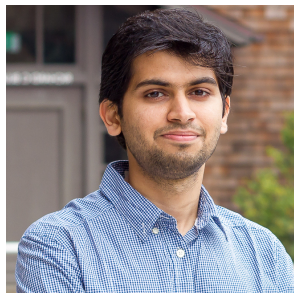


Self-Supervised Exploration via Disagreement



Deepak Pathak*
UC Berkeley



Dhiraj Gandhi*
CMU



Abhinav Gupta
CMU, FAIR

ICML 2019

* equal contribution



Exploration – a major challenge!

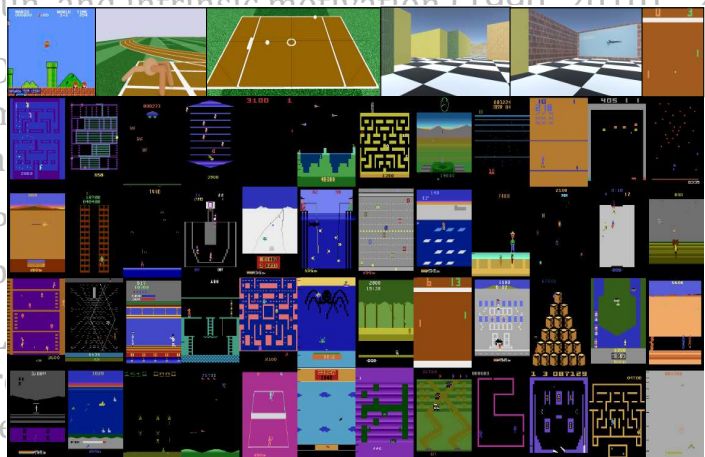
Exploration – a major challenge!

- Schmidhuber, Jurgen. “A possibility for implementing curiosity and boredom in model building neural controllers”, 1991.
- Schmidhuber, Jurgen. “Formal theory of creativity, fun, and intrinsic motivation (1990–2010)”, 2010.
- Oudeyer, P.-Y. and Kaplan, F. What is intrinsic motivation? a typology of computational approaches. *Frontiers in neurorobotics*, 2009.
- Poupart *et.al.* “An analytic solution to discrete bayesian reinforcement learning”. *ICML*, 2006.
- Lopes *et.al.* “Exploration in model-based reinforcement learning by empirically estimating learning progress”. *NIPS*, 2012.
- Bellemare *et.al.* “Unifying count-based exploration and intrinsic motivation”. *NIPS*, 2016.
- Mohamed *et.al.* “Variational information maximisation for intrinsically motivated reinforcement learning”. *NIPS*, 2015.
- Houthoofd *et.al.* “VIME: Variational information maximizing exploration”. *NIPS*, 2016.
- Gregor *et.al.* “Variational intrinsic control”. *ICLR Workshop*, 2017.
- Pathak *et.al.* “Curiosity-driven Exploration by Self-supervised Exploration”. *ICML* 2017
- Ostrovski *et.al.* “Count-based exploration with neural density models”. *ICML*, 2017.
- Burda*, Edwards*, Pathak* *et.al.* “Large-Scale Study of Curiosity-driven Learning”. *ICLR* 2019
- Eysenbach et al. “Diversity is all you need: Learn skills without a reward function”. *ICLR* 2019.
- Savinov et al. “Episodic curiosity through reachability”. *ICLR* 2019.

Exploration – a major challenge!

- Schmidhuber, Jurgen. “A possibility for implementing curiosity and boredom in model building neural controllers”, 1991.
- Schmidhuber, Jurgen. “Formal theory of creativity, fun, and intrinsic motivation (1990–2010)”, 2010.
- Oudeyer, P.-Y. and Kaplan, F. What is intrinsic motivation? a typology of computational approaches. *Frontiers in neurorobotics*, 2009.
- Poupart *et.al.* “An analytic solution to discrete bayesian reinforcement learning”. *ICML*, 2006.
- Lopes *et.al.* “Exploration in model-based reinforcement learning by empirically estimating learning progress”. *NIPS*, 2012.
- Bellemare *et.al.* “Unifying count-based exploration and intrinsic motivation”. *NIPS*, 2016.
- Mohamed *et.al.* “Variational information maximisation for intrinsically motivated reinforcement learning”. *NIPS*, 2015.
- Houthoofd *et.al.* “VIME: Variational information maximizing exploration”. *NIPS*, 2016.
- Gregor *et.al.* “Variational intrinsic control”. *ICLR Workshop*, 2017.
- Pathak *et.al.* “Curiosity-driven Exploration by Self-supervised Exploration”. *ICML 2017*
- Ostrovski *et.al.* “Count-based exploration with neural density models”. *ICML*, 2017.
- Burda*, Edwards*, Pathak* *et.al.* “Large-Scale Study of Curiosity-driven Learning”. *ICLR 2019*
- Eysenbach et al. “Diversity is all you need: Learn skills without a reward function”. *ICLR 2019*.
- Savinov et al. “Episodic curiosity through reachability”. *ICLR 2019*.

Exploration – a major challenge!



- Schmidhuber, Jürgen. "Formal theory of creativity, fun, and intrinsic motivation (1990, 2010)". 2010.
- Coates, Andrew. "Curiosity-driven exploration in reinforcement learning". 2010.
- Pathak, Dhruv, et al. "Curiosity-driven Exploration by Self-supervised Exploration". ICML 2017.
- Ostrovski, Alexander, et al. "Count-based exploration with neural density models". ICML, 2017.
- Burda, Lasse, et al. "Large-Scale Study of Curiosity-driven Learning". ICLR 2019.
- Eysenbach, Samuel, et al. "Diversity is all you need: Learn skills without a reward function". ICLR 2019.
- Savinov, Alexander, et al. "Episodic curiosity through reachability". ICLR 2019.
- Bellemare, Aaron G. "Unifying count-based exploration and intrinsic motivation". NIPS, 2016.

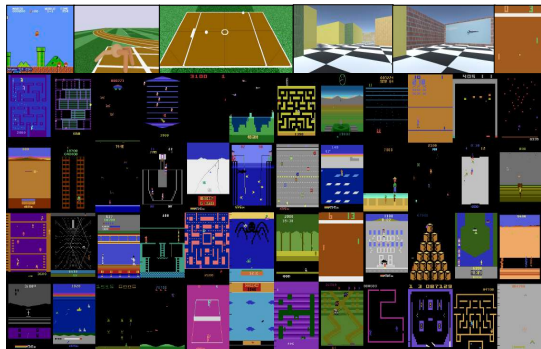
- Mohamed *et al.* "Variational information maximisation for intrinsically motivated reinforcement learning". NIPS, 2015.
- Houthoofd *et al.* "VIME: Variational information maximizing exploration". NIPS, 2016.
- Gregor *et al.* "Variational intrinsic control". ICLR Workshop, 2017.
- Pathak *et al.* "Curiosity-driven Exploration by Self-supervised Exploration". ICML 2017
- Ostrovski *et al.* "Count-based exploration with neural density models". ICML, 2017.
- Burda*, Edwards*, Pathak* *et al.* "Large-Scale Study of Curiosity-driven Learning". ICLR 2019
- Eysenbach et al. "Diversity is all you need: Learn skills without a reward function". ICLR 2019.
- Savinov et al. "Episodic curiosity through reachability". ICLR 2019.

Exploration – a major challenge!

- Schmidhuber, Jurgen. "A possibility for implementing curiosity and boredom in model building neural controllers", 1991.
- Schmidhuber, Jurgen. "Formal theory of creativity, fun, and intrinsic motivation (1990–2010)", 2016.
- Oudeyer, P.-Y. and Kaplan, F. What is intrinsic motivation? a typology of computational approaches. *Frontiers in neurorobotics*, 2017.
- Poupart *et.al.* "An analytic solution to bayesian reinforcement learning". *ICML*, 2007.
- Lopes *et.al.* "Exploration in reinforcement learning programs". *ICML*, 2008.
- Bellemare *et.al.* "Count-based exploration and intrinsic motivation". *NIPS*, 2016.
- Mohamed *et.al.* "Variational information maximisation for motivated reinforcement learning". *NIPS*, 2015.
- Oudeyer *et.al.* "Variational information maximisation for motivated exploration". *NIPS*, 2016.
- Oudeyer *et.al.* "Variational intrinsic control". *ICLR*, 2017.
- Sutton *et.al.* "Curiosity-driven Exploration by Self-supervised Exploration". *ICML* 2017
- Ostrovski *et.al.* "Count-based exploration with neural density models". *ICML*, 2017.
- Burda*, Edwards*, Pathak* *et.al.* "Large-Scale Study of Curiosity-driven Learning". *ICLR* 2019
- Eysenbach et al. "Diversity is all you need: Learn skills without a reward function". *ICLR* 2019.
- Savinov et al. "Episodic curiosity through reachability". *ICLR* 2019.

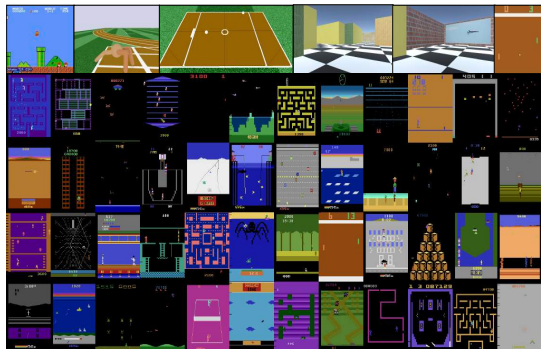
**Sample Inefficient
[millions of samples]**

Sample Inefficient

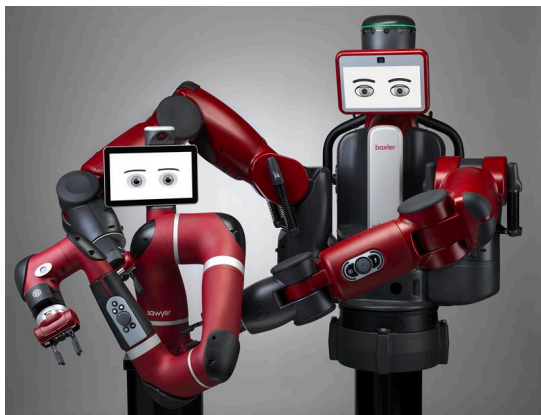


Simulation

Sample Inefficient

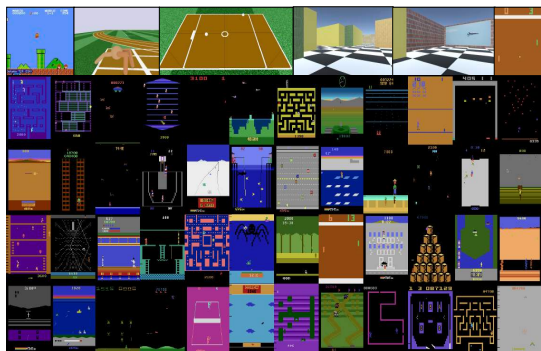


Simulation

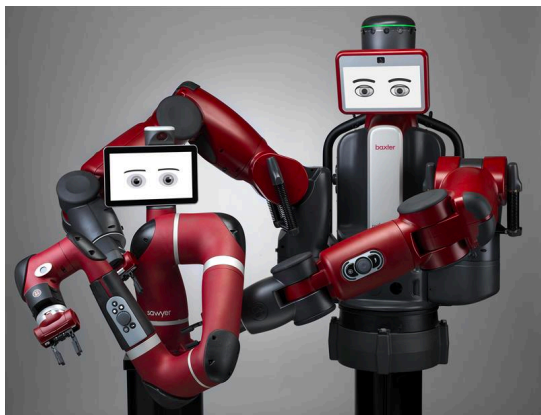


Real Robots

Sample Inefficient



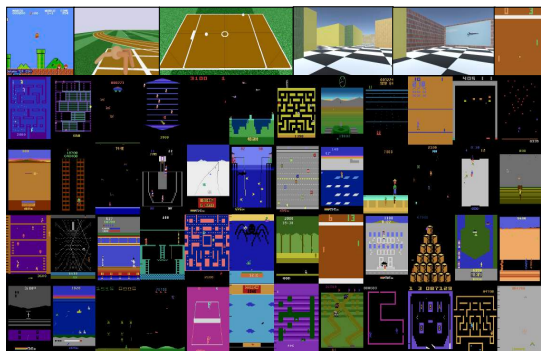
Simulation



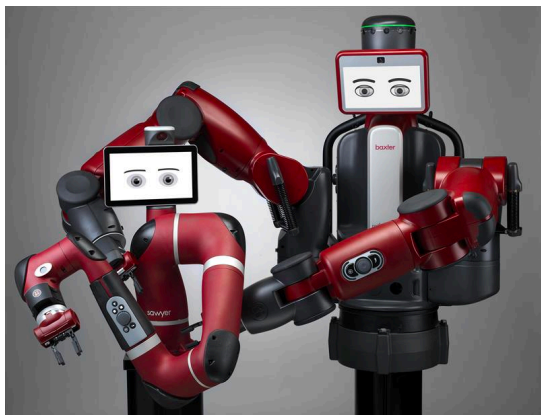
Real Robots

“Stuck” in Stochastic Envs

Sample Inefficient



Simulation



Real Robots

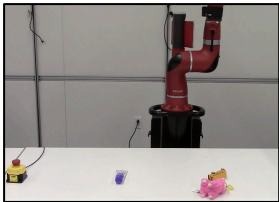
“Stuck” in Stochastic Envs



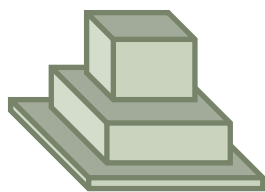
Curiosity Exploration
w/ Noisy TV & Remote

[Burda*, Edwards*, Pathak* et. al. ICLR'19]
[Juliani et.al., ArXiv'19]

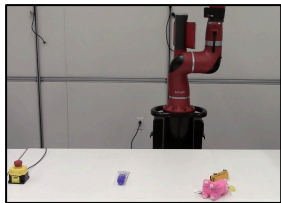
Why inefficient?



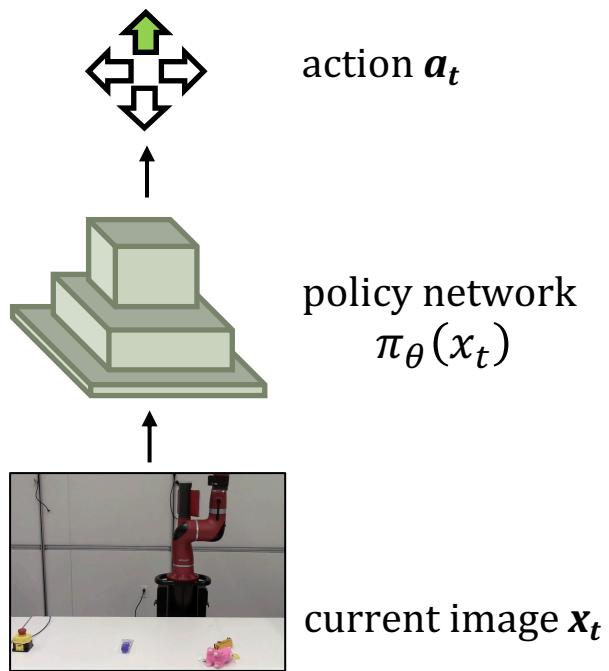
current image \mathbf{x}_t

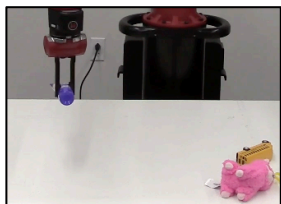


policy network
 $\pi_{\theta}(x_t)$



current image x_t

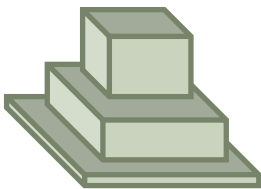




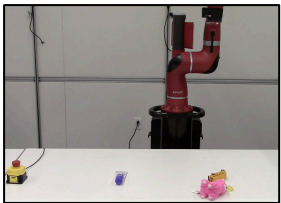
next image x_{t+1}



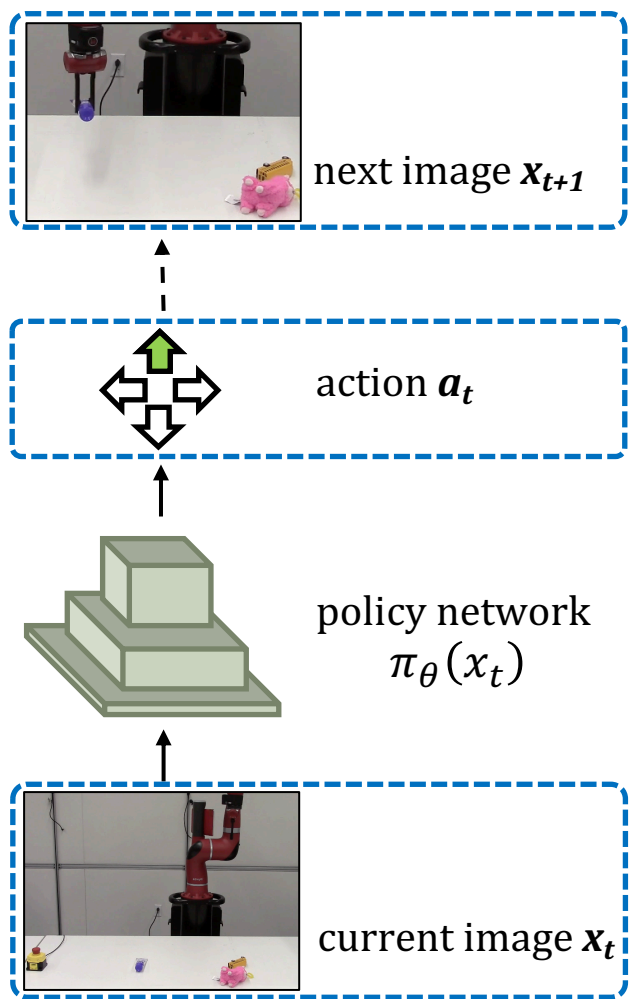
action a_t

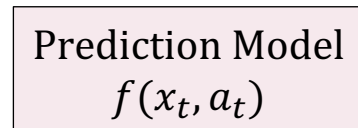
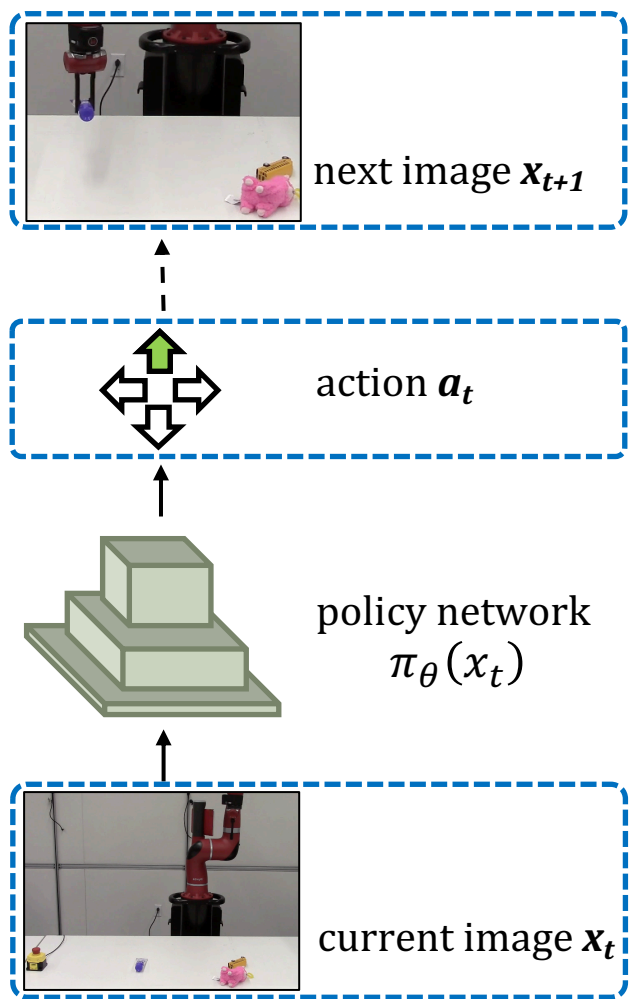


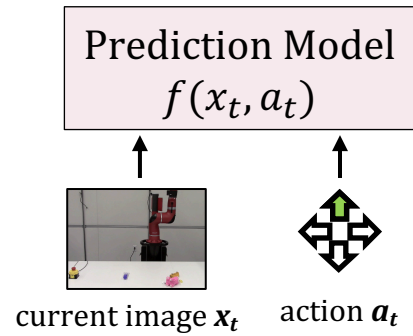
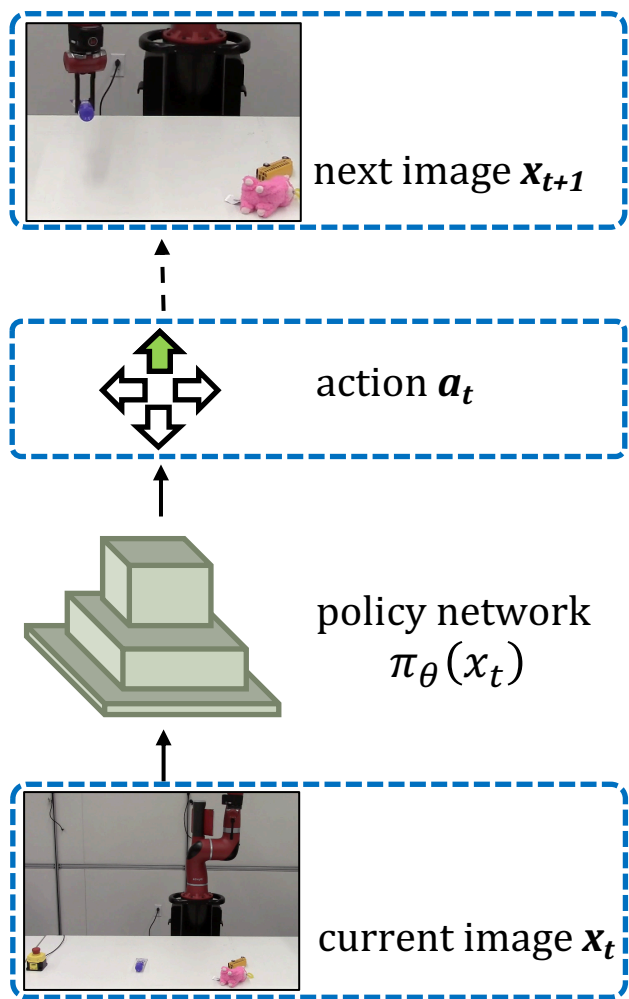
policy network
 $\pi_{\theta}(x_t)$

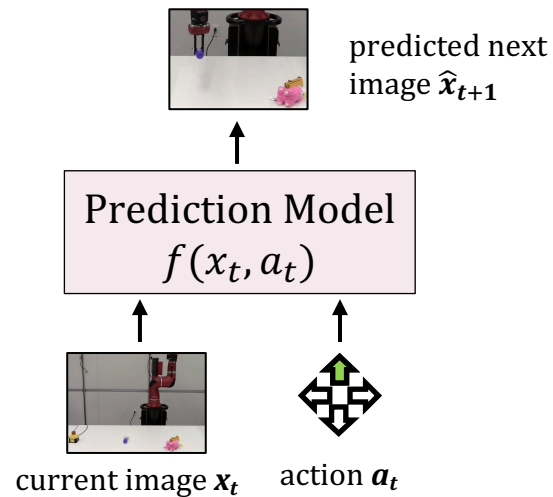
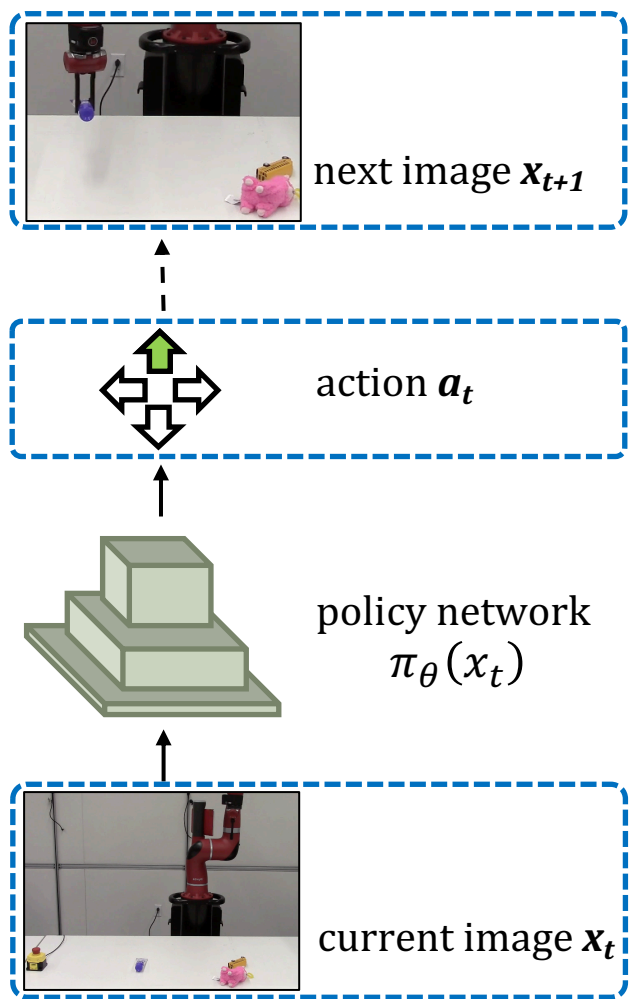


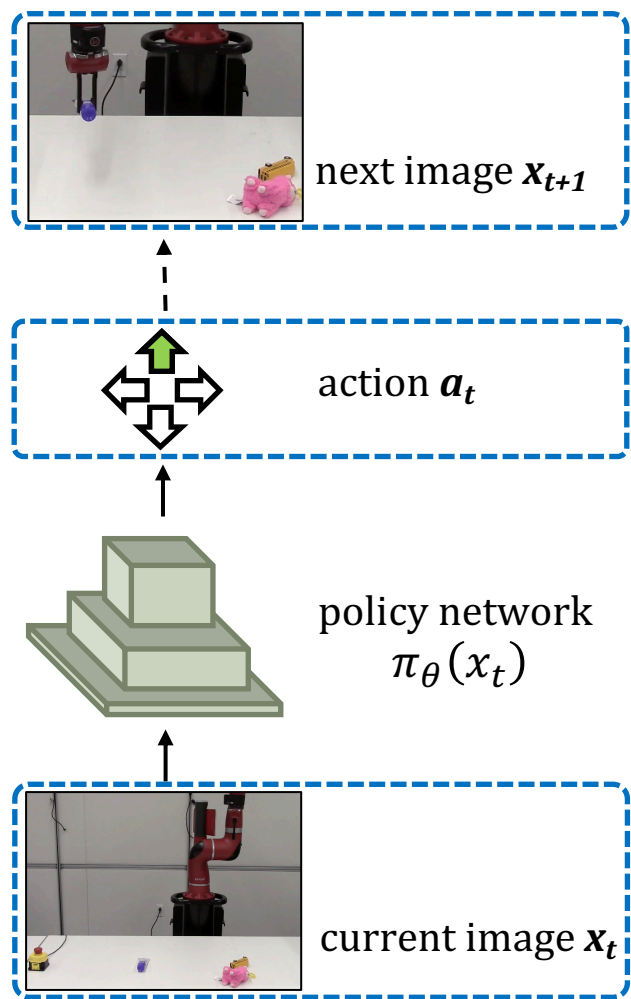
current image x_t



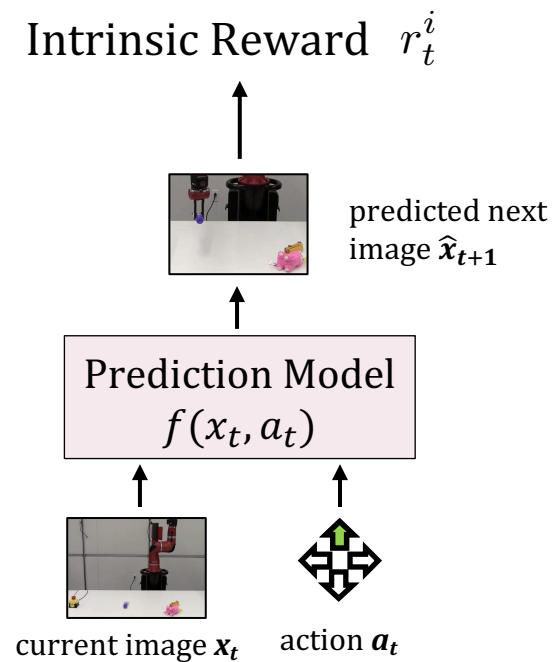


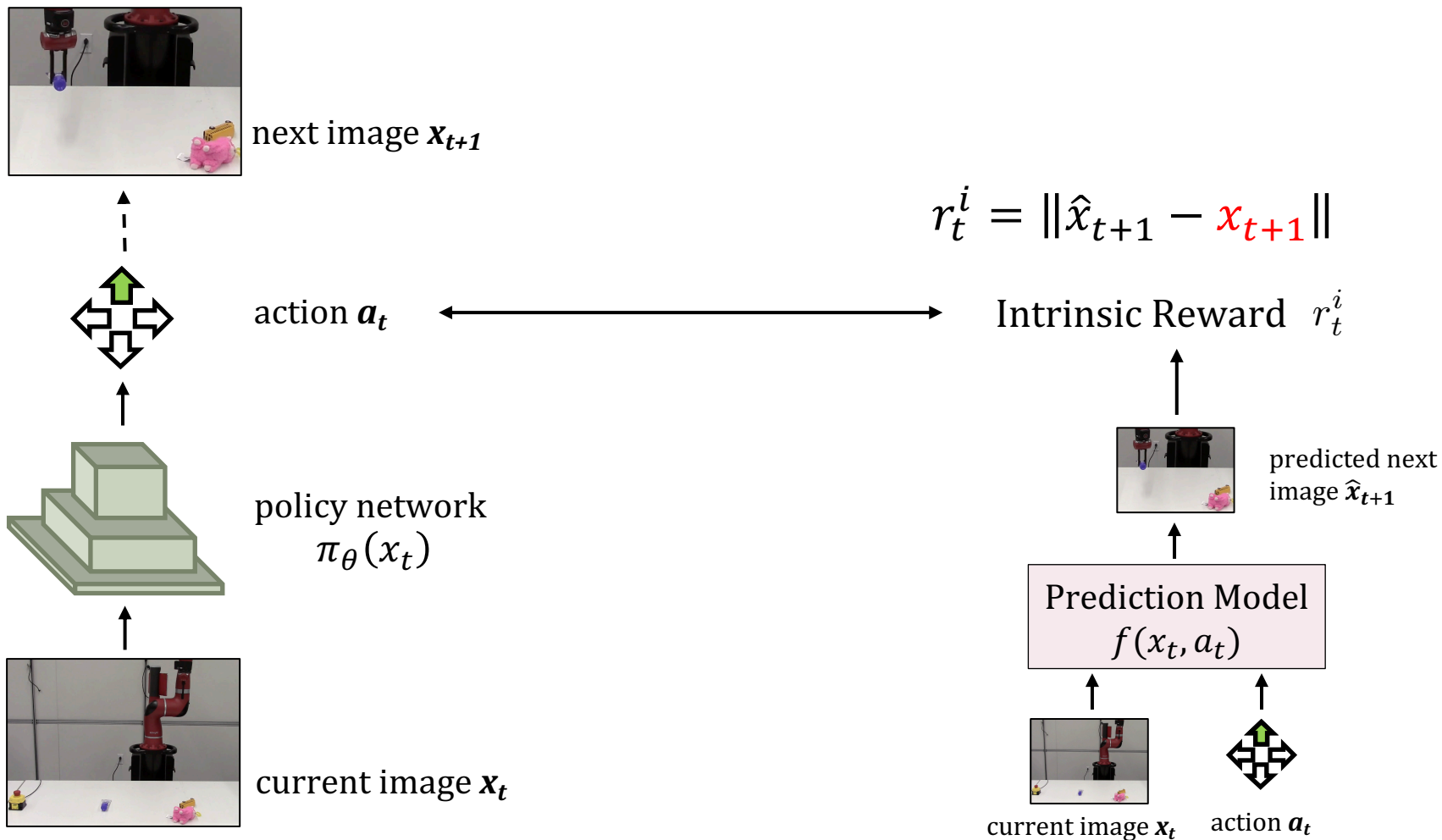




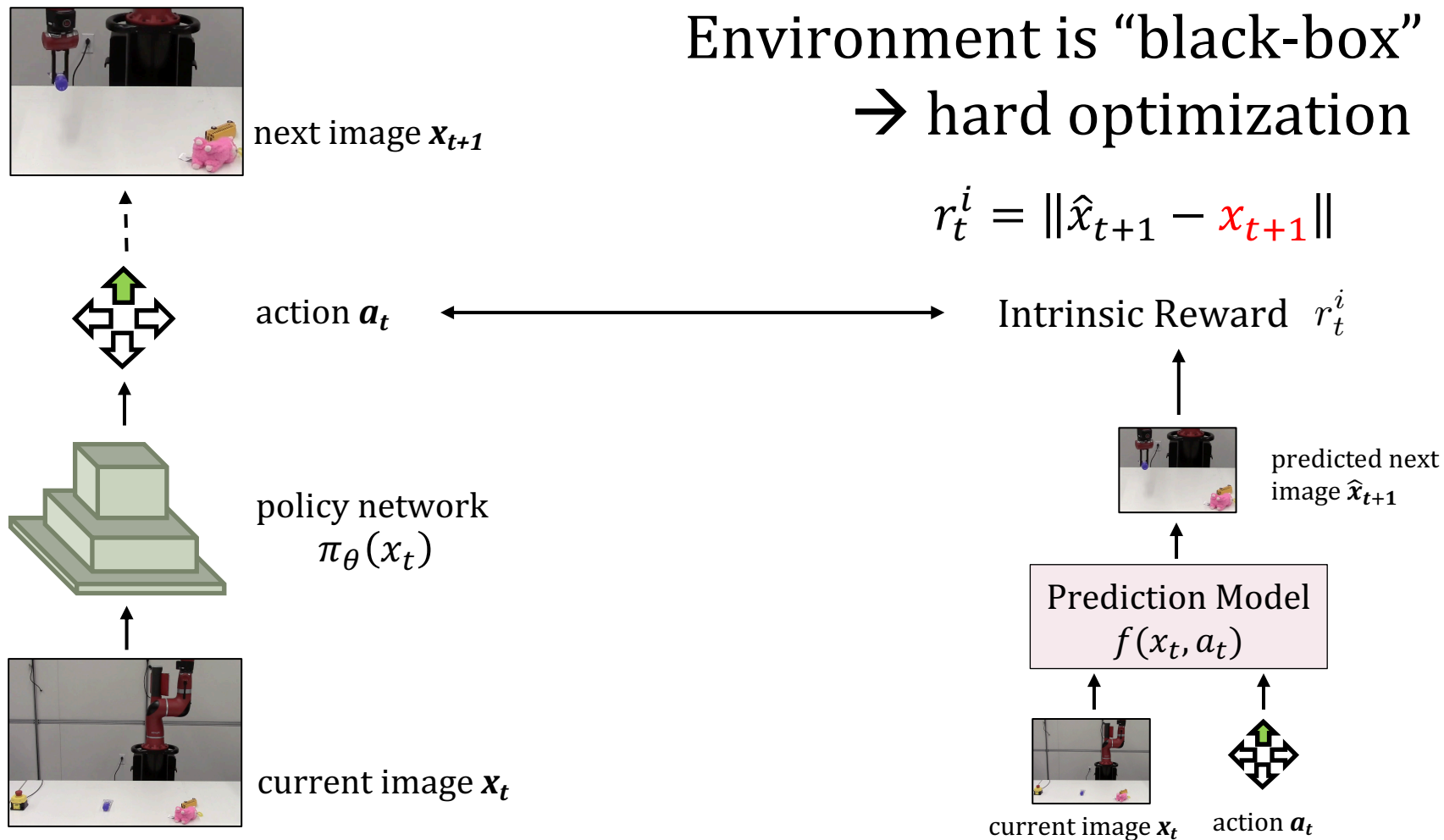


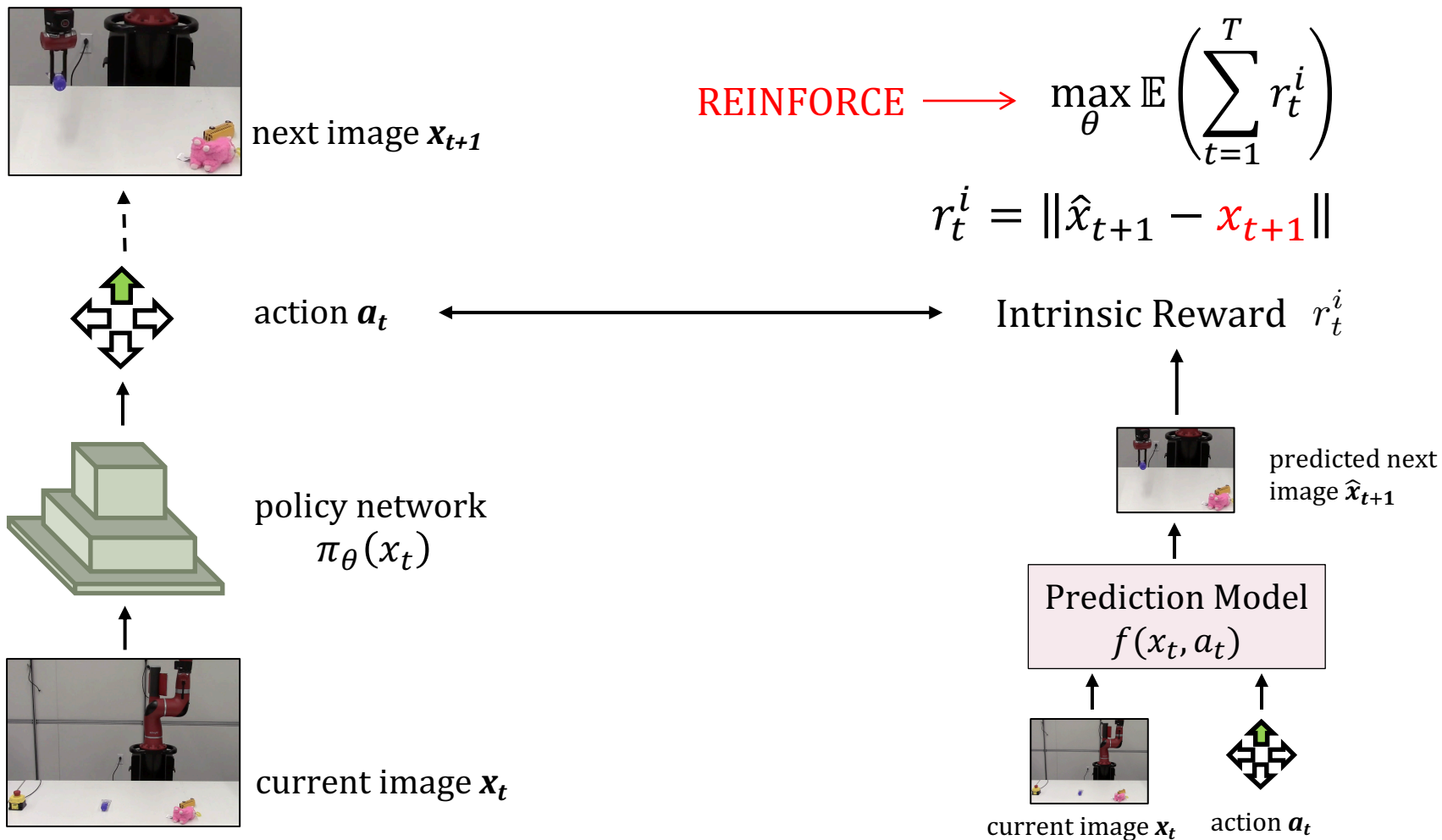
$$r_t^i = \|\hat{x}_{t+1} - x_{t+1}\|$$

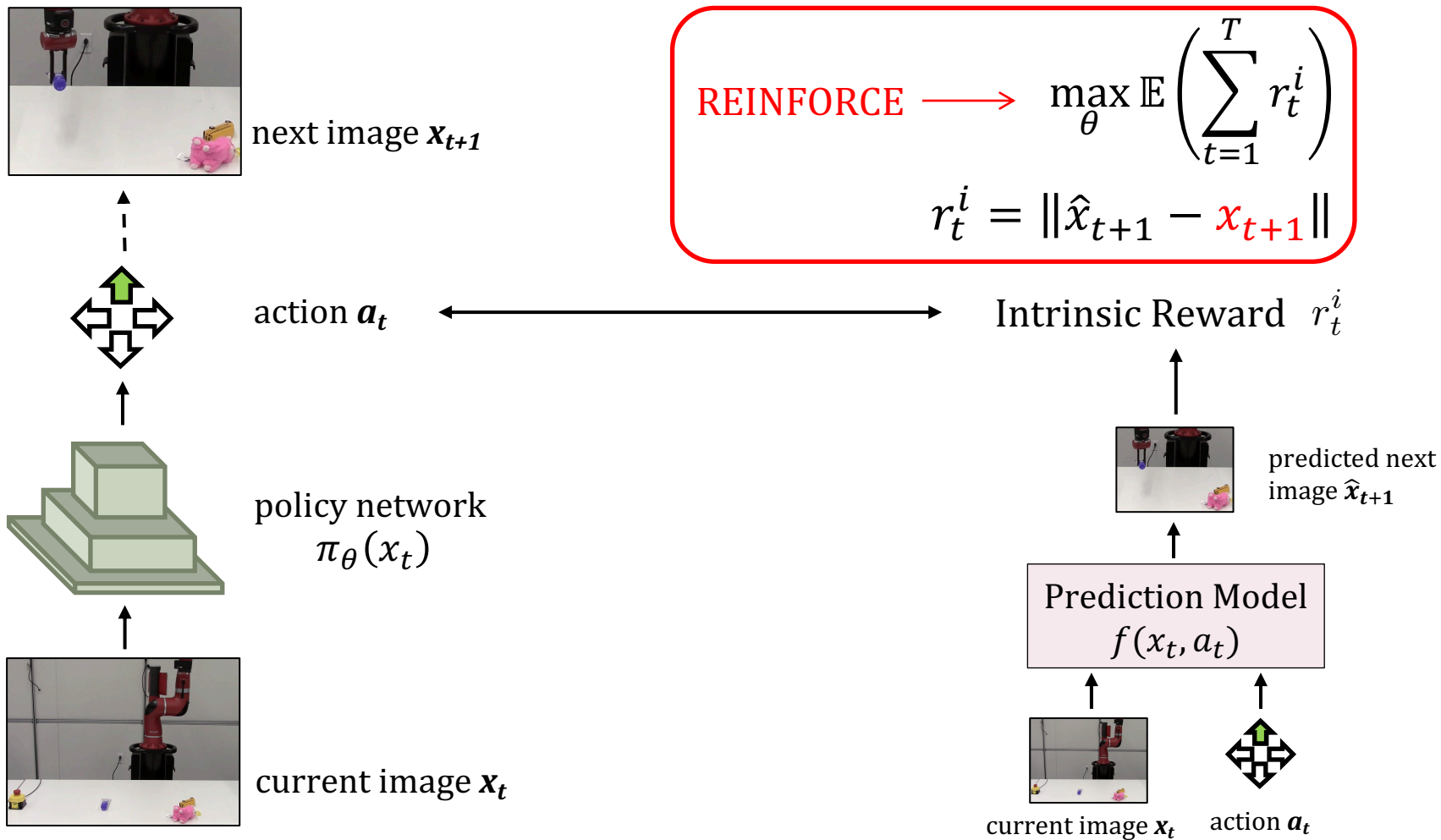


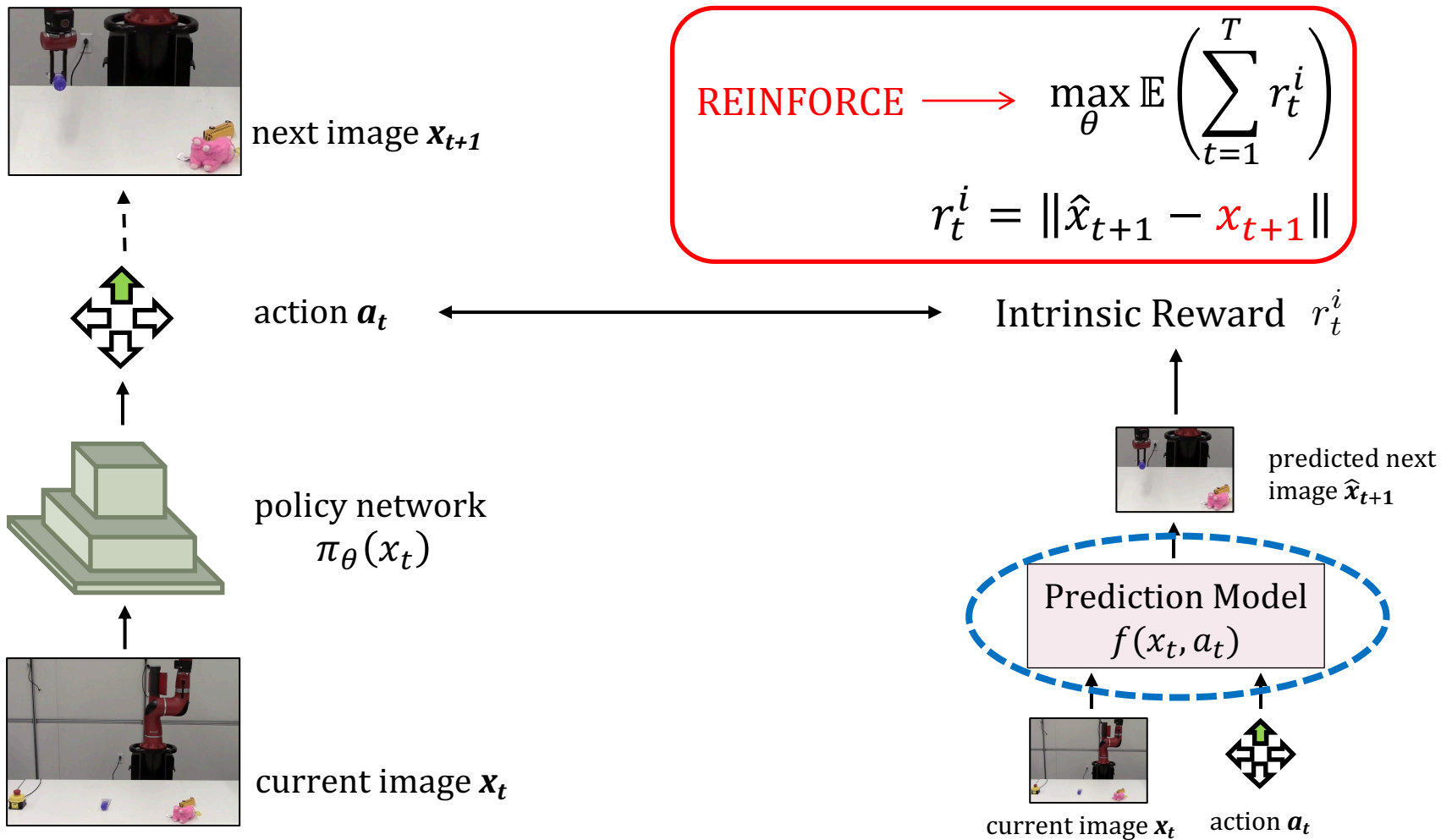


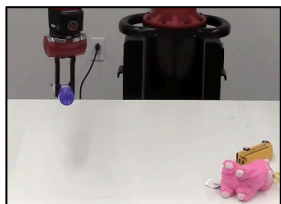
Environment is “black-box”
→ hard optimization



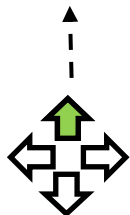




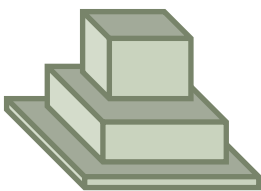




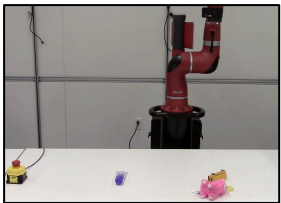
next image \mathbf{x}_{t+1}



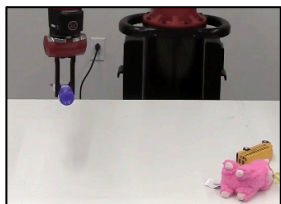
action \mathbf{a}_t



policy network
 $\pi_{\theta}(\mathbf{x}_t)$



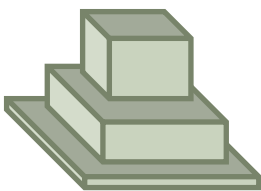
current image \mathbf{x}_t



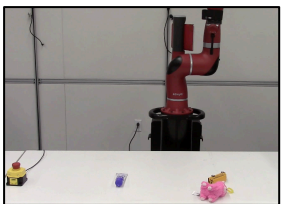
next image x_{t+1}



action a_t



policy network
 $\pi_{\theta}(x_t)$



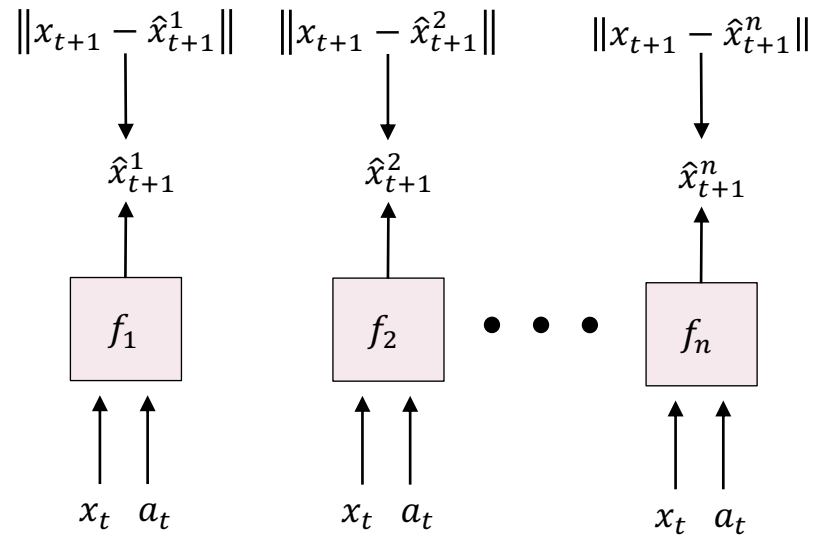
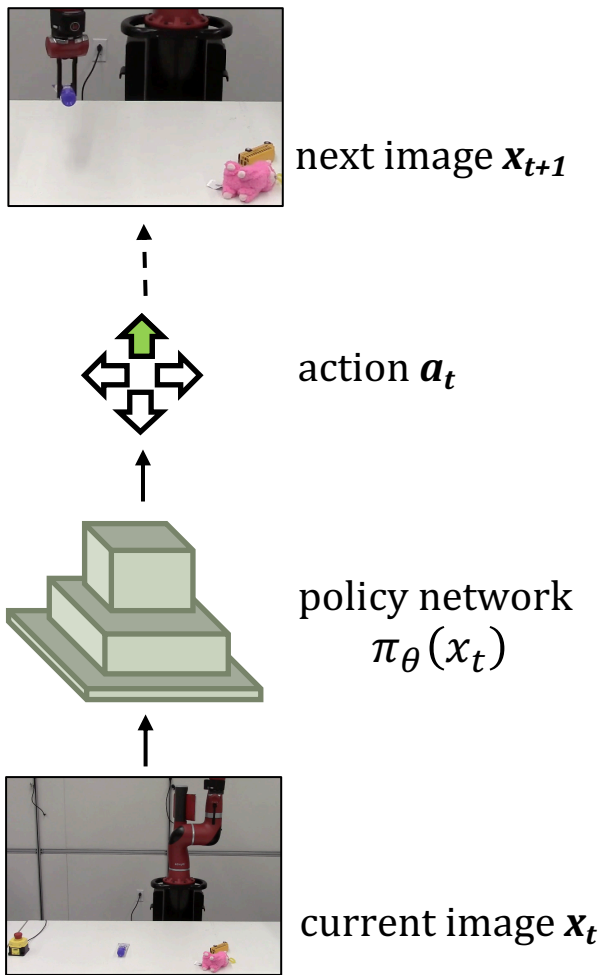
current image x_t

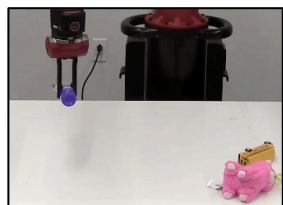
$$\|x_{t+1} - \hat{x}_{t+1}^1\|$$

$$\hat{x}_{t+1}^1$$

$$f_1$$

$$x_t \quad a_t$$

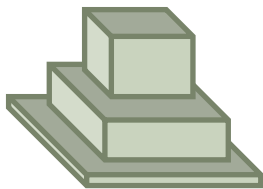




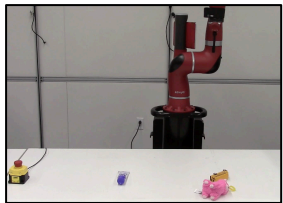
next image x_{t+1}



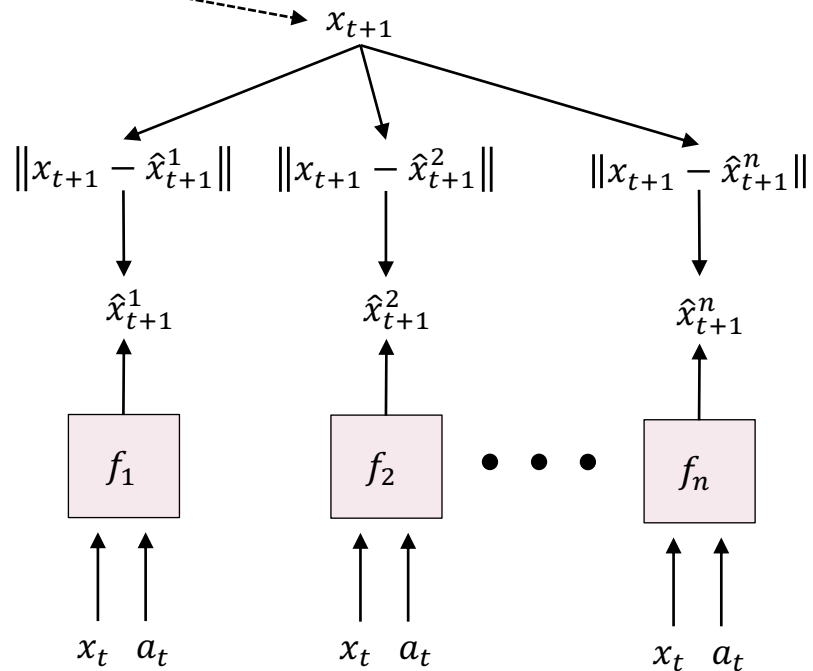
action a_t

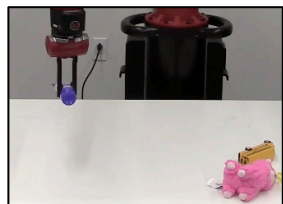


policy network
 $\pi_{\theta}(x_t)$



current image x_t

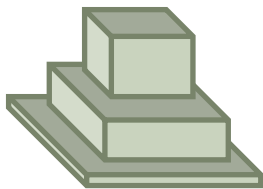




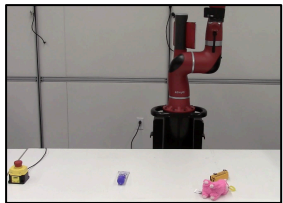
next image x_{t+1}



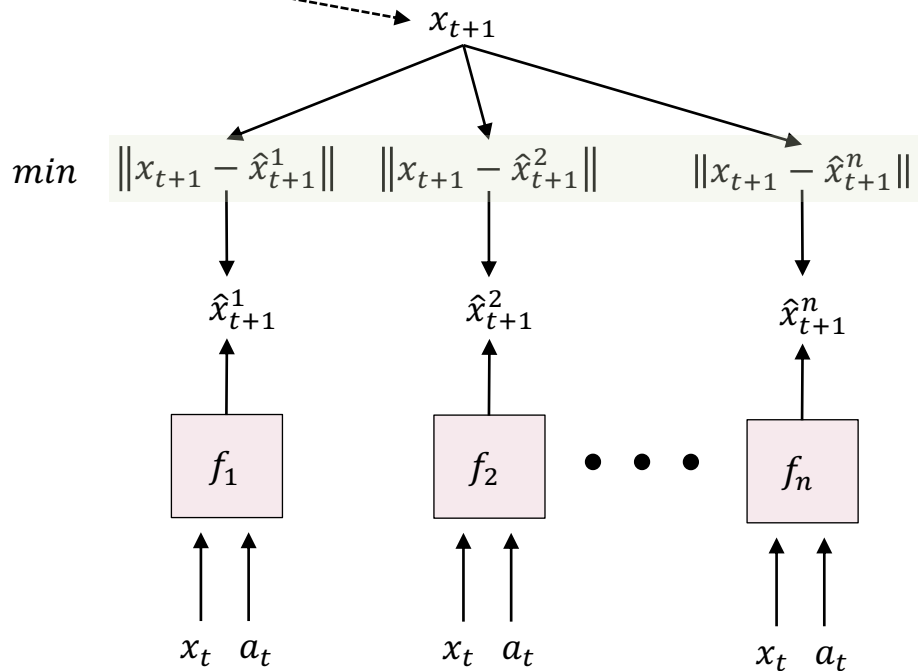
action a_t

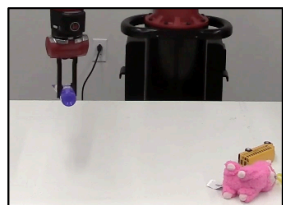


policy network
 $\pi_{\theta}(x_t)$



current image x_t

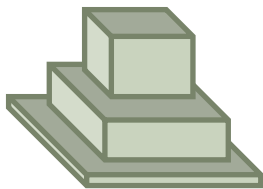




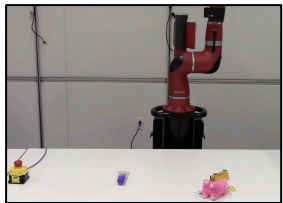
next image x_{t+1}



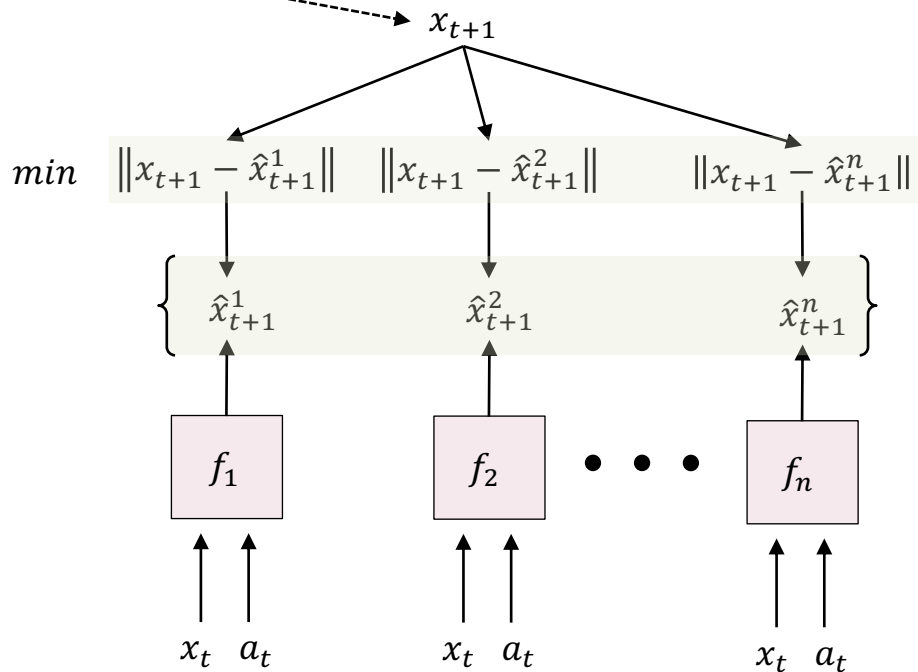
action a_t

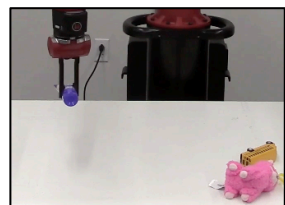


policy network
 $\pi_{\theta}(x_t)$



current image x_t

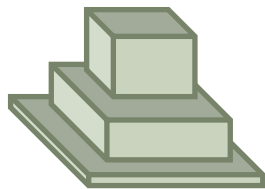




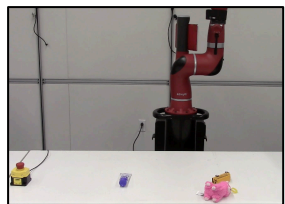
next image x_{t+1}



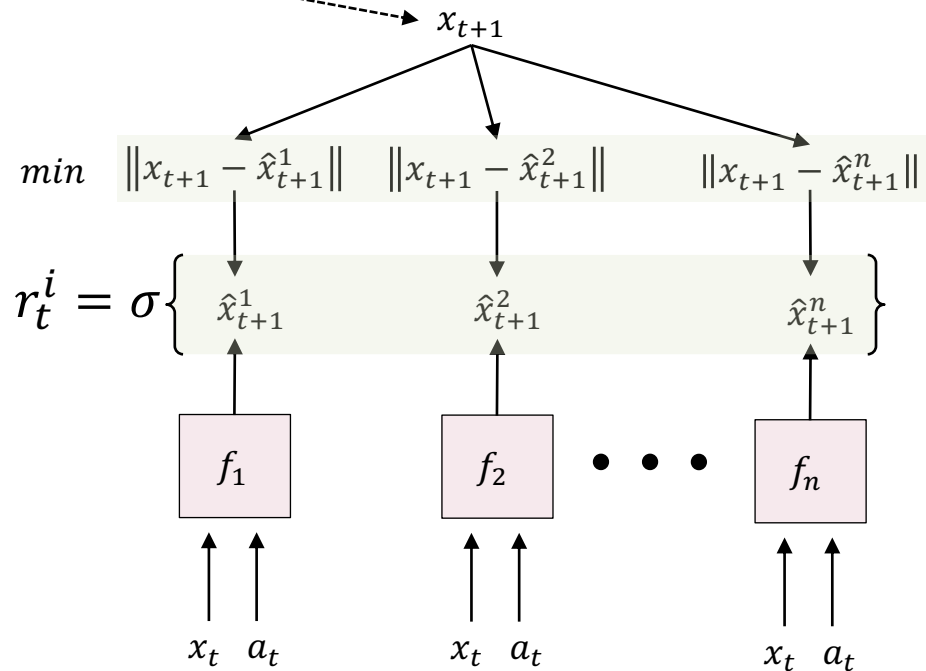
action a_t

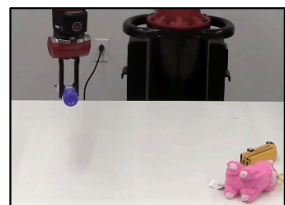


policy network
 $\pi_{\theta}(x_t)$



current image x_t

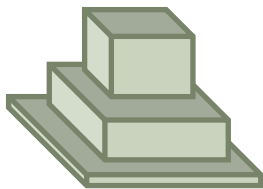




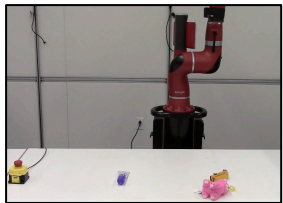
next image x_{t+1}



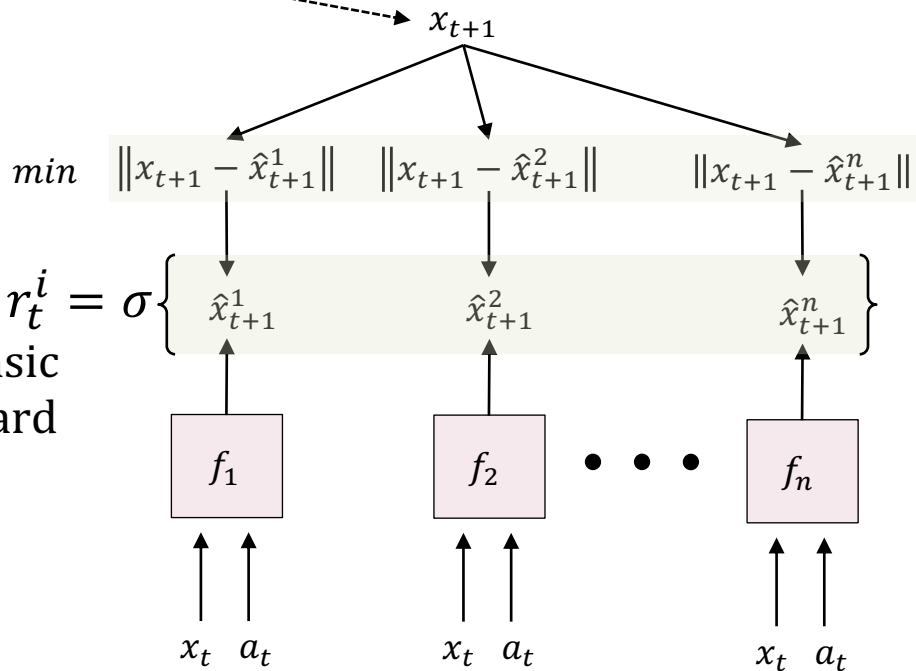
action a_t

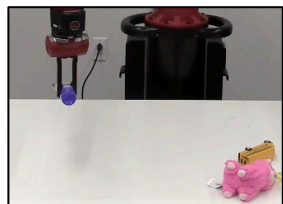


policy network
 $\pi_\theta(x_t)$



current image x_t

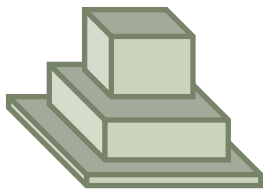




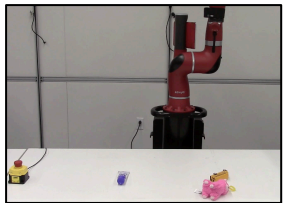
next image x_{t+1}



action a_t

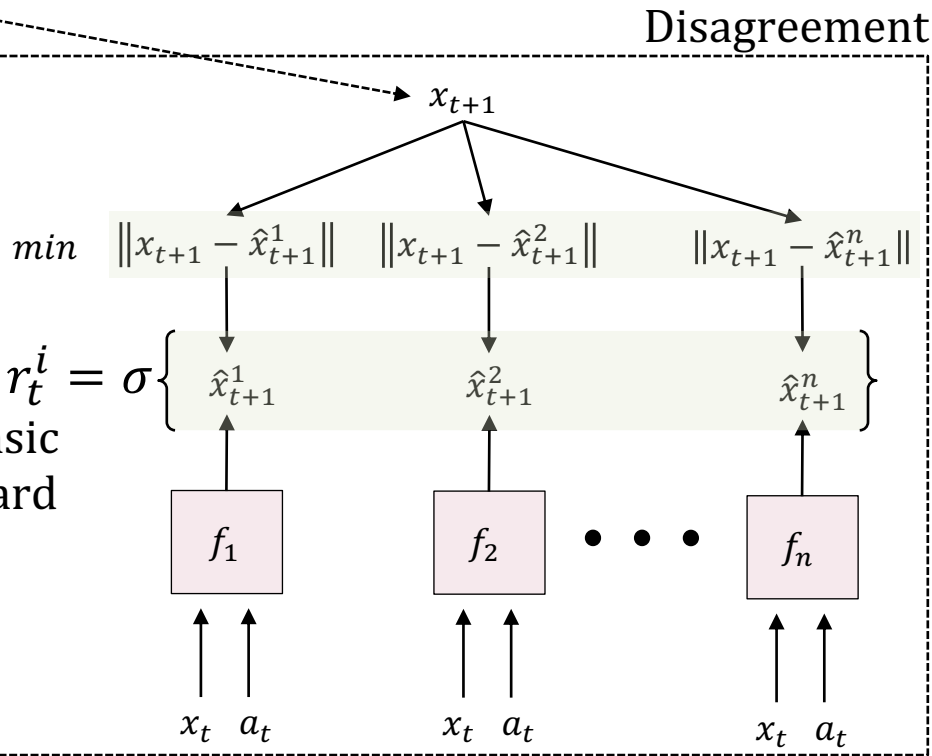


policy network
 $\pi_\theta(x_t)$



current image x_t

Intrinsic
Reward



Deterministic Environments

performs as well as state-of-the-art methods

Deterministic Environments

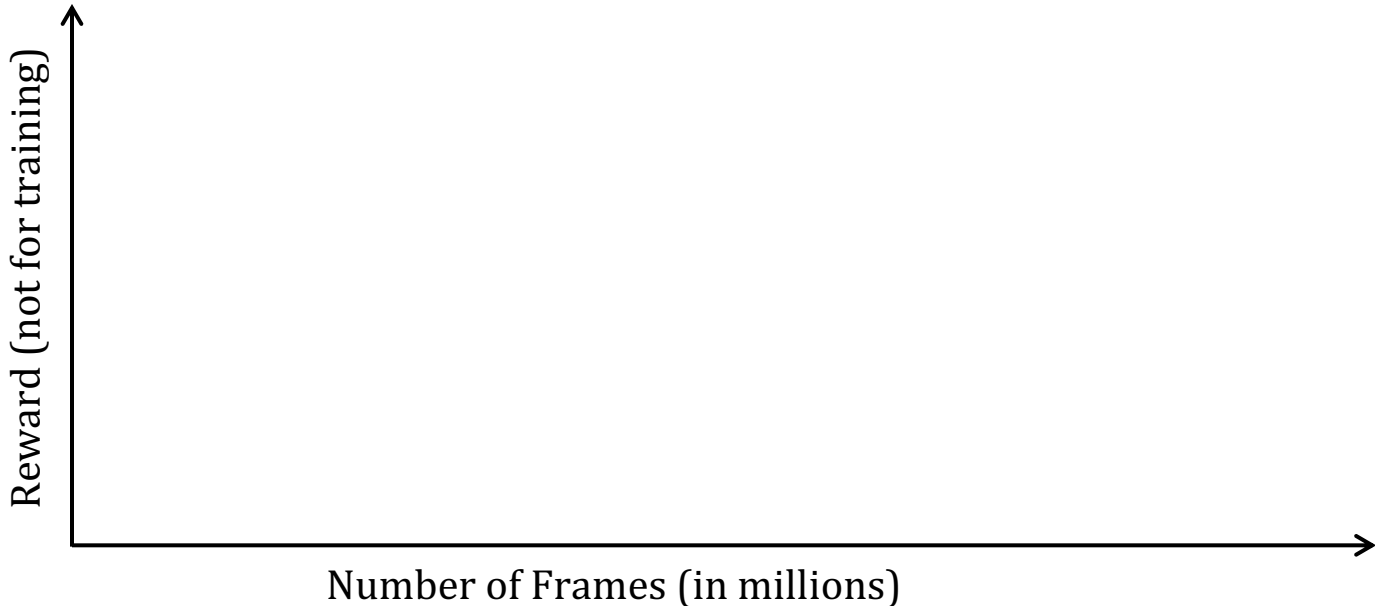
performs as well as state-of-the-art methods

Reward (not for training)



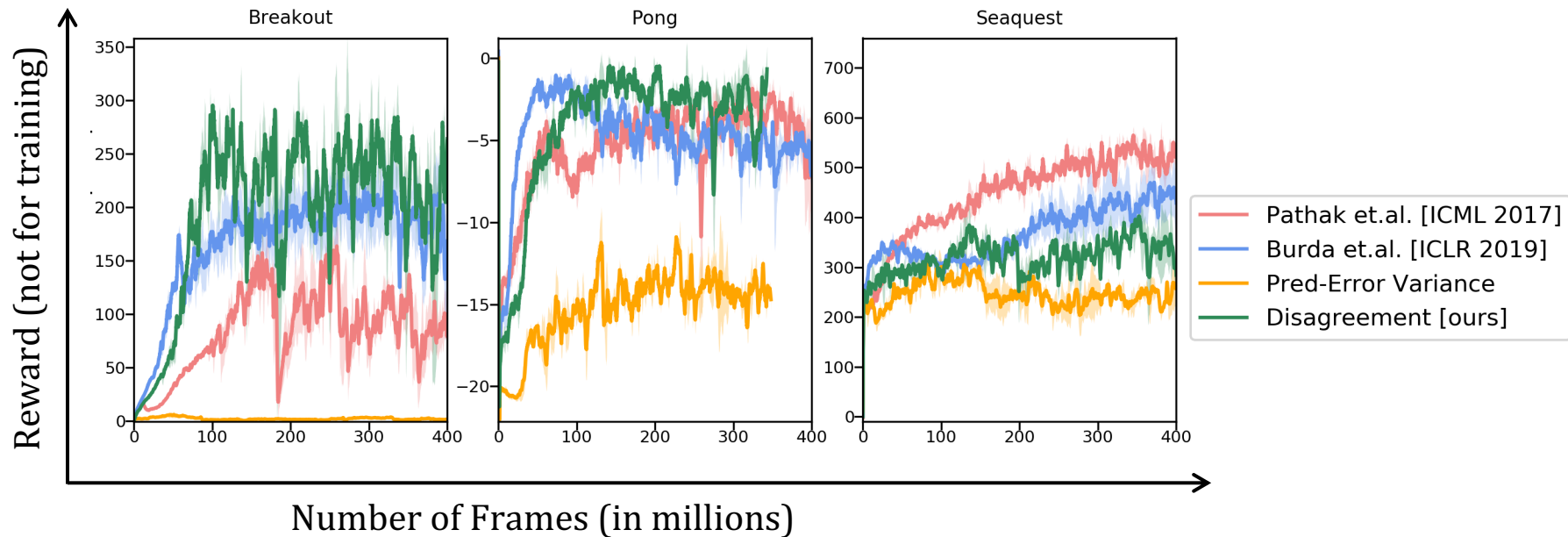
Deterministic Environments

performs as well as state-of-the-art methods



Deterministic Environments

performs as well as state-of-the-art methods



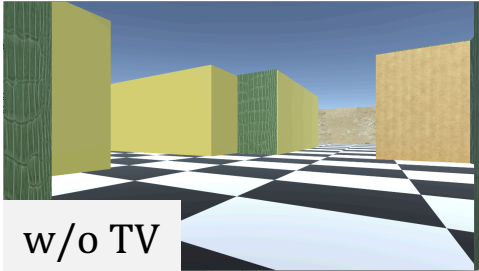
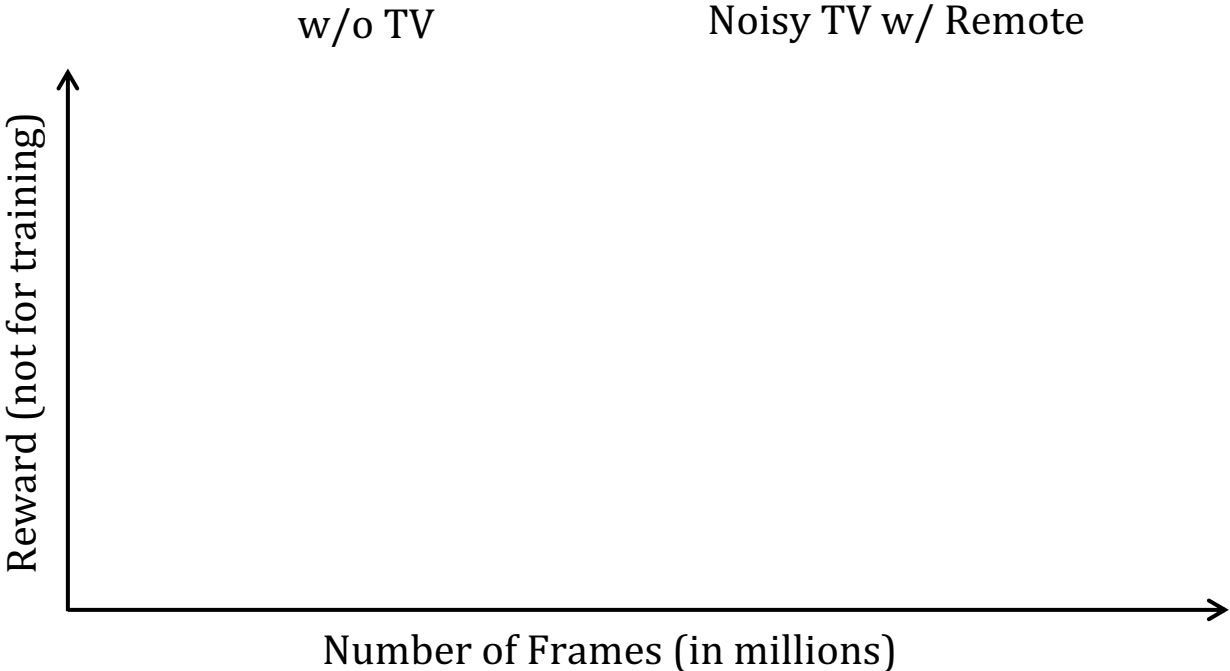
Stochastic Environments

Stochastic Environments

Every model's goes to mean \rightarrow variance drops \rightarrow unstuck

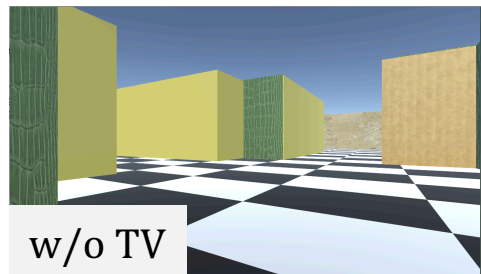
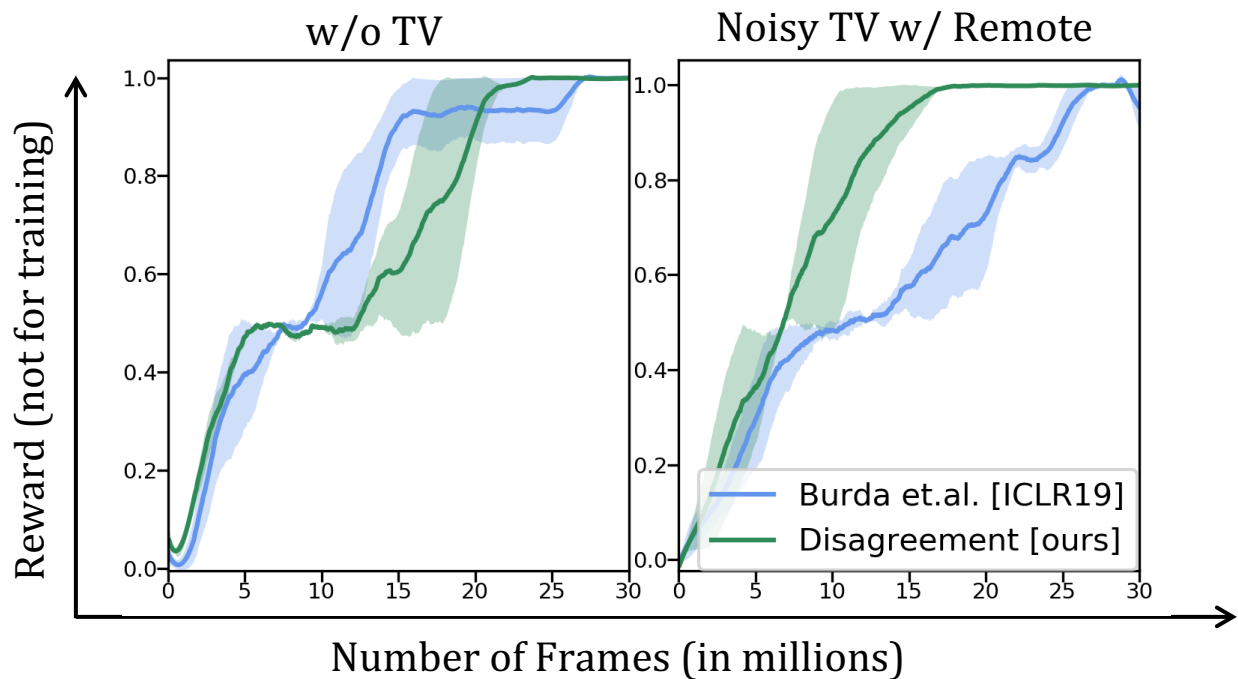
Stochastic Environments: 3D Navigation

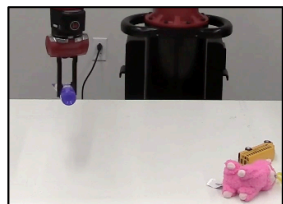
Every model's goes to mean \rightarrow variance drops \rightarrow unstuck



Stochastic Environments: 3D Navigation

Every model's goes to mean \rightarrow variance drops \rightarrow unstuck

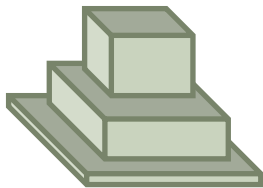




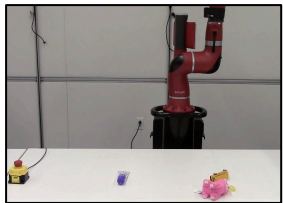
next state x_{t+1}



action a_t

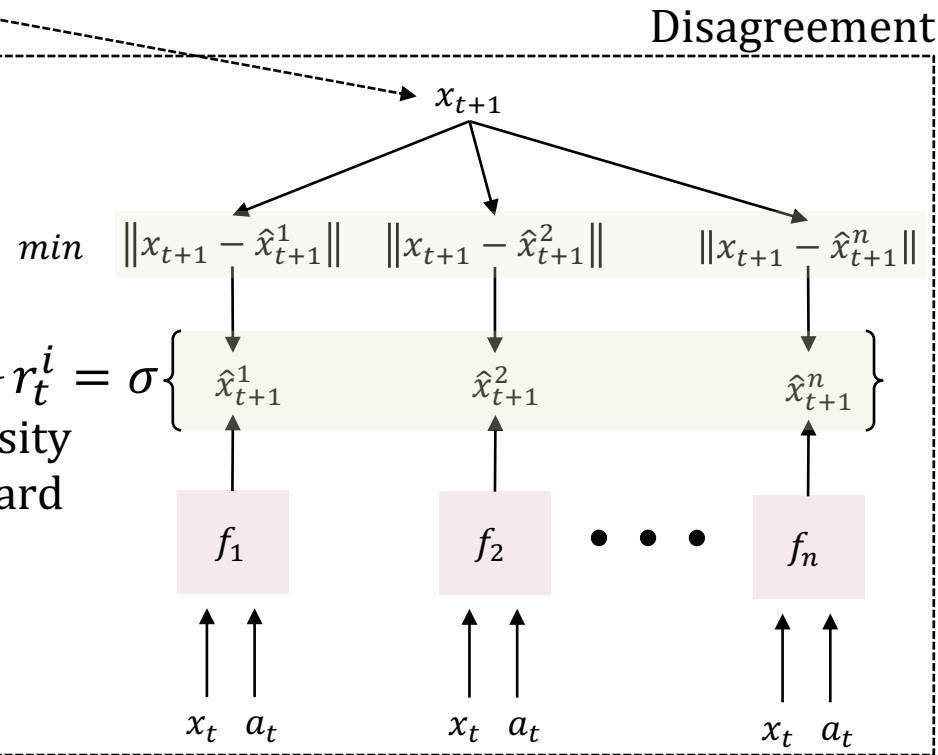


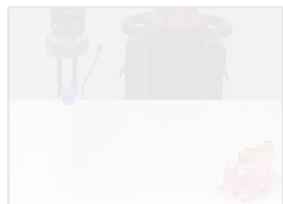
policy network
 $\pi_\theta(x_t)$



current state x_t

Curiosity
Reward





next state x_{t+1}

Disagreement

$$r_t^i \triangleq \mathbb{E}_\theta \left[\left\| f(x_t, a_t; \theta) - \mathbb{E}_\theta[f(x_t, a_t; \theta)] \right\|_2^2 \right]$$

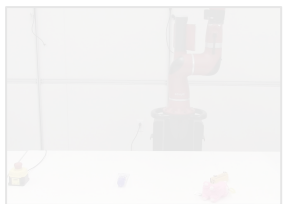
$x_{t+1} - \hat{x}_{t+1}^n$



policy network
 $\pi_\theta(x_t)$

Curiosity
Reward

$$r_t^i = \sigma \left\{ \hat{x}_{t+1}^1, \hat{x}_{t+1}^2, \dots, \hat{x}_{t+1}^n \right\}$$



current state x_t



x_t

a_t



x_t

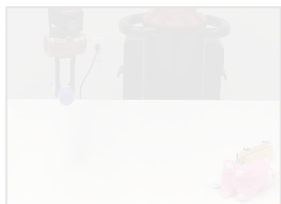
a_t

...



x_t

a_t



next state x_{t+1}

Disagreement

$$r_t^i \triangleq \mathbb{E}_\theta \left[\left\| f(x_t, a_t; \theta) - \mathbb{E}_\theta[f(x_t, a_t; \theta)] \right\|_2^2 \right]$$

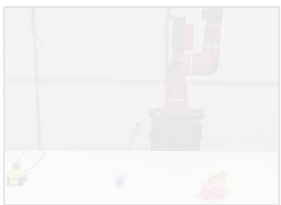
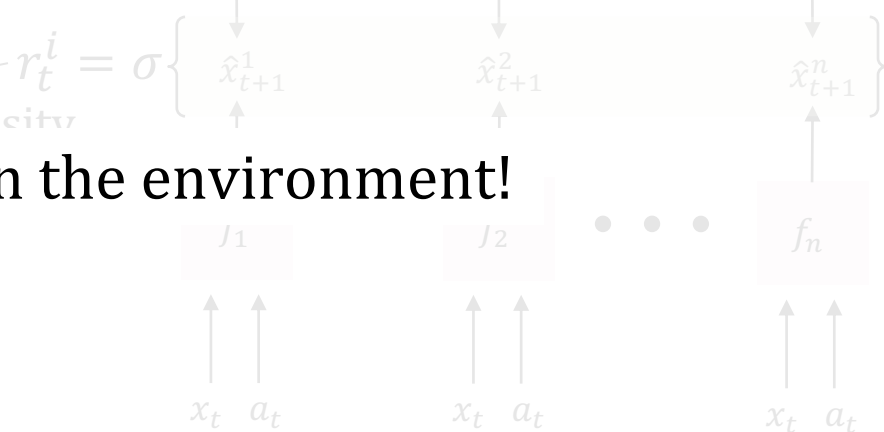
$x_{t+1} - \hat{x}_{t+1}^n$



policy network

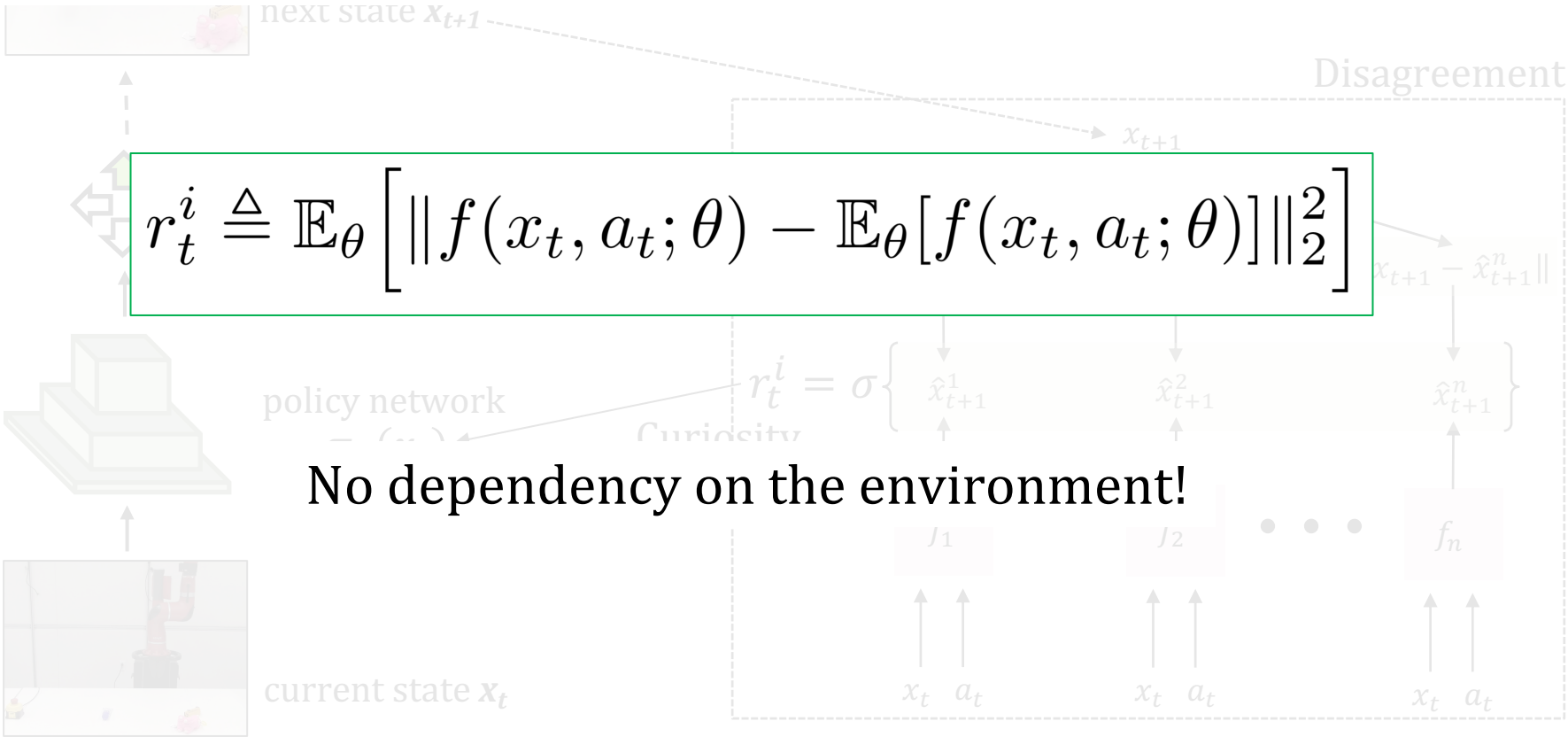
Curiosity

No dependency on the environment!



current state x_t

Differentiable Exploration



No dependency on the environment!

Differentiable Exploration

Differentiable Exploration

Model
Optimization

$$\min_{\theta_1, \dots, \theta_k} \sum_{i=1}^k \|f_{\theta_i}(x_t, \pi(x_t; \theta_P)) - x_{t+1}\|_2$$

Differentiable Exploration

Model
Optimization

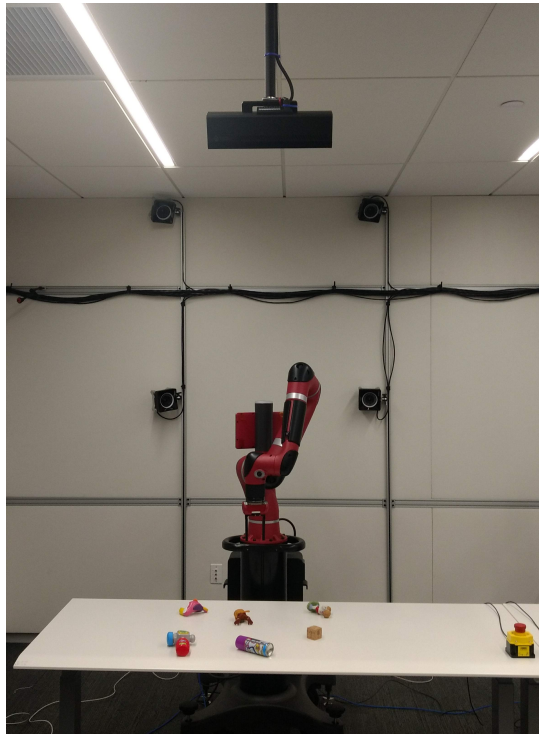
$$\min_{\theta_1, \dots, \theta_k} \sum_{i=1}^k \|f_{\theta_i}(x_t, \pi(x_t; \theta_P)) - x_{t+1}\|_2$$

Policy
Optimization

$$\max_{\theta_P} \sum_{i=1}^k \left\| f_{\theta_i}(x_t, \pi(x_t; \theta_P)) - \left(\frac{1}{k}\right) \sum_{j=1}^k f_{\theta_j}(x_t, \pi(x_t; \theta_P)) \right\|_2$$

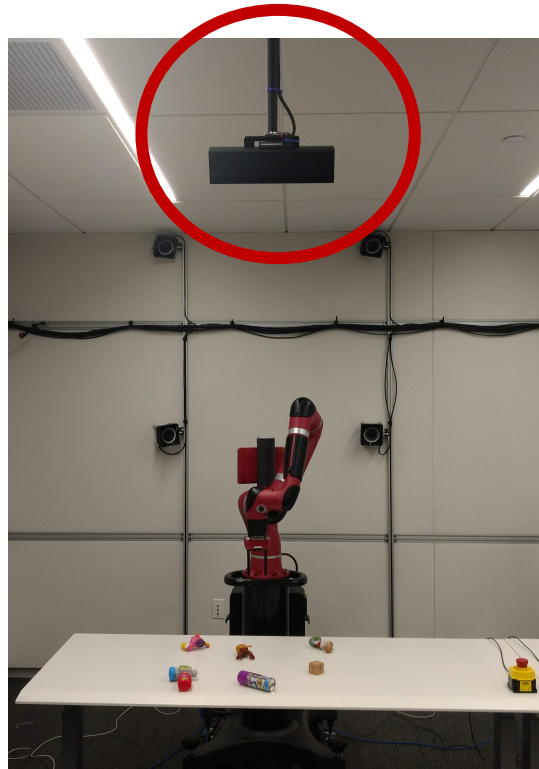
Differentiable Exploration

Differentiable Exploration

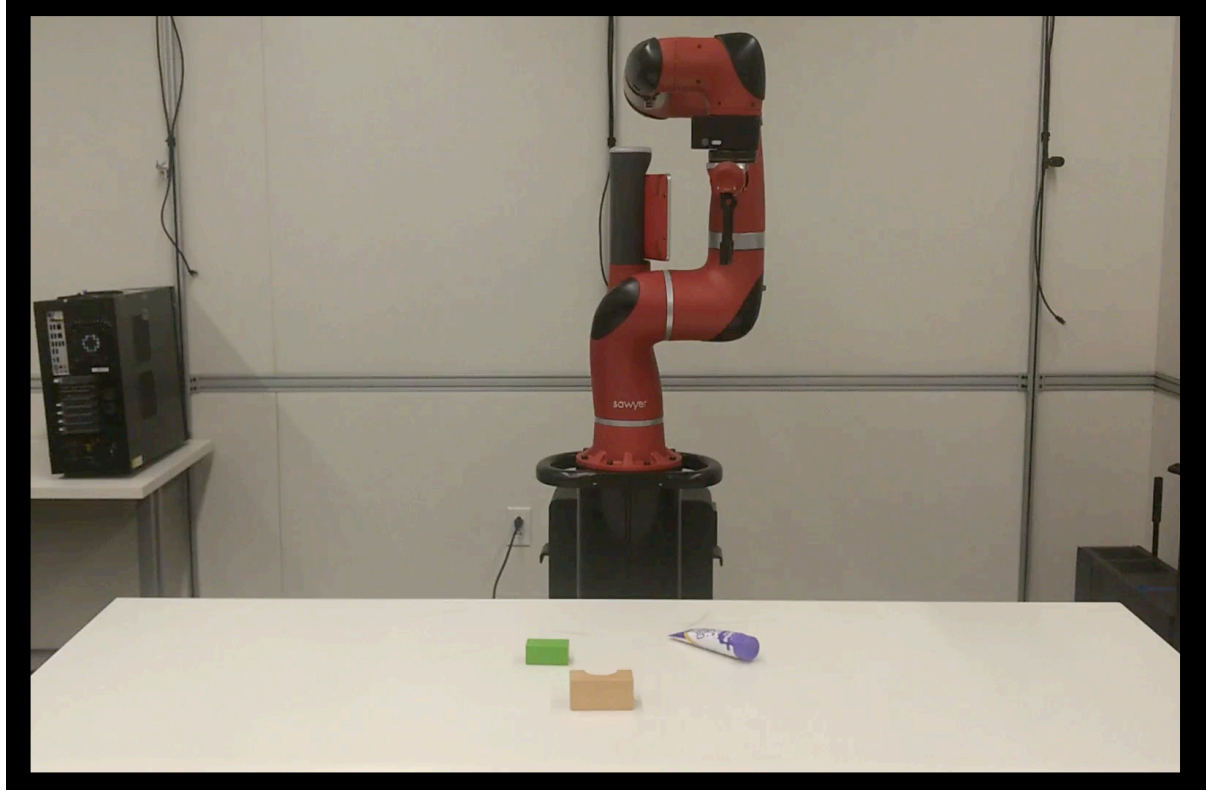


Pathak*, **Gandhi***, **Gupta**. "Self-Supervised Exploration via Disagreement". ICML, 2019.

Differentiable Exploration



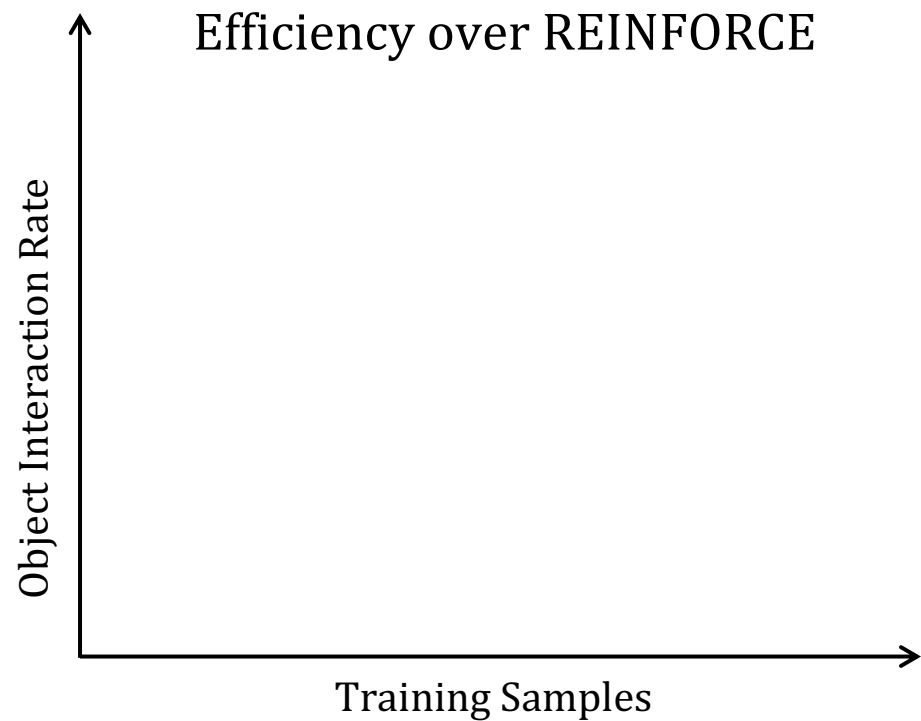
Differentiable Exploration



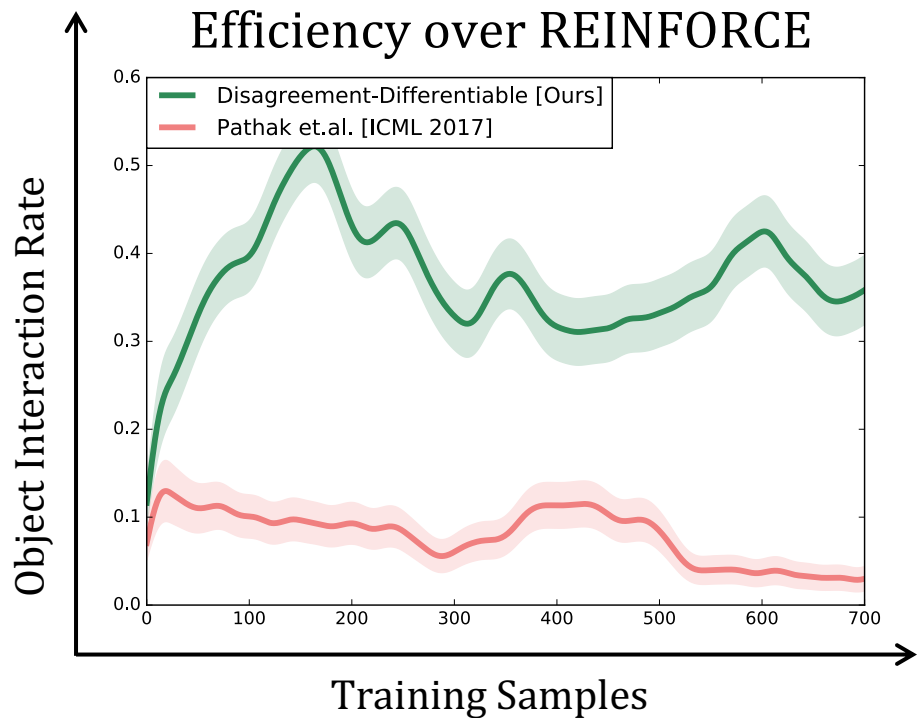
Position Control:

1. Position
2. Direction
3. Gripper Angle
4. Gripper Distance

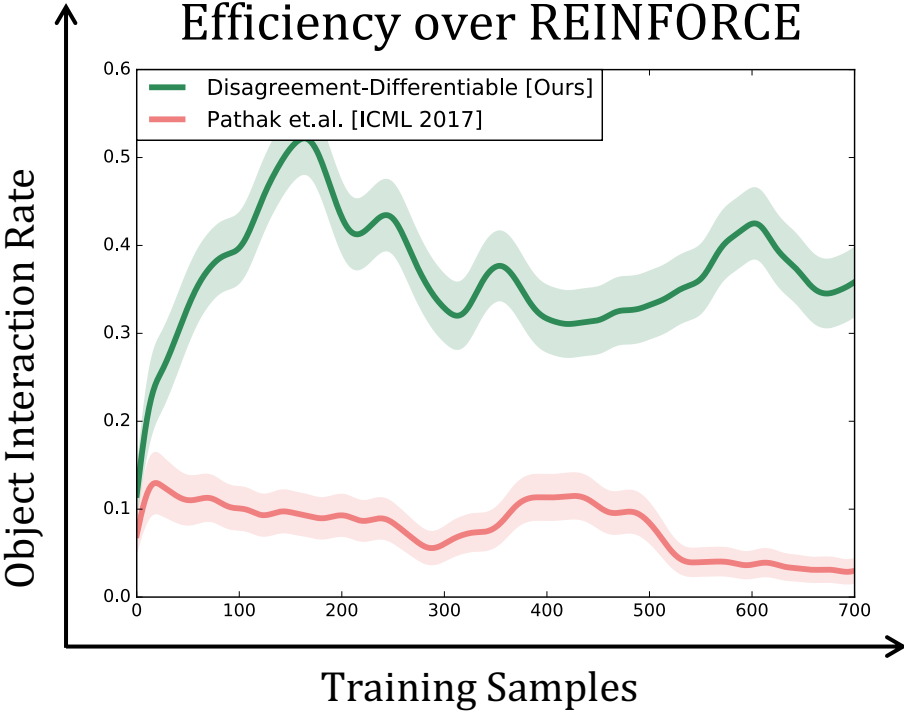
Differentiable Exploration



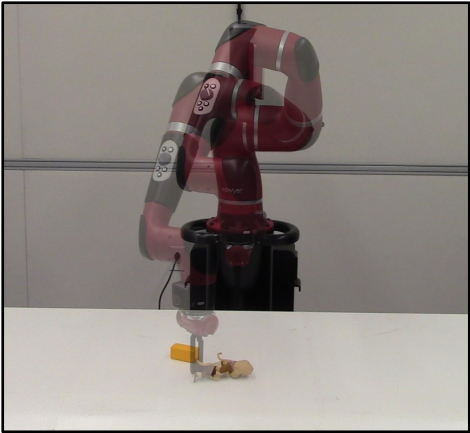
Differentiable Exploration



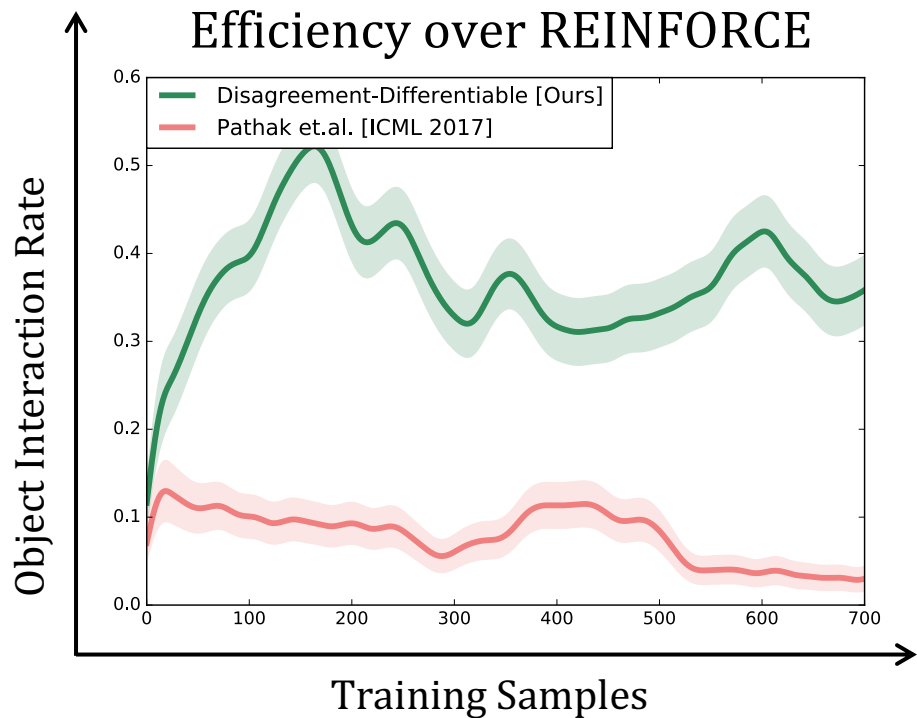
Differentiable Exploration



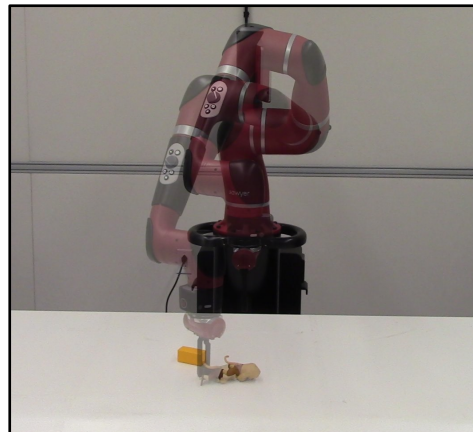
Pushing
skill



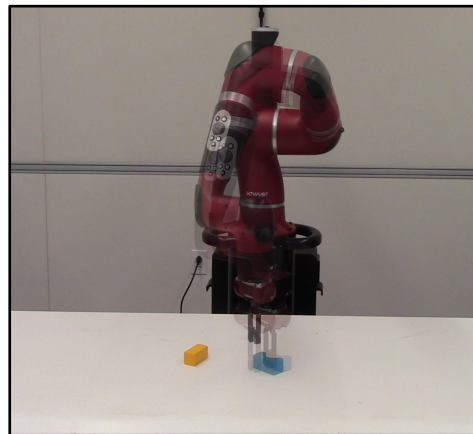
Differentiable Exploration



Pushing
skill



Picking
skill



Summary: Exploration via Disagreement

Summary: Exploration via Disagreement

- Similar to state-of-the-art in deterministic envs
(Atari games)

Summary: Exploration via Disagreement

- Similar to state-of-the-art in deterministic envs
(Atari games)
- Does not get stuck in stochastic scenarios
(Stochastic Atari; Unity-TV)

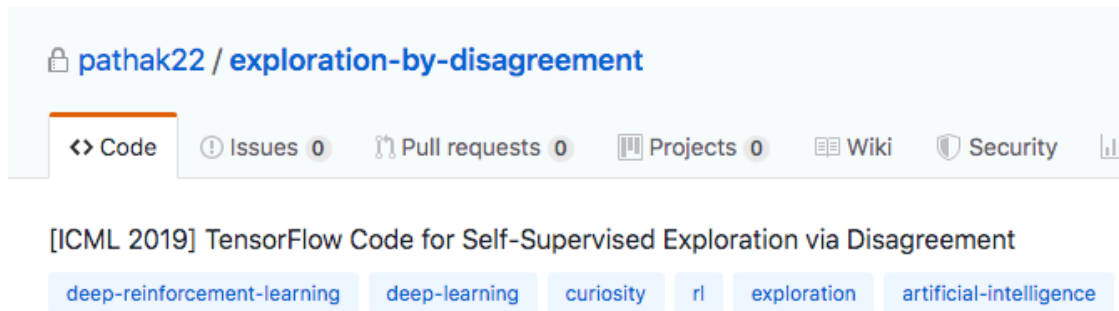
Summary: Exploration via Disagreement

- Similar to state-of-the-art in deterministic envs
(Atari games)
- Does not get stuck in stochastic scenarios
(Stochastic Atari; Unity-TV)
- Differentiable reformulation for real robots
(Sawyer Robot)

Code Available



<https://pathak22.github.io/exploration-by-disagreement/>



Poster # 39

(today)



Thank you!