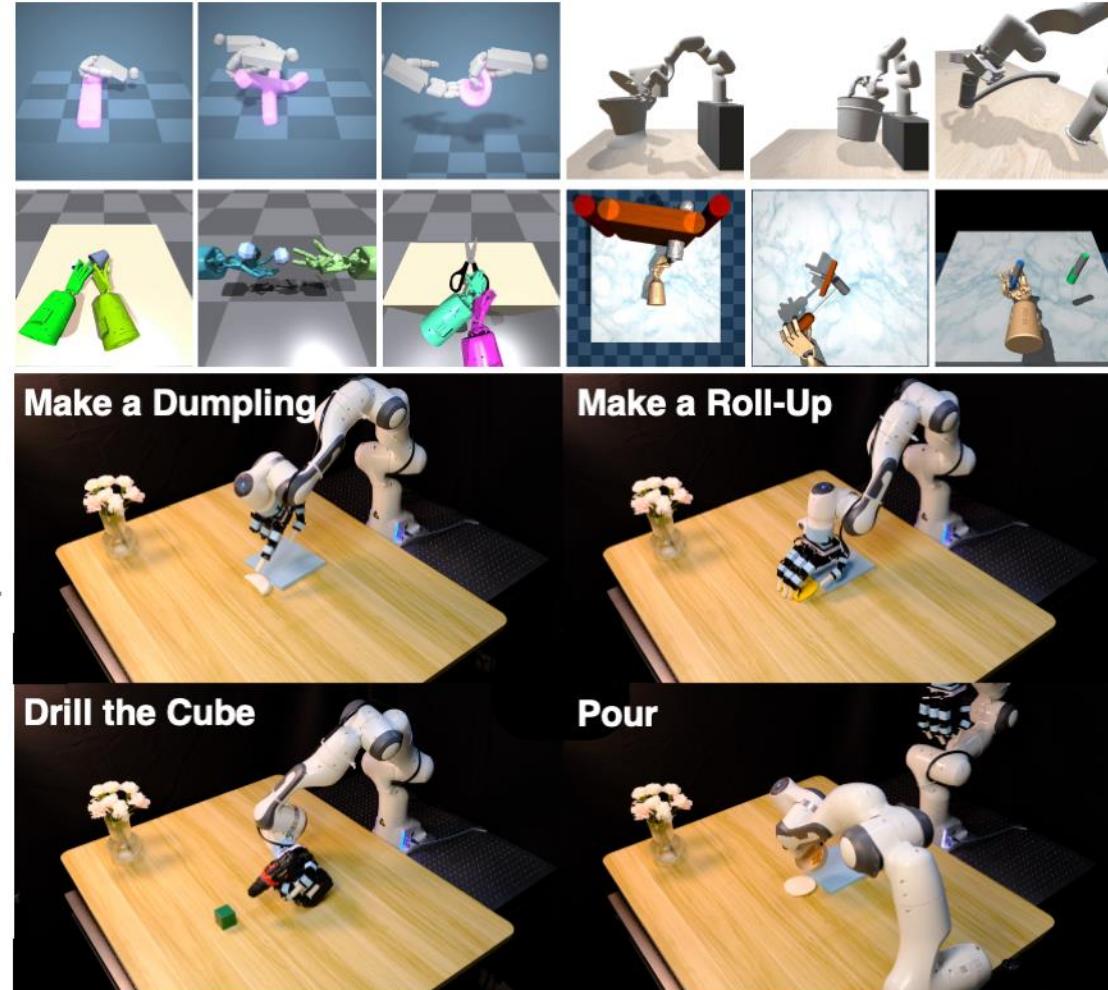


Checkpoint selection requires  
expensive realworld evaluation

# DP3

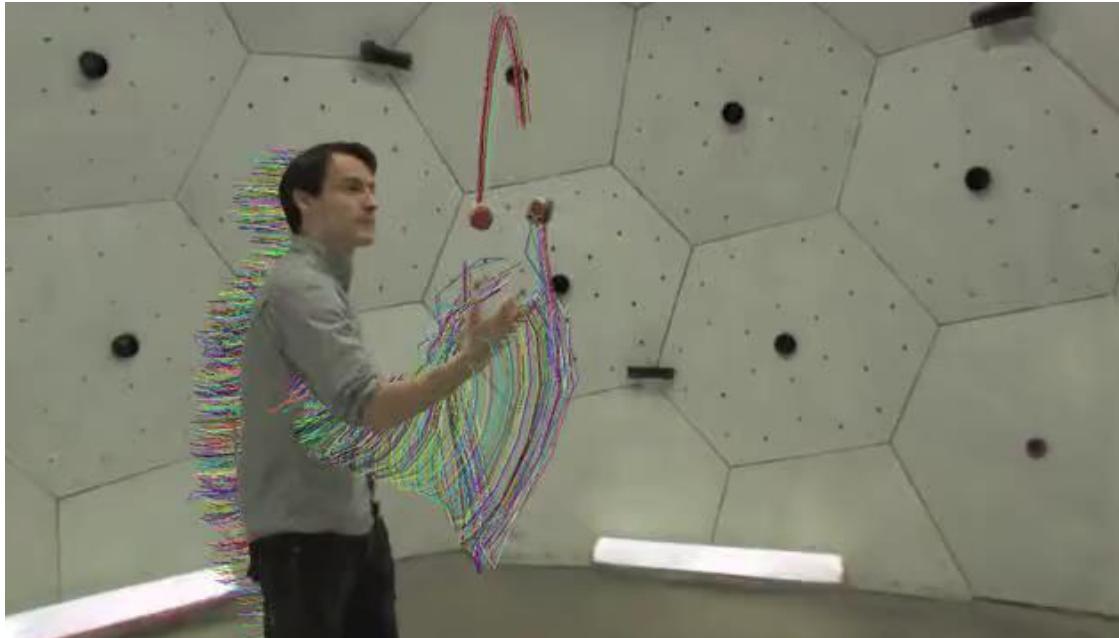


(a) 3D Diffusion Policy (DP3) vs. Diffusion Policy: **Better, Faster, Stronger.**



# What could be the Missing Point?

The world itself is 3D.



Luiten et al., 2023

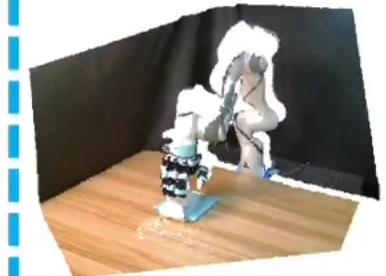
# 3D Diffusion Policy (DP3)

## Diffusion Policy + Very Simple 3D Representations

### Perception: Compact 3D Representations from Point Clouds

#### (a) Point Cloud Processing

Single-view Point Cloud



# Why Simple 3D Representations?

## EXP1: Point Cloud $\approx$ Oracle State > Voxel/Depth/Image

Repr.	H	D	P	A	DA	SP	Average
Oracle State	99 $\pm$ 2	61 $\pm$ 2	44 $\pm$ 3	94 $\pm$ 1	72 $\pm$ 7	91 $\pm$ 8	76.8
<b>Point cloud</b>	<b>100<math>\pm</math>0</b>	<b>62<math>\pm</math>4</b>	43 $\pm$ 6	<b>99<math>\pm</math>1</b>	69 $\pm$ 4	<b>97<math>\pm</math>4</b>	<b>78.3</b>
Image	48 $\pm$ 17	50 $\pm$ 5	25 $\pm$ 4	15 $\pm$ 1	43 $\pm$ 7	63 $\pm$ 3	40.7
Depth	39 $\pm$ 15	49 $\pm$ 1	12 $\pm$ 3	15 $\pm$ 4	15 $\pm$ 2	62 $\pm$ 3	32.0
RGB-D	57 $\pm$ 14	47 $\pm$ 5	14 $\pm$ 2	15 $\pm$ 3	14 $\pm$ 1	61 $\pm$ 3	34.7
Voxel	54 $\pm$ 5	33 $\pm$ 3	18 $\pm$ 2	10 $\pm$ 2	17 $\pm$ 1	62 $\pm$ 6	32.3

## EXP2: DP3 Encoder >> Other Point Encoders (Why?)

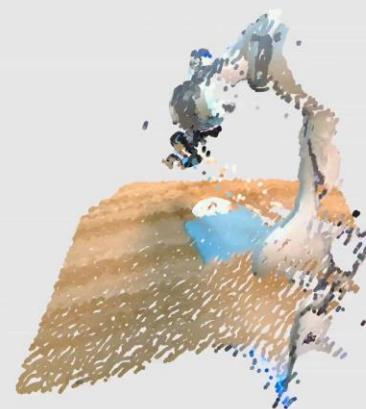
Encoders	H	D	P	A	DA	SP	Average
<b>DP3 Encoder</b>	<b>100<math>\pm</math>0</b>	<b>62<math>\pm</math>4</b>	<b>43<math>\pm</math>6</b>	<b>99<math>\pm</math>1</b>	<b>69<math>\pm</math>4</b>	<b>97<math>\pm</math>4</b>	<b>78.3</b>
PointNet	46 $\pm$ 8	34 $\pm$ 8	14 $\pm$ 4	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	15.7
PointNet++	0 $\pm$ 0	0 $\pm$ 0	13 $\pm$ 3	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	2.2
PointNeXt	0 $\pm$ 0	0 $\pm$ 0	14 $\pm$ 3	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	2.3
Point Transformer	0 $\pm$ 0	0 $\pm$ 0	6 $\pm$ 5	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	1.0
PointNet++ (pre-trained)	5 $\pm$ 9	19 $\pm$ 12	17 $\pm$ 6	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	6.8
PointNeXt (pre-trained)	0 $\pm$ 0	36 $\pm$ 13	17 $\pm$ 6	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	8.8

# Evaluate DP3 on Challenging Real-World Tasks

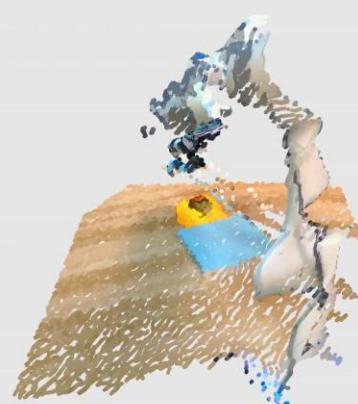
Drill



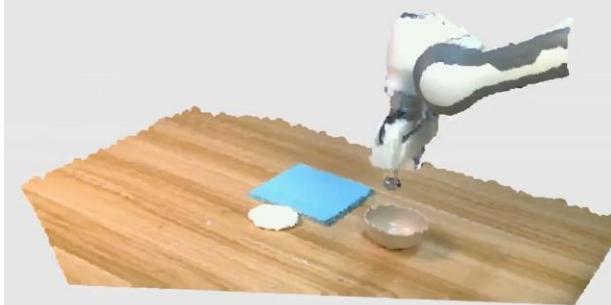
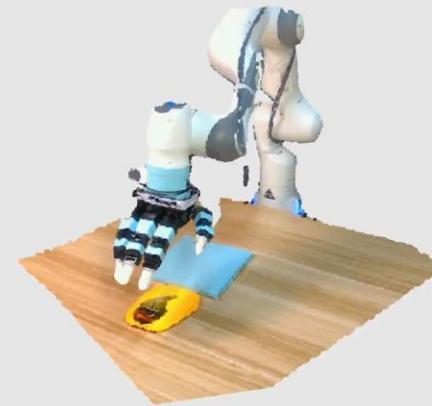
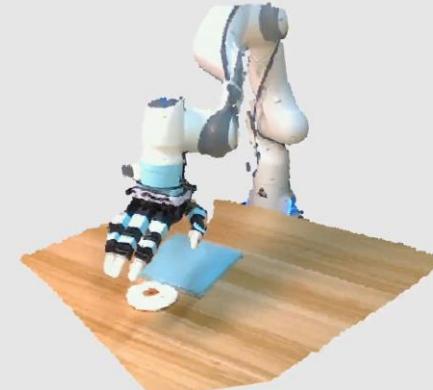
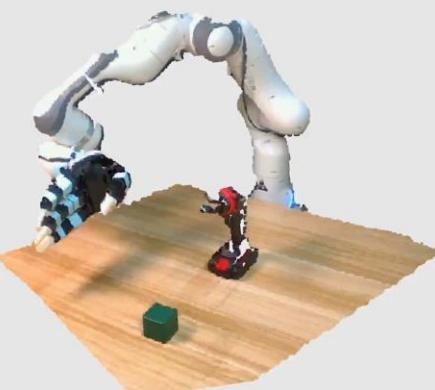
Dumpling



Roll



Pour

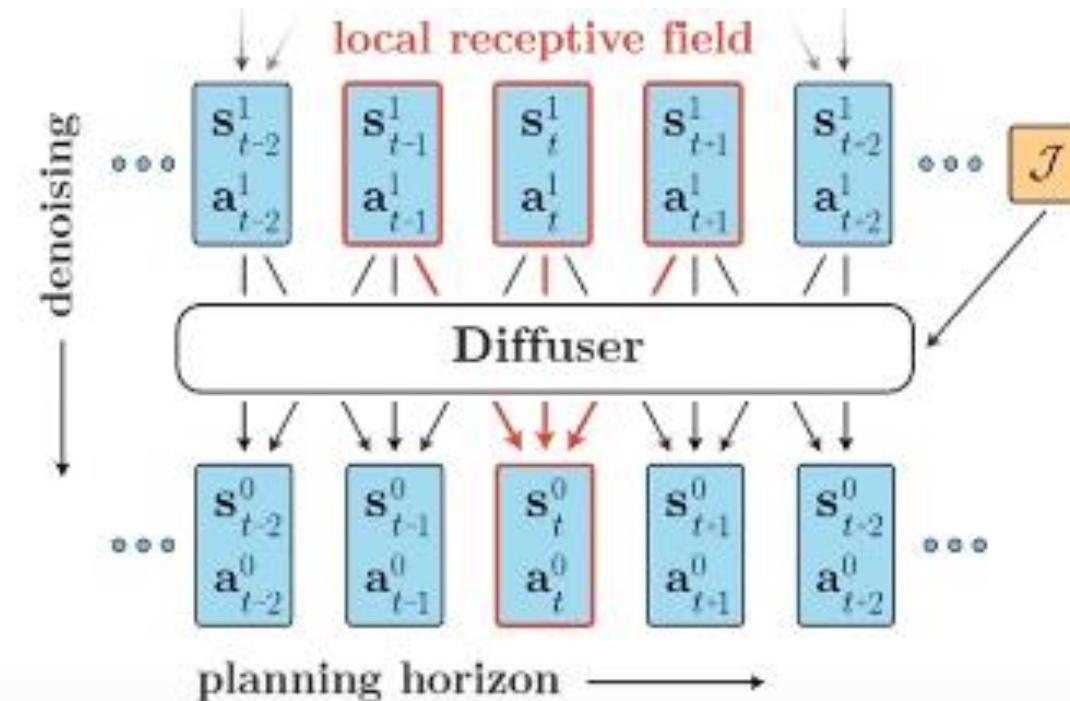


# Diffusion Planner

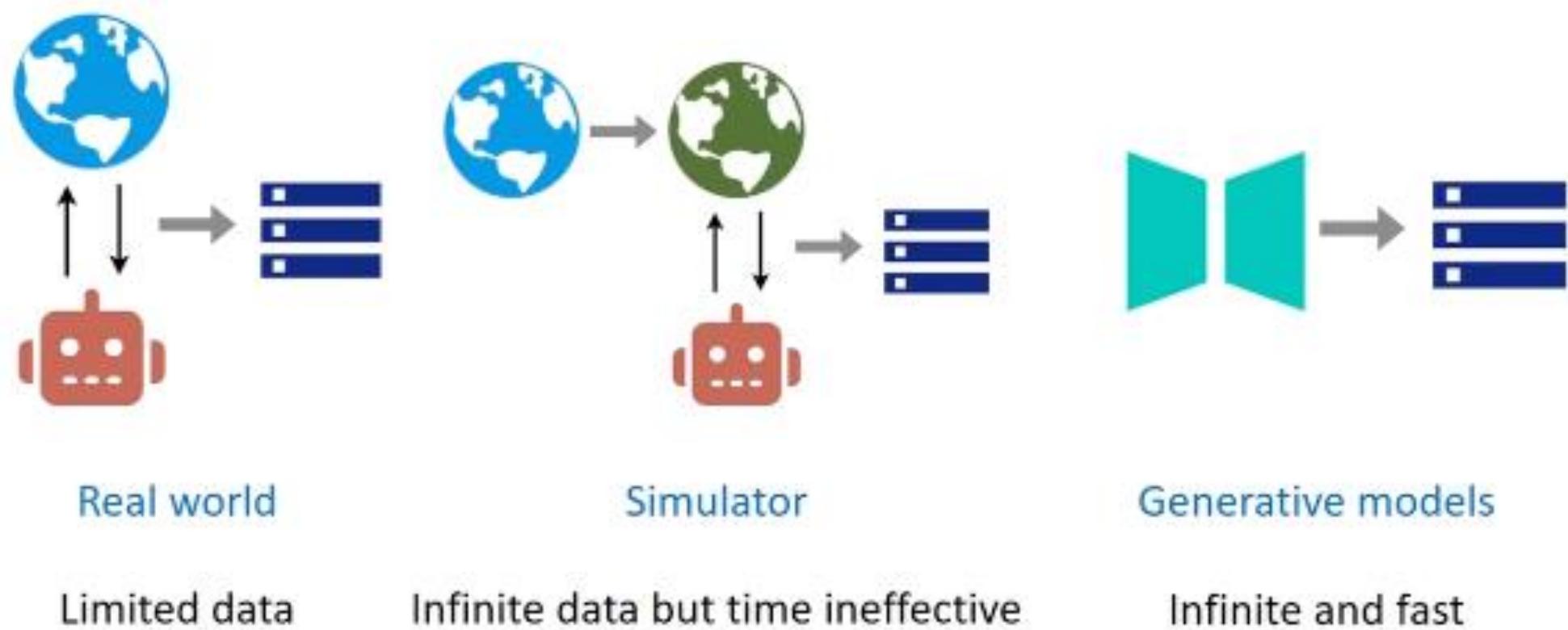
- Diffusion is just a data generator. So basically you have a lot of choices.
- Given  $s$  predict  $s'$
- Given  $s$  predict  $a$
- Given  $s$  predict  $s, a$
- Given  $s$  predict  $s, a, r$
- Given  $s$  predict a series of the previous variables

# Why is diffusion better at all?

- Key insight:
  - Action chunking: regularize locally, propagate globally
  - If we have an end goal, things would be much easier!



# Diffusion for Data Generation!



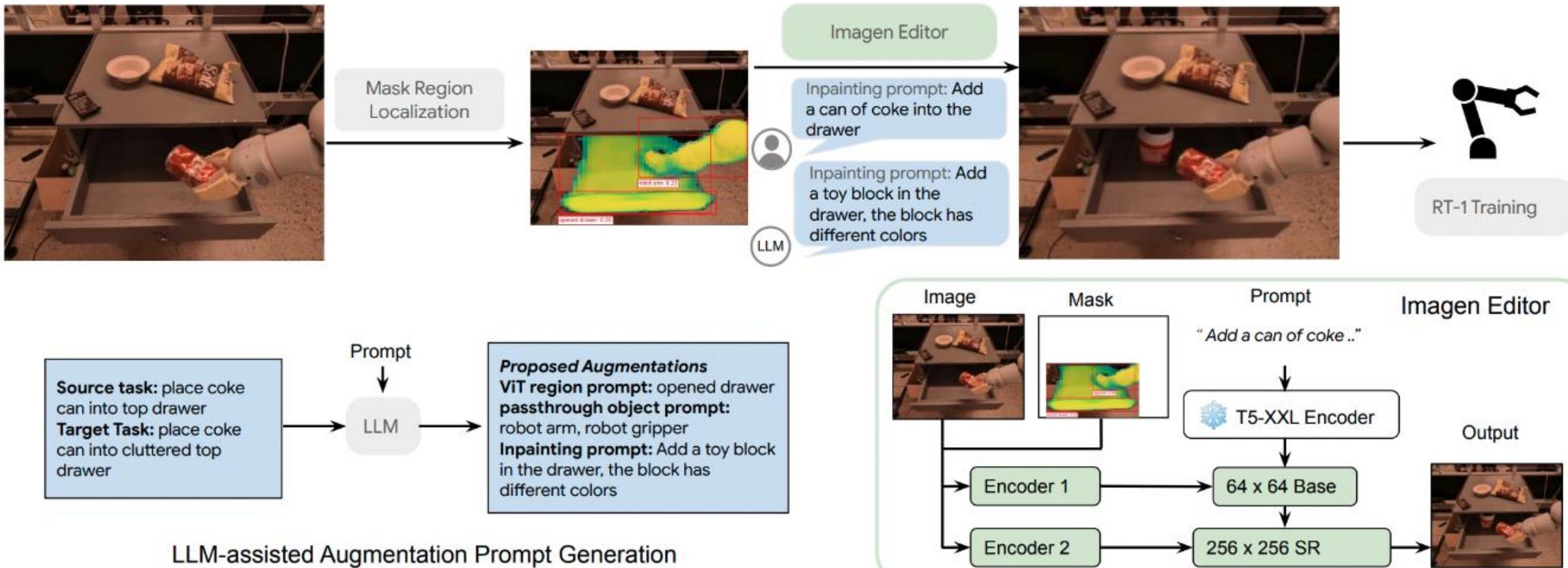
# Scaling Robot Learning with Semantically Imagined Experience

Tianhe Yu    Ted Xiao    Austin Stone    Jonathan Tompson  
Anthony Brohan    Su Wang    Jaspiar Singh    Clayton Tan    Dee M  
Jodilyn Peralta    Brian Ichter    Karol Hausman    Fei Xia



13 robots, 17 months,  
to collect 130k demonstrations

+ Generative models?







ROSIE Augmented Data



Real-world Rollout

*robot putting objects in drawer* → *robot putting objects in sink*



# In Lec9

1 Diffusion Models

2 Diffusion in Policy Learning

3 Flow Matching



# $\pi_0$ : Our First Generalist Policy

Published

October 31, 2024

Email

[research@physicalintelligence.company](mailto:research@physicalintelligence.company)

Paper

[\$\pi\_0\$ .pdf](#)

VLMs effectively transfer semantic knowledge from the web, but they are trained to output only discrete language tokens. Dexterous robot manipulation requires  $\pi_0$  to output motor commands at a high frequency, up to 50 times per second. To provide this level of dexterity, we developed a novel method to augment pre-trained VLMs with continuous action outputs via **flow matching**, a variant of diffusion models. Starting from diverse robot data and a VLM pre-trained on Internet-scale data, we train our vision-language-action **flow matching** model, which we can then post-train on high-quality robot data to solve a range of downstream tasks.

# Imagine a ball in the 3D space

- Starting position at  $x_0$
- Speed  $\mathbf{u}_t(\mathbf{x})$
- $$\frac{d\mathbf{x}}{dt} = \mathbf{u}_t(\mathbf{x}), \quad \phi_0(\mathbf{x}) = \mathbf{x}$$

# Now imagine a distribution of balls

- There initial distribution  $p_0$  at  $t = 0$ .
- There final distribution  $p_1$  at  $t = 1$ .
- If  $p_0$  is very easy to sample, and  $p_1$  is the target distribution.
- Then everything is solved!
- **Continuous Normalizing Flow (CNF)**

# Probability Path

$$p : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}_{>0}$$

- $p_t(x)$  is the probability density at time  $t$  and position  $x$ .
- $t \mapsto p_t$  is the path
- $\int p_t(\mathbf{x}) d\mathbf{x} = 1$  It is continuous. It is normalized.

# Continuity equation

- $\mathbf{u} : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  velocity field.
- $\mathbf{u}_t(\mathbf{x})$  is the velocity at time  $t$  and position  $x$ .
- The continuity equation:

$$\frac{\partial}{\partial t} p_t(\mathbf{x}) + \nabla \cdot (p_t(\mathbf{x}) \mathbf{u}_t(\mathbf{x})) = 0$$

- Intuition:
  - First term: particles amount changes
  - Second term: particles flow in/out
- Question:
  - Is  $p$  a one-to-one mapping to  $u$ ?
  - Considering adding a Solenoidal Vector Field to  $u$ ?

# Flow

- The trajectory of particles following  $u$  is called flow  $\phi : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$

- $$\frac{\partial}{\partial t} \phi_t(\mathbf{x}) = \mathbf{u}_t(\phi_t(\mathbf{x}))$$
$$\phi_0(\mathbf{x}) = \mathbf{x}$$

- you can do the integral

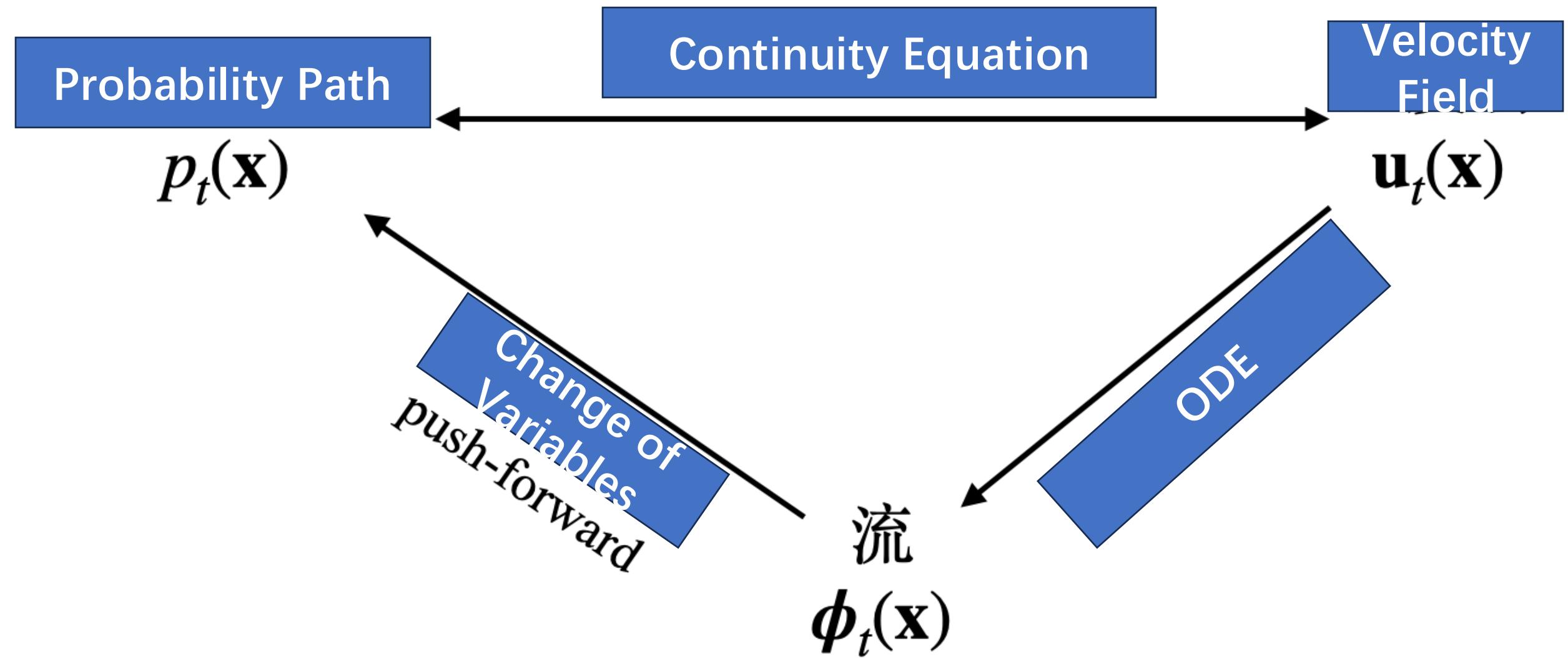
$$\phi_t(\mathbf{x}) - \mathbf{x} = \int_0^t \mathbf{u}_s(\phi_s(\mathbf{x})) ds$$

# Change of variables formula

- $\phi_t^{-1}(\mathbf{x})$  is the starting position of the particles at position  $x$  and time  $t$ .
- Following change of variables formula:

$$p_t(\mathbf{x}) = p_0(\phi_t^{-1}(\mathbf{x})) \left| \det \left[ \frac{\partial \phi_t^{-1}}{\partial \mathbf{x}}(\mathbf{x}) \right] \right|$$

- Push-forward equation



# Flow Matching

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, p_t(\mathbf{x})} \left[ \left\| \mathbf{v}_t^\theta(\mathbf{x}) - \mathbf{u}_t(\mathbf{x}) \right\|^2 \right]$$

- But  $u_t(x)$  is unknown.
- But we might be able to know the velocity given samples.

# Conditional Flow Matching

- $p_t(\mathbf{x} \mid \mathbf{x}_1)$  conditional probability path.
- Then we have:

$$p_t(\mathbf{x}) = \int p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1) d\mathbf{x}_1 = \mathbb{E}_{q(\mathbf{x}_1)} [p_t(\mathbf{x} \mid \mathbf{x}_1)]$$

- $\mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1)$  conditional velocity field.

$$\mathbf{u}_t(\mathbf{x}) = \int \mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1) \frac{p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1)}{p_t(\mathbf{x})} d\mathbf{x}_1 = \mathbb{E}_{p(\mathbf{x}_1 \mid \mathbf{x})} [\mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1)]$$

$$\begin{aligned} \frac{\partial p_t(\mathbf{x})}{\partial t} &= \int \left( \frac{\partial}{\partial t} p_t(\mathbf{x} \mid \mathbf{x}_1) \right) q(\mathbf{x}_1) d\mathbf{x}_1 \\ &= - \int \nabla \cdot (p_t(\mathbf{x} \mid \mathbf{x}_1) \mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1)) q(\mathbf{x}_1) d\mathbf{x}_1 \\ &= - \nabla \cdot \left( \int p_t(\mathbf{x} \mid \mathbf{x}_1) \mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1) d\mathbf{x}_1 \right) \\ &= - \nabla \cdot (p_t(\mathbf{x}) \mathbf{u}_t(\mathbf{x})) \end{aligned}$$

# CFM LOSS

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t,q(\mathbf{x}_1),p_t(\mathbf{x}|\mathbf{x}_1)} \left[ \left\| \mathbf{v}_t^\theta(\mathbf{x}) - \mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1) \right\|^2 \right]$$

$$\begin{aligned}\mathcal{L}_{\text{FM}}(\theta) &= \mathbb{E}_{t,p_t(\mathbf{x})} \left[ \left\| \mathbf{v}_t^\theta(\mathbf{x}) \right\|^2 \right] + \mathbb{E}_{t,p_t(\mathbf{x})} \left[ \left\| \mathbf{u}_t(\mathbf{x}) \right\|^2 \right] - 2\mathbb{E}_{t,p_t(\mathbf{x})} \left[ \langle \mathbf{v}_t^\theta(\mathbf{x}), \mathbf{u}_t(\mathbf{x}) \rangle \right] \\ \mathbb{E}_{p_t(\mathbf{x})} \left[ \langle \mathbf{v}_t^\theta(\mathbf{x}), \mathbf{u}_t(\mathbf{x}) \rangle \right] &= \int p_t(\mathbf{x}) \left\langle \mathbf{v}_t^\theta(\mathbf{x}), \int \mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1) \frac{p_t(\mathbf{x} \mid \mathbf{x}_1)q(\mathbf{x}_1)}{p_t(\mathbf{x})} d\mathbf{x}_1 \right\rangle d\mathbf{x} \\ &= \int \left\langle \mathbf{v}_t^\theta(\mathbf{x}), \int \mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1) p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1) d\mathbf{x}_1 \right\rangle d\mathbf{x} \\ &= \iint \langle \mathbf{v}_t^\theta(\mathbf{x}), \mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1) \rangle p_t(\mathbf{x} \mid \mathbf{x}_1) q(\mathbf{x}_1) d\mathbf{x}_1 d\mathbf{x} \\ &= \mathbb{E}_{q(\mathbf{x}_1),p_t(\mathbf{x}|\mathbf{x}_1)} \left[ \langle \mathbf{v}_t^\theta(\mathbf{x}), \mathbf{u}_t(\mathbf{x} \mid \mathbf{x}_1) \rangle \right]\end{aligned}$$

$$\mathcal{L}_{\text{CFM}}'(\theta) = \mathbb{E}_{t,q(\mathbf{x}_1),p_t(\mathbf{x}|\mathbf{x}_1)}\left[\left\|\mathbf{v}_t^{\theta}(\mathbf{x})\right\|^2 - 2\left\langle\mathbf{v}_t^{\theta}(\mathbf{x}), \mathbf{u}_t\left(\mathbf{x} \mid \mathbf{x}_1\right)\right\rangle\right]$$

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t,q(\mathbf{x}_1),p_t(\mathbf{x}|\mathbf{x}_1)}\left[\left\|\mathbf{v}_t^{\theta}(\mathbf{x}) - \mathbf{u}_t\left(\mathbf{x} \mid \mathbf{x}_1\right)\right\|^2\right]$$

# General form of conditional velocity field

$$p_t(\mathbf{x} \mid \mathbf{x}_1) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_t(\mathbf{x}_1), \sigma_t^2(\mathbf{x}_1)\mathbf{I})$$

$$p_0(\mathbf{x} \mid \mathbf{x}_1) = p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$$

$$p_1(\mathbf{x} \mid \mathbf{x}_1) = \mathcal{N}(\mathbf{x}; \mathbf{x}_1, \sigma_{\min}^2 \mathbf{I})$$

$$\boldsymbol{\psi}_t(\mathbf{x}) = \sigma_t(\mathbf{x}_1)\mathbf{x} + \boldsymbol{\mu}_t(\mathbf{x}_1)$$

$$\begin{aligned}
p_0\left(\boldsymbol{\psi}_t^{-1}(\mathbf{x}) \mid \mathbf{x}_1\right)\left|\det \left[\frac{\partial \boldsymbol{\psi}_t^{-1}}{\partial \mathbf{x}}(\mathbf{x})\right]\right| &= p\left(\boldsymbol{\psi}_t^{-1}(\mathbf{x})\right)\left|\det \left[\frac{\partial \boldsymbol{\psi}_t^{-1}}{\partial \mathbf{x}}(\mathbf{x})\right]\right| \\
&= \mathcal{N}\left(\frac{\mathbf{x}-\boldsymbol{\mu}_t\left(\mathbf{x}_1\right)}{\sigma_t\left(\mathbf{x}_1\right)} ; \mathbf{0}, \mathbf{I}\right)\left|\frac{1}{\sigma_t^d\left(\mathbf{x}_1\right)}\right| \\
&= \mathcal{N}\left(\mathbf{x} ; \boldsymbol{\mu}_t\left(\mathbf{x}_1\right), \sigma_t^2\left(\mathbf{x}_1\right) \mathbf{I}\right) \\
&= p_t\left(\mathbf{x} \mid \mathbf{x}_1\right)
\end{aligned}$$

$$\frac{\partial}{\partial t} \boldsymbol{\psi}_t(\mathbf{x}) = \mathbf{u}_t\left(\boldsymbol{\psi}_t(\mathbf{x}) \mid \mathbf{x}_1\right)$$

$$\begin{aligned}
\mathcal{L}_{\text{CFM}}(\theta) &= \mathbb{E}_{t, q(\mathbf{x}_1), p(\mathbf{x}_0)}\left[\left\|\mathbf{v}_t^{\theta}\left(\boldsymbol{\psi}_t\left(\mathbf{x}_0\right)\right)-\mathbf{u}_t\left(\boldsymbol{\psi}_t\left(\mathbf{x}_0\right) \mid \mathbf{x}_1\right)\right\|^2\right] \\
&= \mathbb{E}_{t, q(\mathbf{x}_1), p(\mathbf{x}_0)}\left[\left\|\mathbf{v}_t^{\theta}\left(\boldsymbol{\psi}_t\left(\mathbf{x}_0\right)\right)-\frac{\partial}{\partial t} \boldsymbol{\psi}_t\left(\mathbf{x}_0\right)\right\|^2\right] \\
&= \mathbb{E}_{t, q(\mathbf{x}_1), p(\mathbf{x}_0)}\left[\left\|\mathbf{v}_t^{\theta}\left(\boldsymbol{\psi}_t\left(\mathbf{x}_0\right)\right)-\sigma_t'\left(\mathbf{x}_1\right) \mathbf{x}_0-\boldsymbol{\mu}_t'\left(\mathbf{x}_1\right)\right\|^2\right]
\end{aligned}$$

# 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 ☺

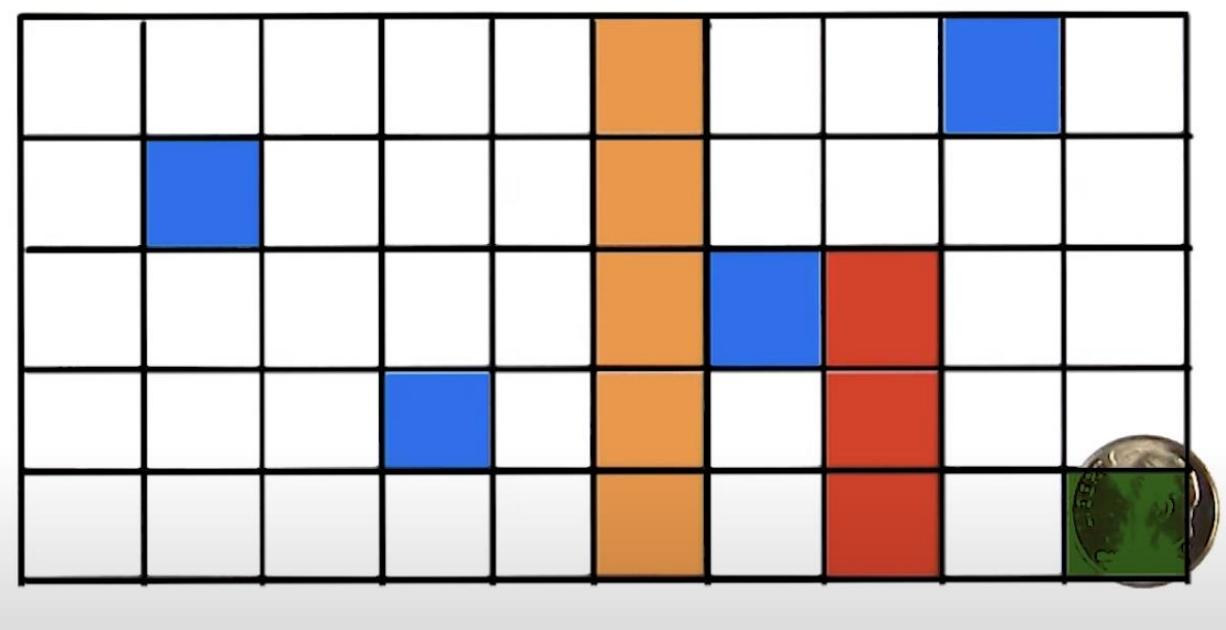
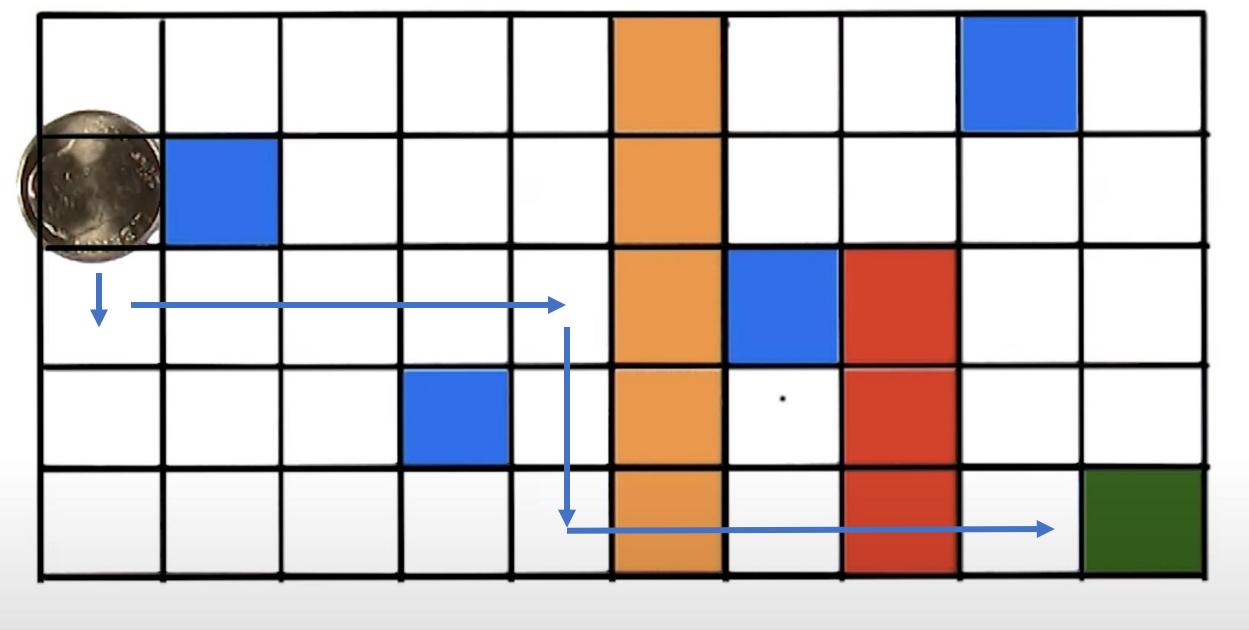


# In Lec9

- 1 Diffusion Models
- 2 Diffusion in Policy Learning
- 3 Flow Matching

Recitation starts!

# 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>

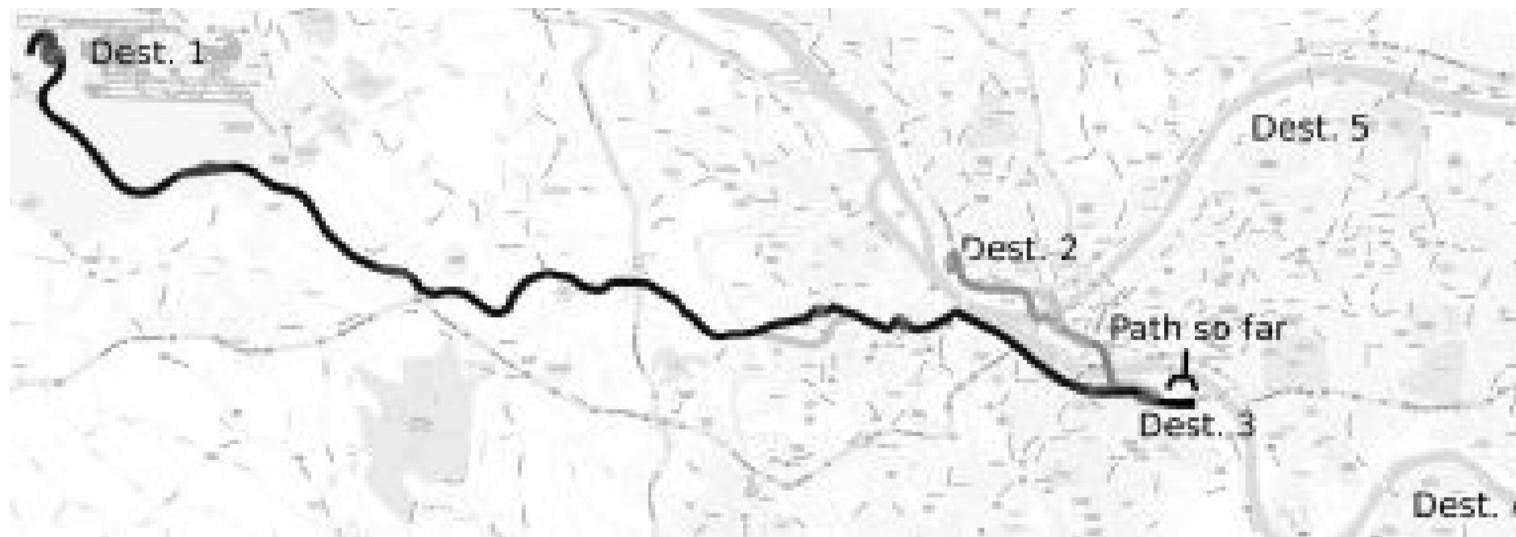
<sup>1</sup>Princeton Neuroscience Institute, Princeton University

<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

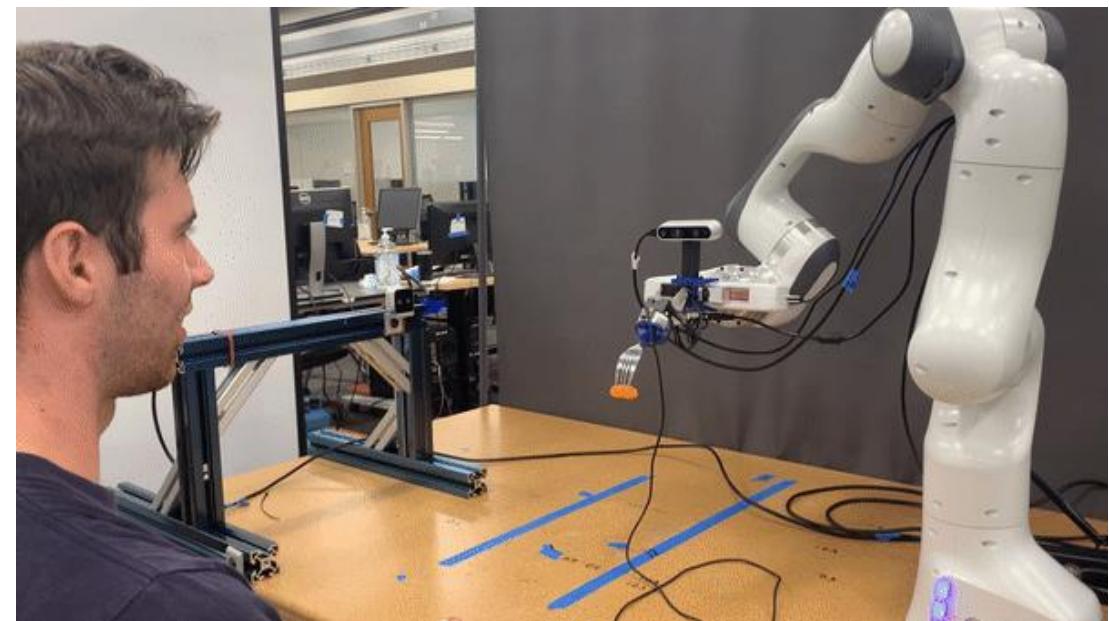
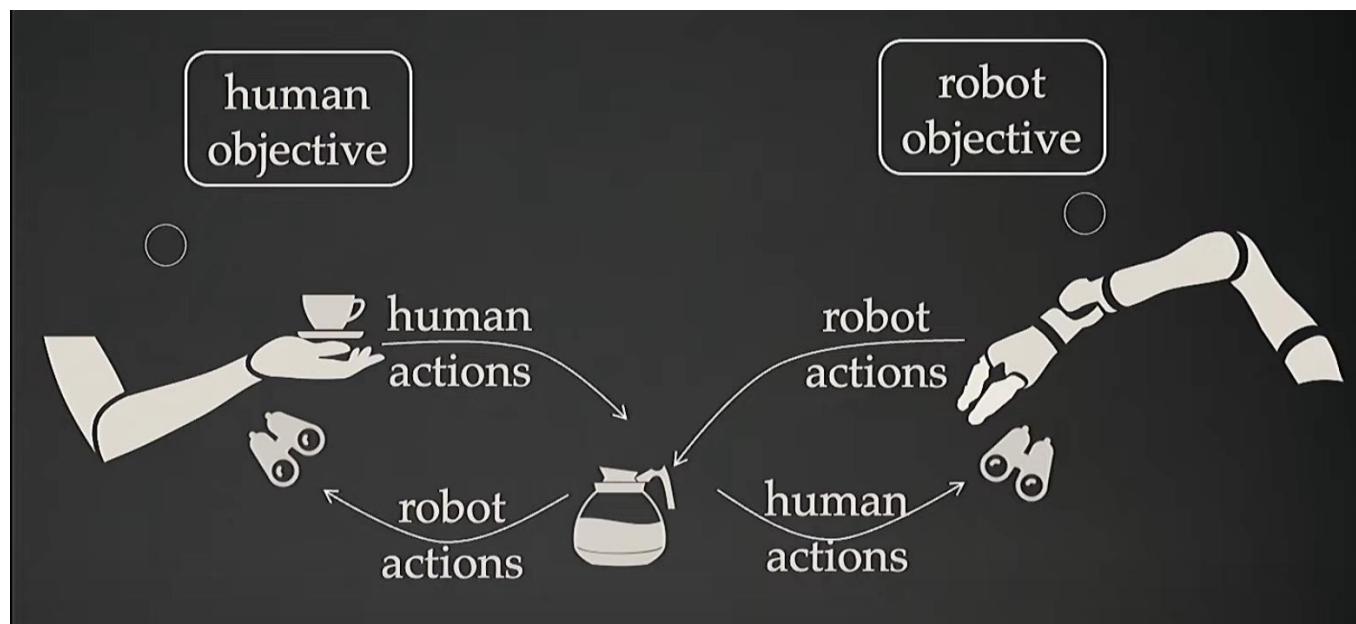
- Infer destination based of partial navigation trajectory



## 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  $\phi : S \rightarrow \mathfrak{R}^n$ .

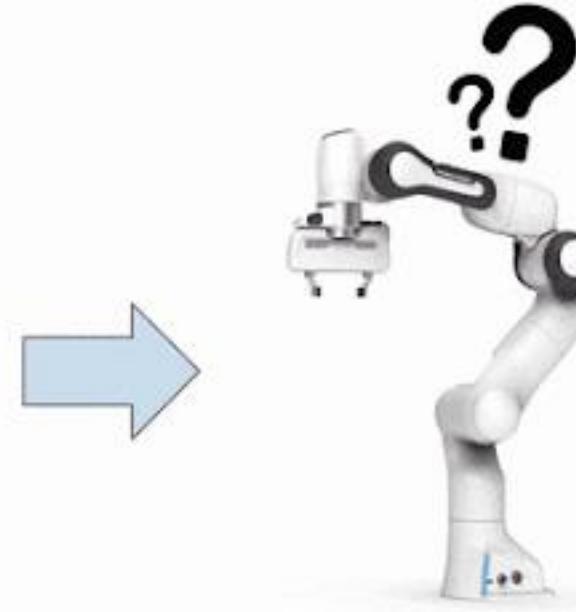
$$\begin{aligned} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi \right] &= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t w^\top \phi(s_t) \mid \pi \right] \\ &= w^\top \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \phi(s_t) \mid \pi \right] \\ &= w^\top \mu(\pi) \end{aligned}$$

- $\mu$  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?

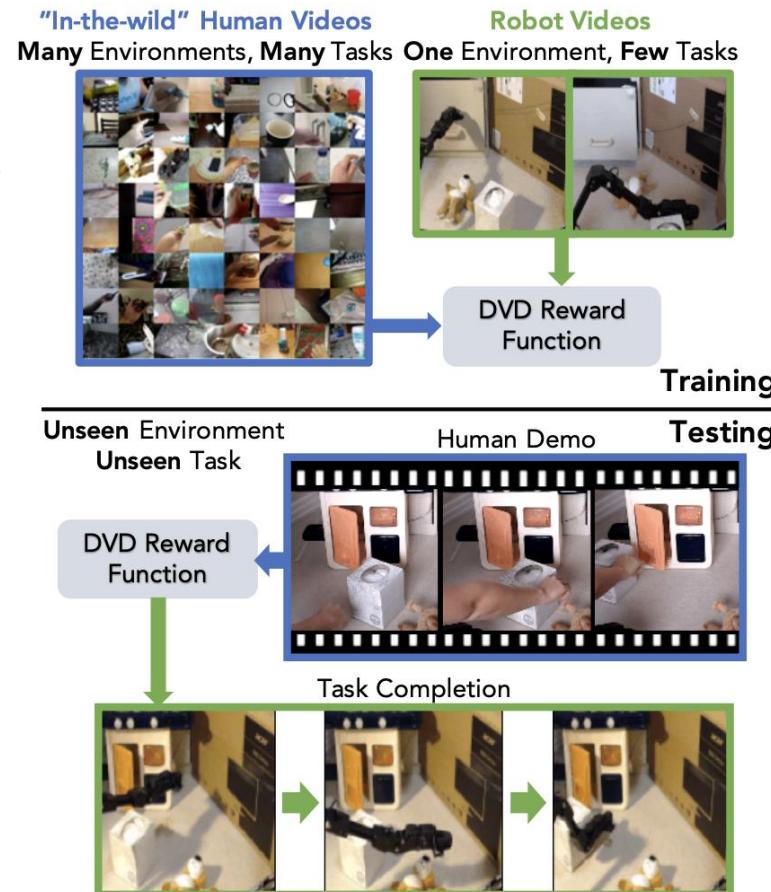
- Different domain!



# 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)

---

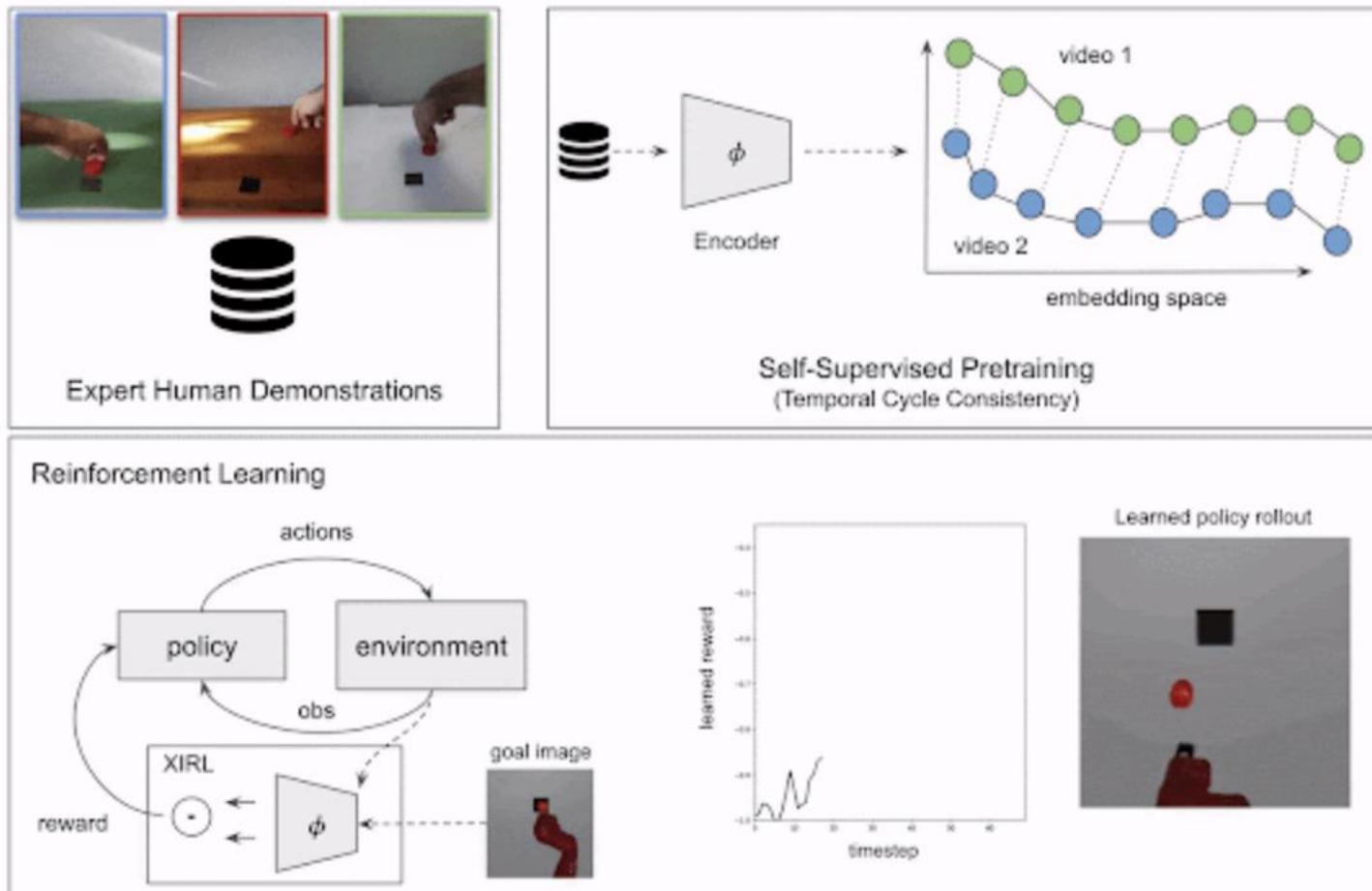
```
1: // Training DVD
2: Require:  $\mathcal{D}^h$  human demonstration data for  $N$  tasks  $\{\mathcal{T}_n\}$ 
3: Require:  $\mathcal{D}^r$  robot demonstration data for  $M$  tasks  $\{\mathcal{T}_m\} \subseteq \{\mathcal{T}_n\}$ 
4: Require: Pre-trained video encoder  $f_{enc}$ 
5: Randomly initialize  $\theta$ 
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>

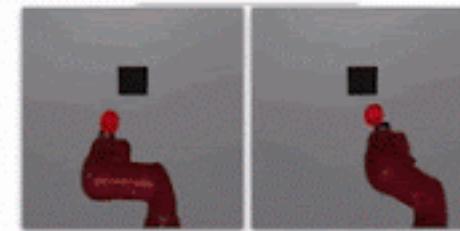
<sup>1</sup>Stanford University, <sup>2</sup>Robotics at Google, <sup>3</sup>UC Berkeley



# From IRL to RL

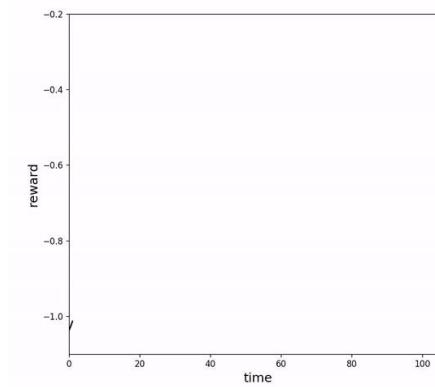
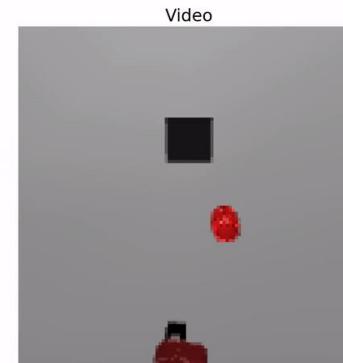
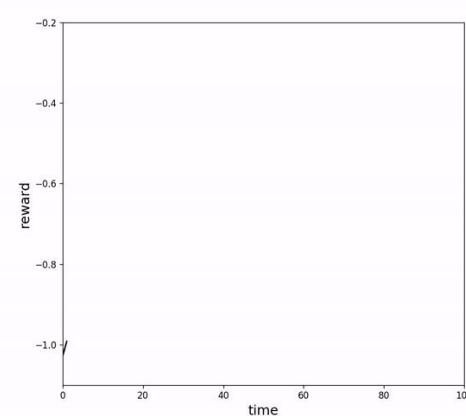
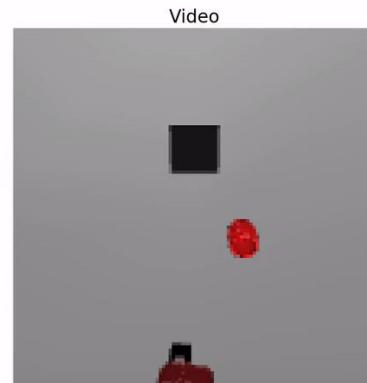
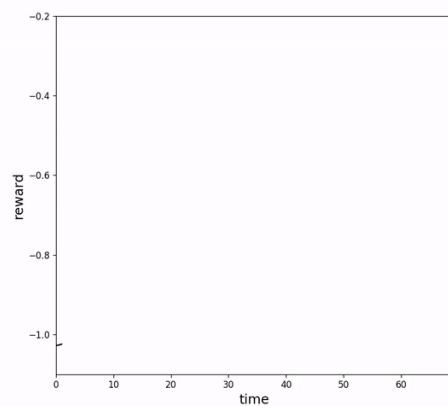
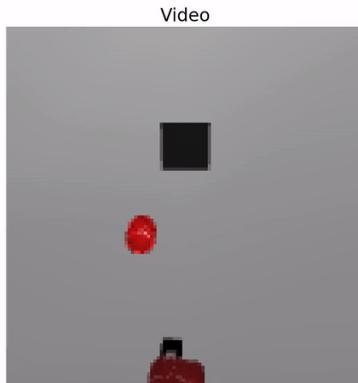


**Reward Learning**



**Reinforcement Learning**

# The XIRL-learned reward



---

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

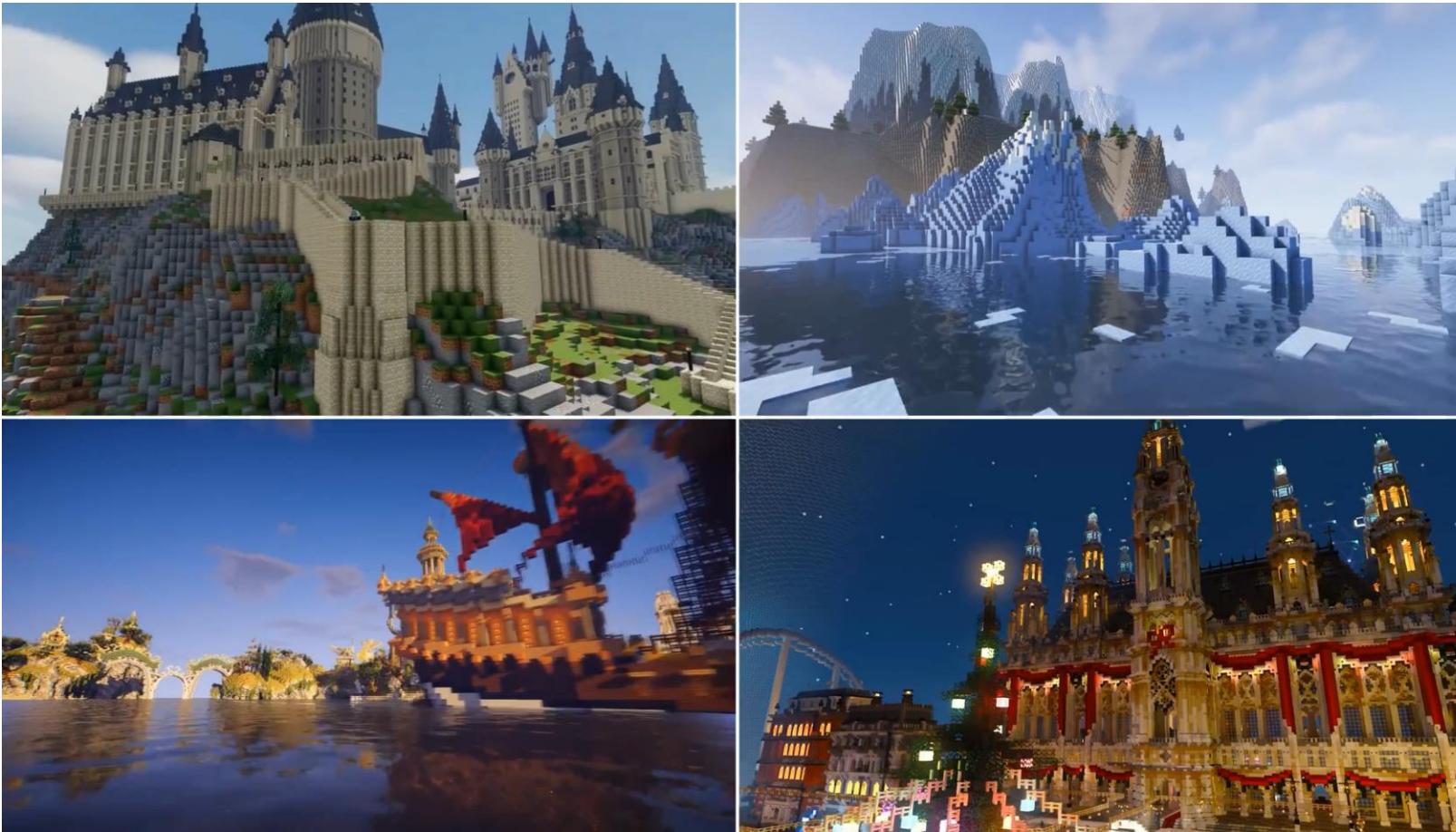
---

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

