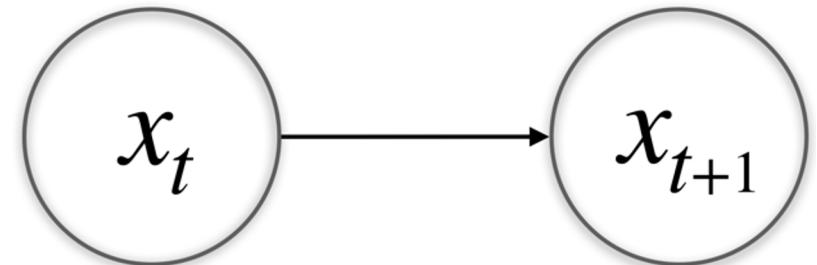
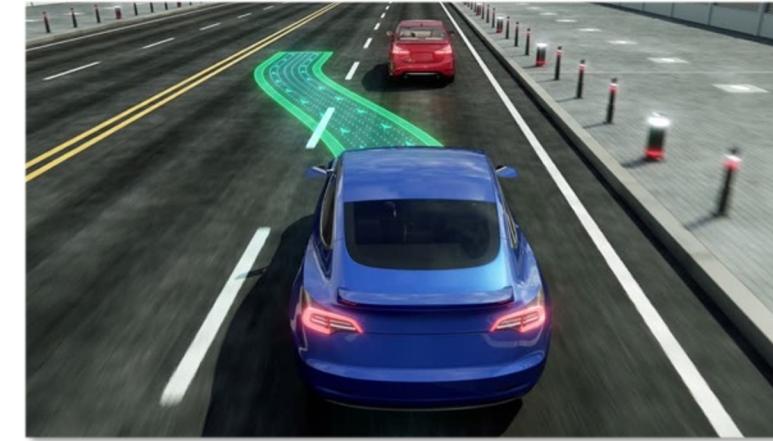
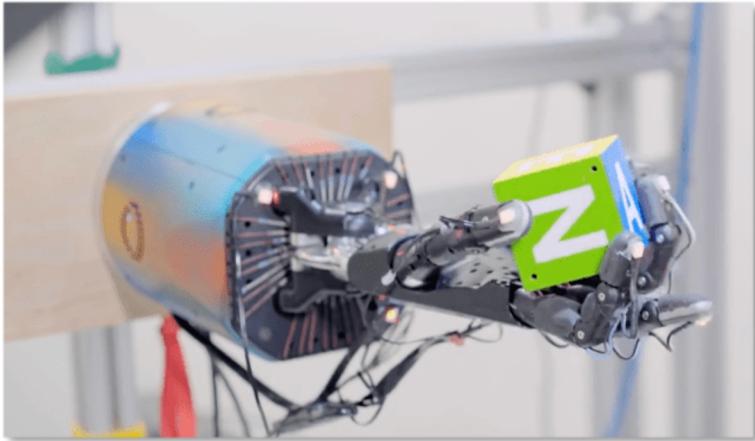


# Visual Model-Based Reinforcement Learning

Oleh Rybkin  
[oleh@seas.upenn.edu](mailto:oleh@seas.upenn.edu)



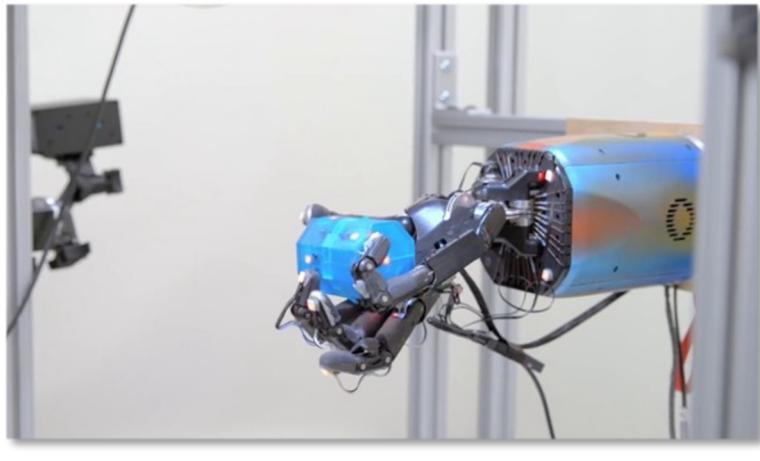
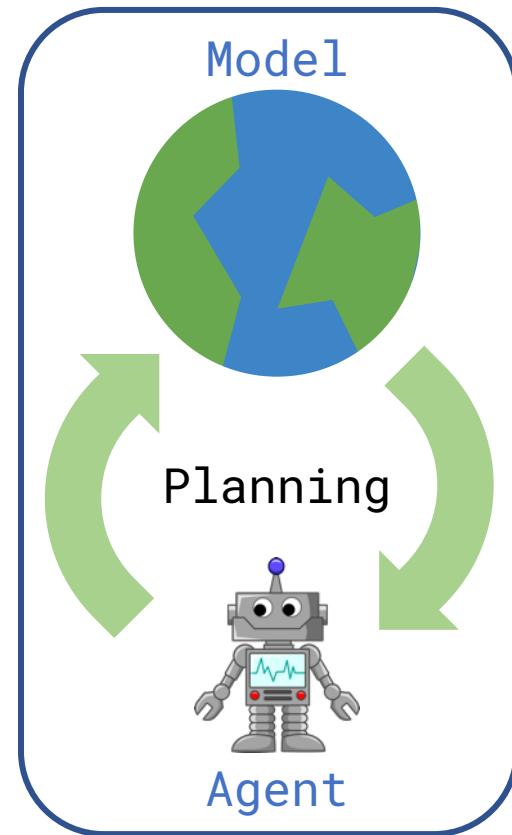
# How do we build intelligent machines?



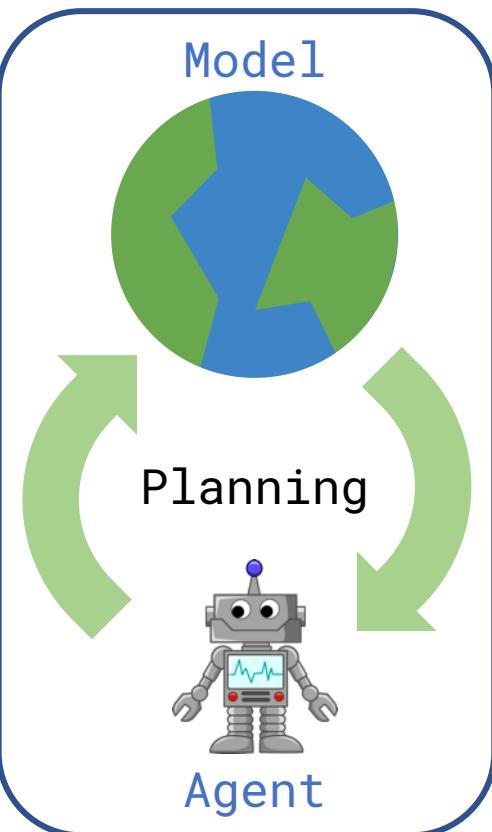
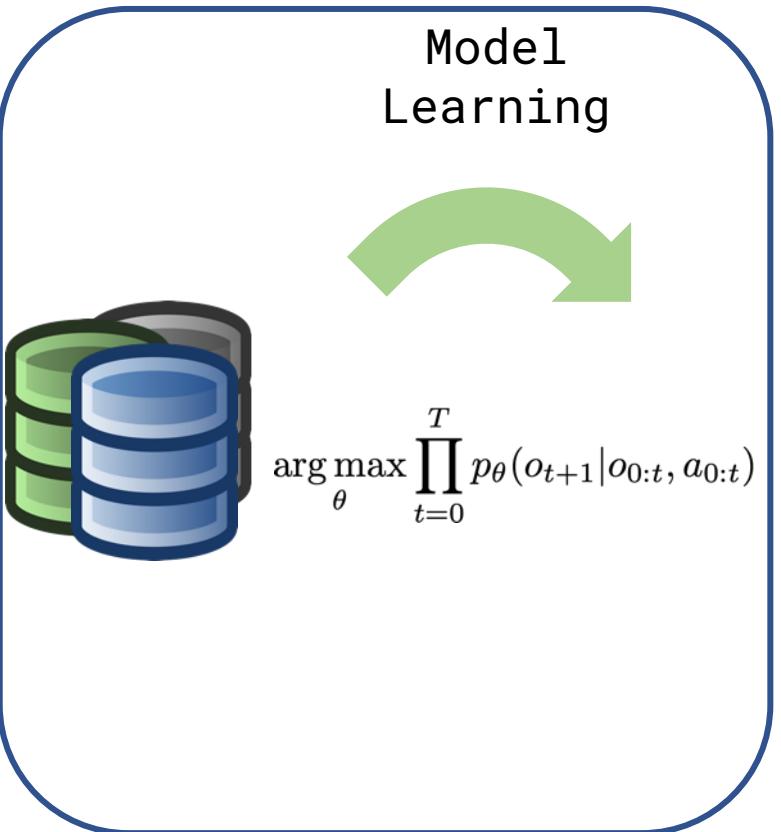
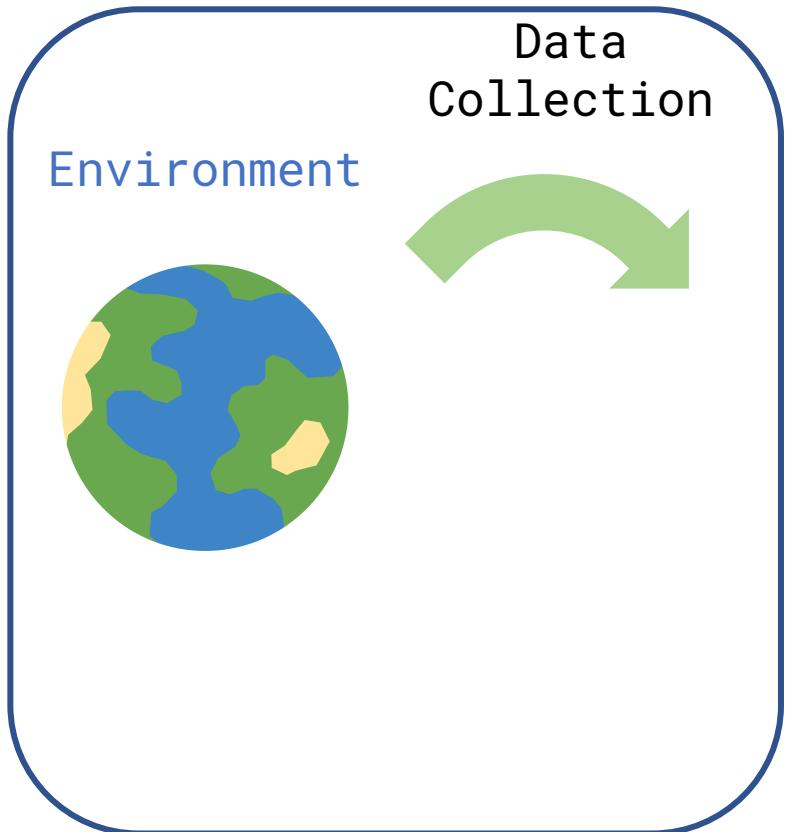
VS.



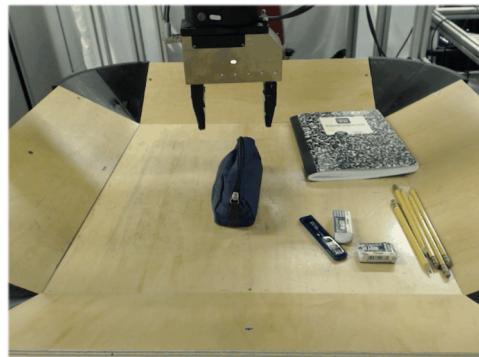
$$p_{\theta}(o_{t+1}|o_{0:t}, a_{0:t})$$



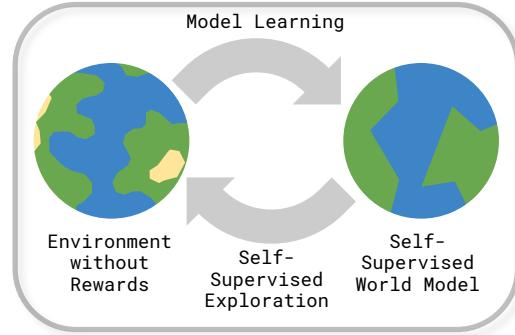
$$\arg \max_{a_{0:T}} \mathbb{E}_{\prod_{t=0}^T p_{\theta}(o_{t+1}|o_{0:t}, a_{0:t})} \left[ \sum_t r(o_t) \right]$$



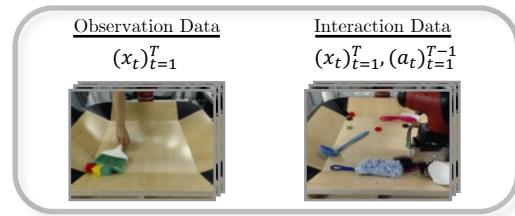
Prediction



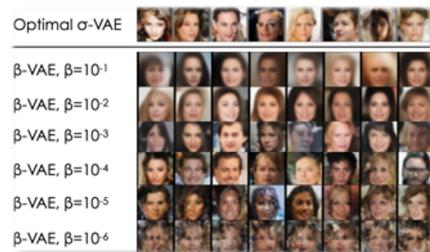
Control



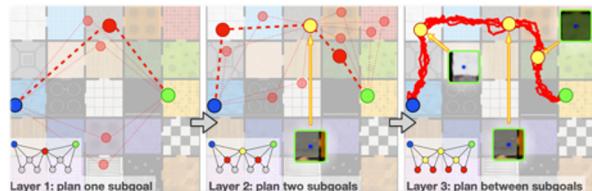
## Self-supervised data collection



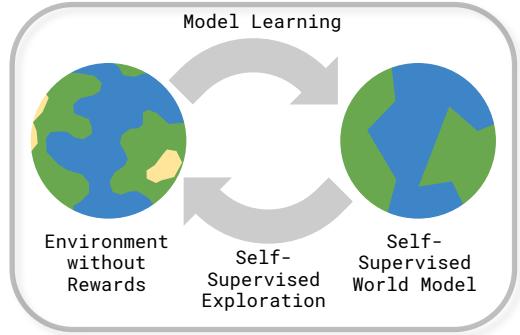
## Leveraging human-collected data



## Probabilistic latent variable models



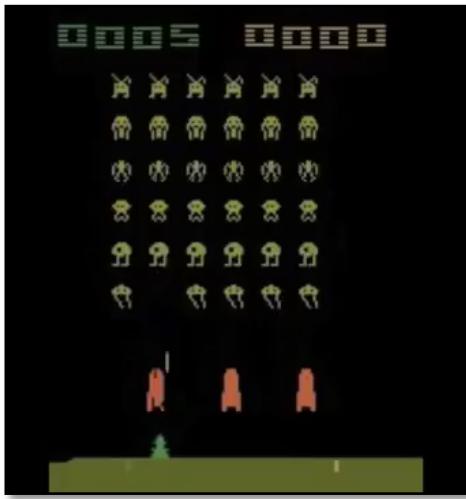
## Long-horizon hierarchical planning



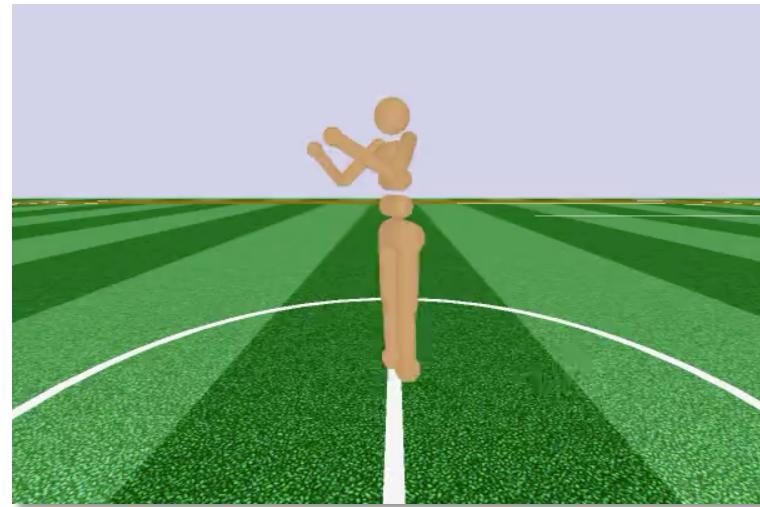
## Self-supervised data collection

Can we learn general and reusable knowledge?

# Reinforcement Learning: task-specific, difficult to generalize



[ Mnih *et al.*, Nature 2015 ]



[ Schulman *et al.*, 2015, 2017 ]



[ Silver *et al.*, Nature 2016 ]



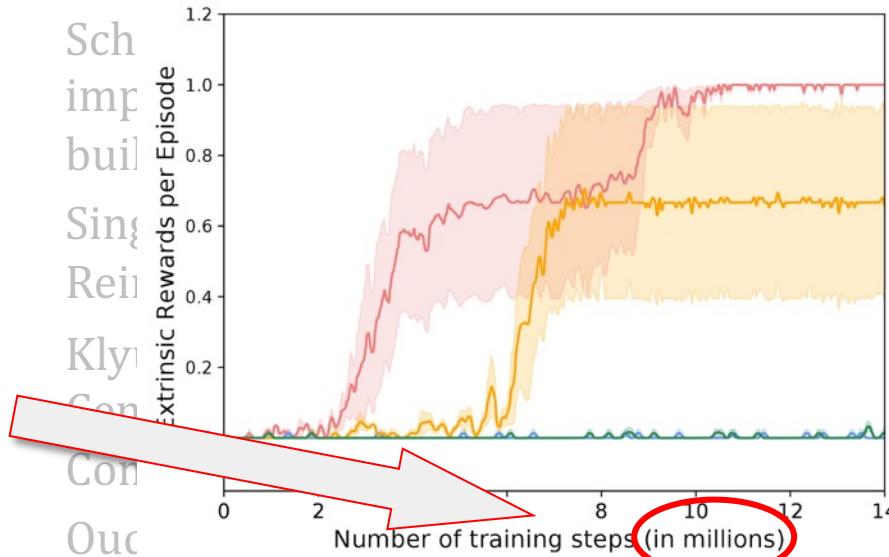
[ Kalashnikov *et al.*, CoRL 2018 ]

# Self-Supervised Reinforcement Learning

- Schmidhuber, “A possibility for implementing curiosity and boredom in model building neural controllers”. SAB, 1991.
- Singh *et al.*, “Intrinsically Motivated Reinforcement Learning”. NeurIPS, 2004.
- Klyubin et al., “Empowerment: A Universal Agent-Centric Measure of Control”. Evolutionary Computation, 2005.
- Oudeyer & Kaplan, “What is intrinsic motivation? A typology of computational approaches”. Frontiers in Neurorobotics, 2009.
- Sun *et al.*, “Planning to Be Surprised: Optimal Bayesian Exploration in Dynamic Environments”. AGI, 2011
- Mohamed & Rezende, “Variational information maximisation for intrinsically motivated reinforcement learning”. NeurIPS, 2015.
- Gregor *et al.*, “Variational intrinsic control”. ICLR Workshop, 2017.
- Pathak *et al.*, “Curiosity-driven Exploration by Self-supervised Exploration”. ICML 2017
- Pathak *et al.*, “Self-Supervised Exploration via Disagreement”. ICML, 2019.
- Hansen *et al.*, “Fast Task Inference with Variational Intrinsic Successor Features”. ICLR 2020.

# Self-Supervised Reinforcement Learning

- Schmidhuber, J., 1991. “Learning to learn by self-improving building blocks”. *Complex Systems*, 5(1), pp. 183–200.
- Singh, S., 1996. “Reinforcement learning: A survey”. *Journal of Artificial Intelligence Research*, 5, pp. 397–418.
- Klymko, C., 1997. “Curiosity-driven exploration”. In: *Proceedings of the 14th National Conference on Artificial Intelligence*. San Jose, CA, USA: AAAI Press, pp. 103–108.
- Condit, C., 1997. “Curiosity-driven exploration”. In: *Proceedings of the 14th National Conference on Artificial Intelligence*. San Jose, CA, USA: AAAI Press, pp. 103–108.
- Oudeyer, Y., 2009. “Intrinsic motivation systems for autonomous mental development”. *Frontiers in Psychology*, 1(1), p. 44.



model

Agent

ivation? A

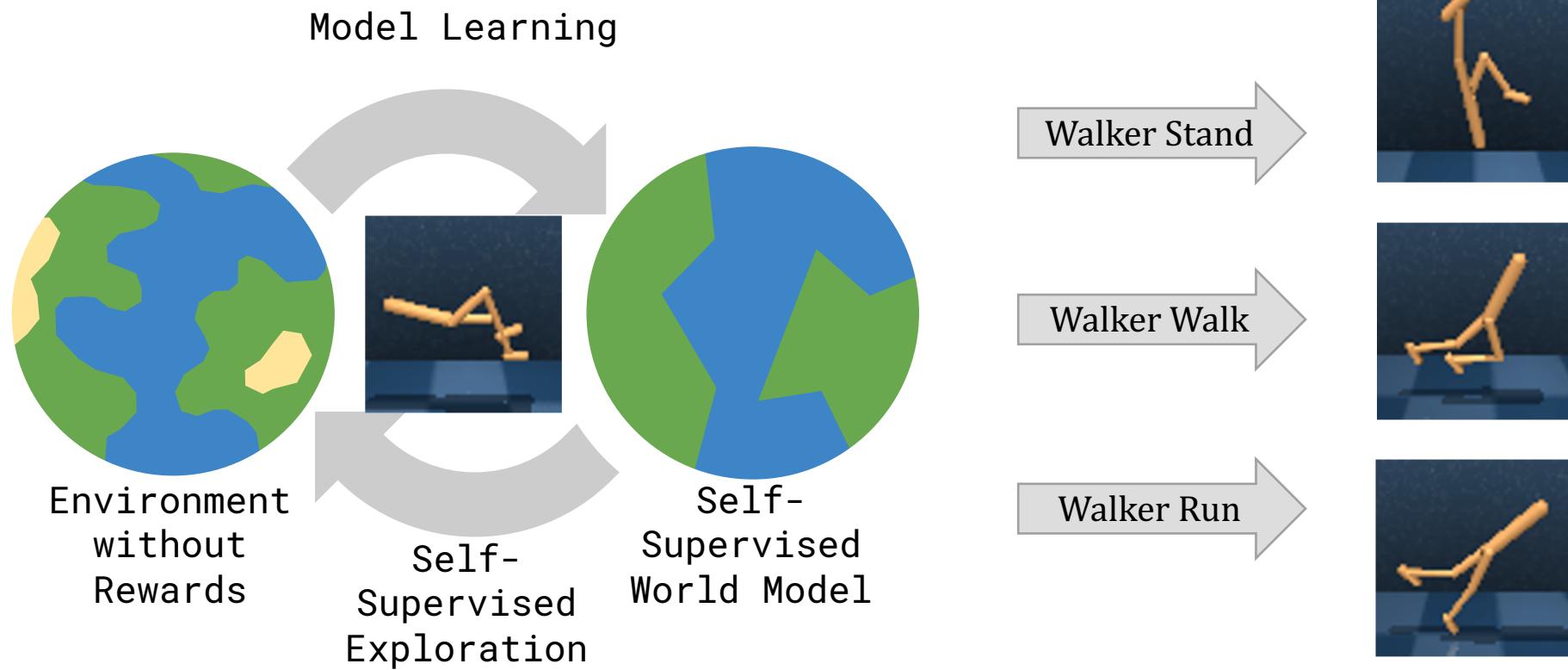
- Gregor et al., “Variational intrinsic control”, *ICML workshop*, 2017.
- Pathak et al., “Curiosity-driven Exploration via Disagreement”, *ICML*, 2017.
- Pathak et al., “Curiosity-driven Exploration via Disagreement”, *ICML*, 2017.
- Han, J., 2020. “Learning to explore via inference with Variational Intrinsic Rewards”. *ICLR 2020*.

Inefficient Adaptation  
[millions of samples]

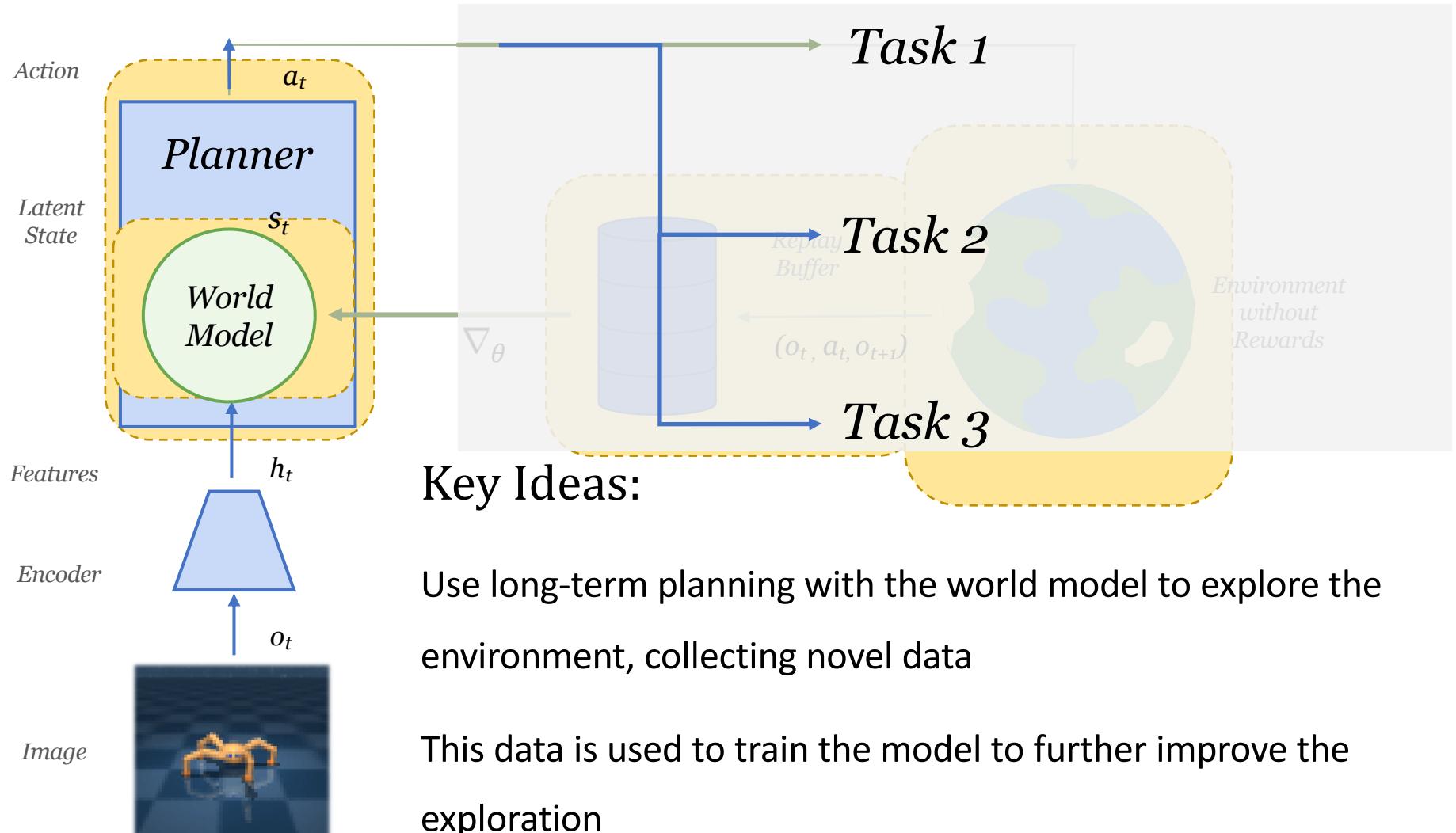


MARIO WORLD TIME

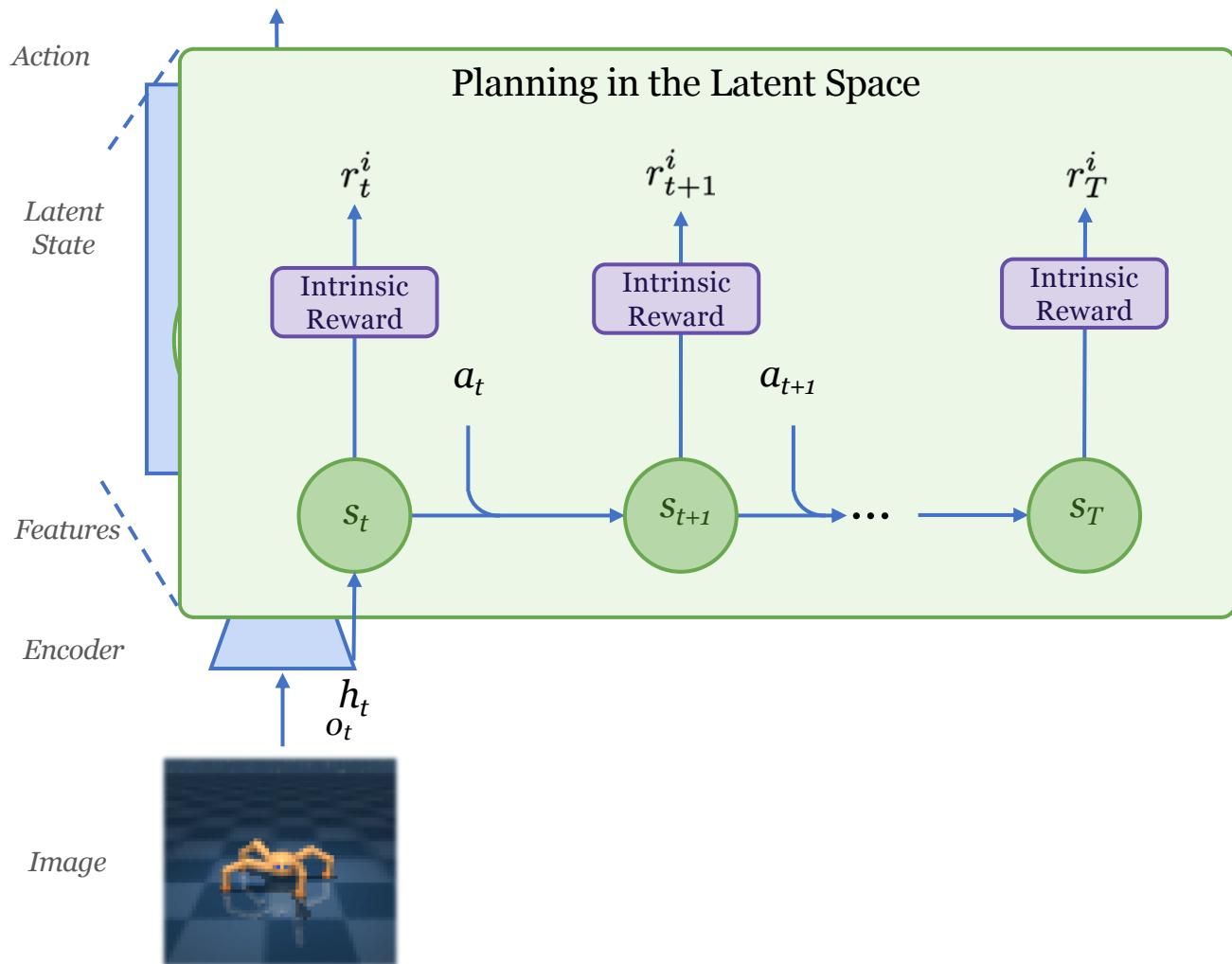
# Model-building exploration



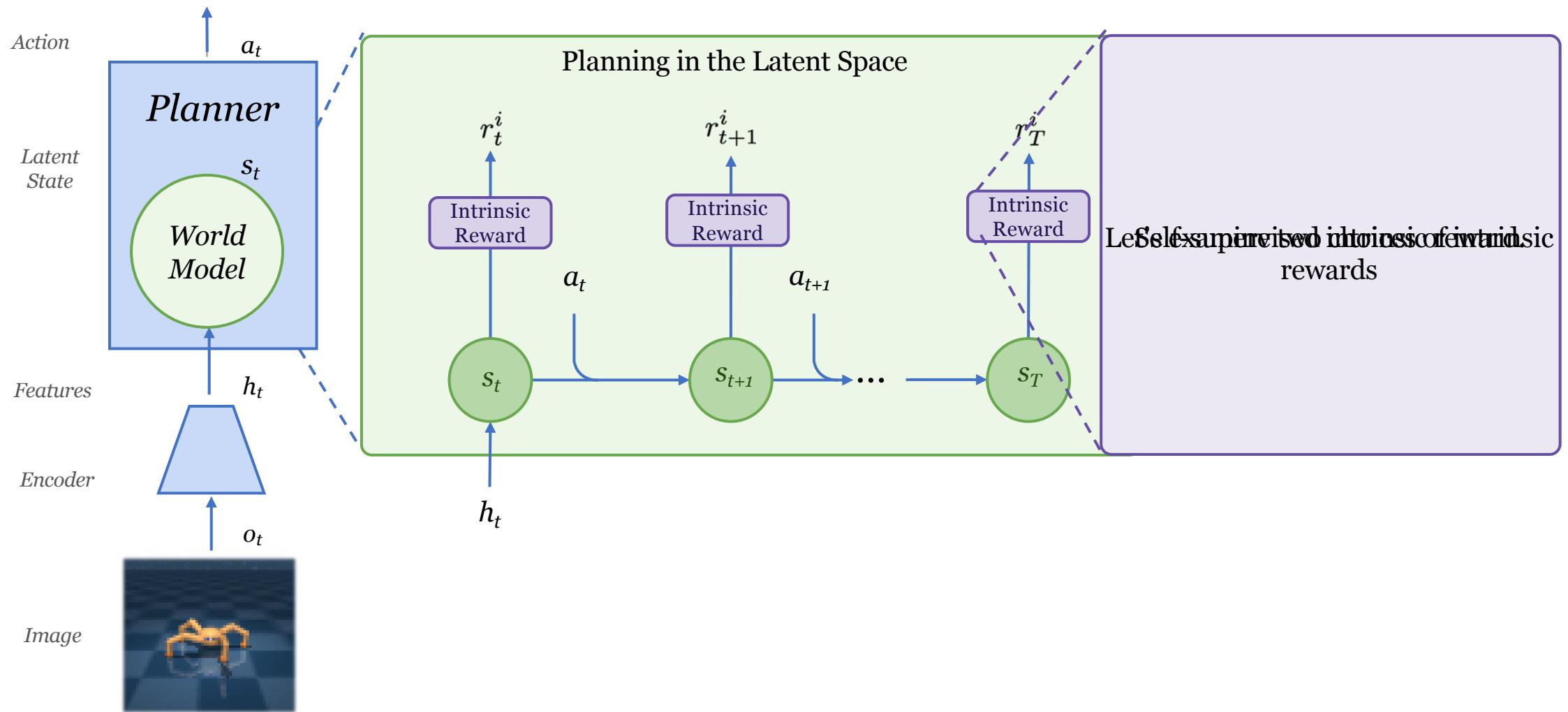
# Planning to Explore



# Planning to Explore

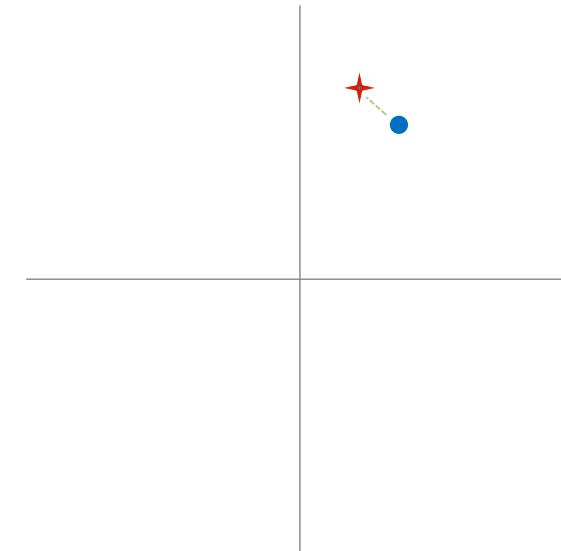


# Planning to Explore

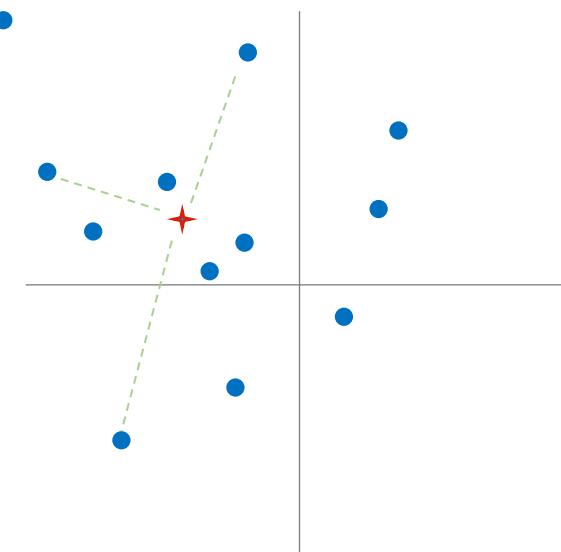


## Model Error

Deterministic  
Environment

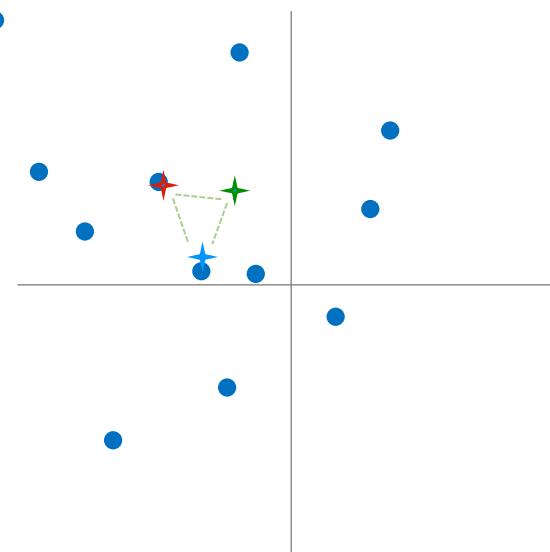
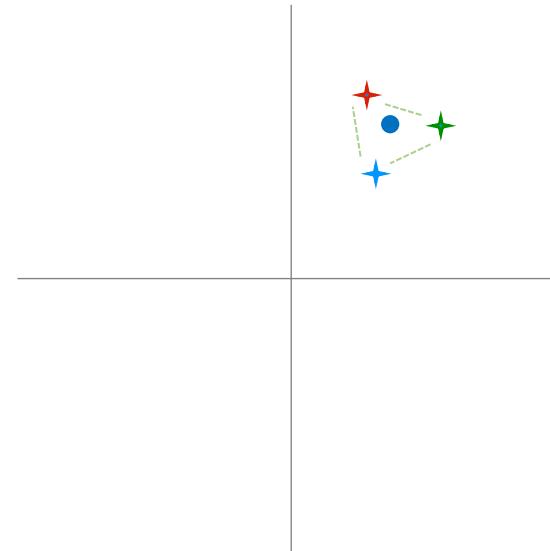


Stochastic  
Environment

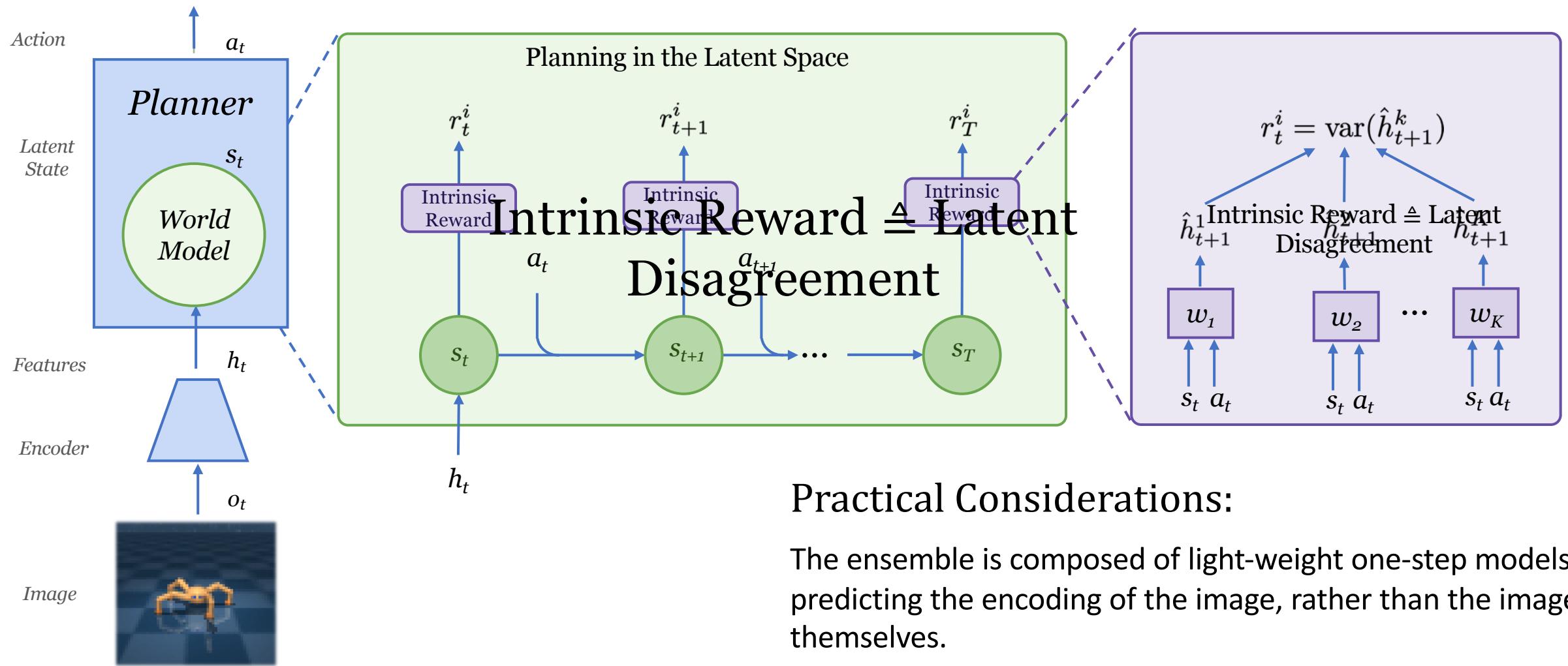


- - Data points
- ★ - Prediction of the model
- ★ - Prediction of the model 2
- ◆ - Prediction of the model 3
- Model Error

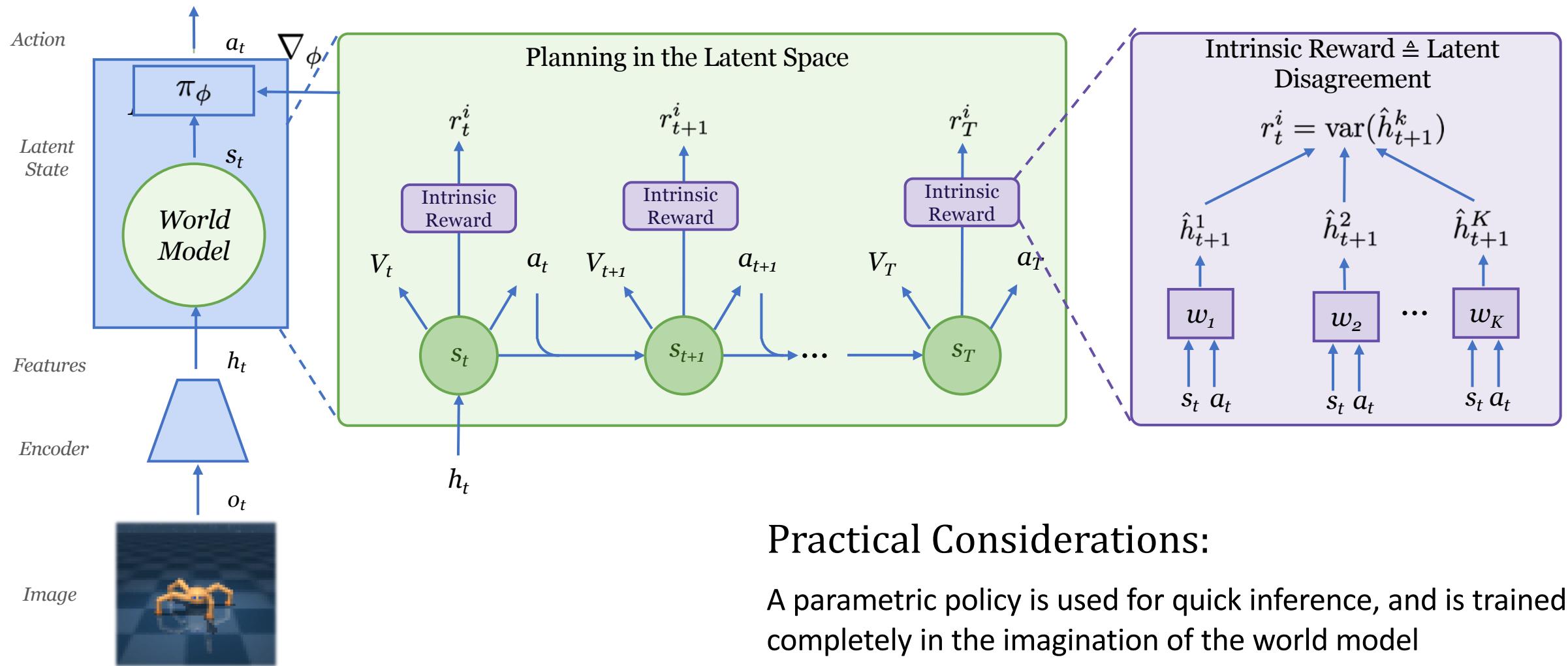
## Model Disagreement



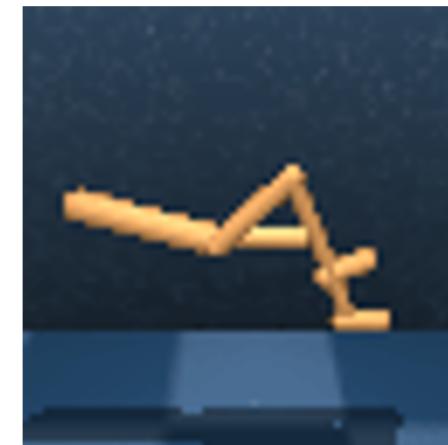
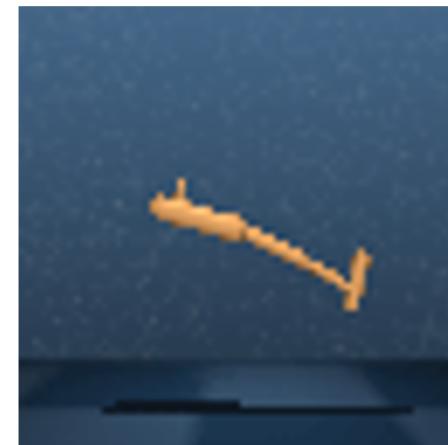
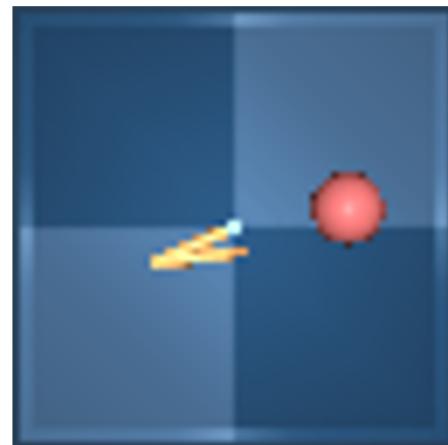
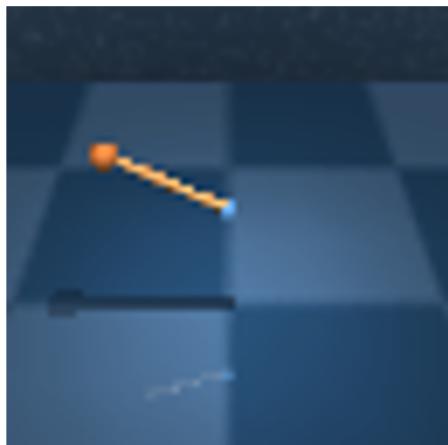
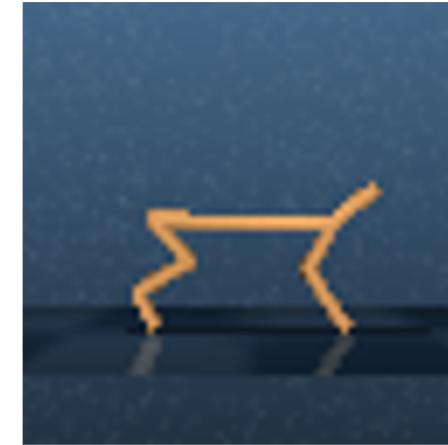
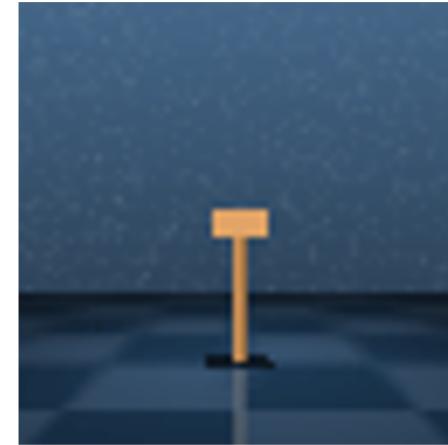
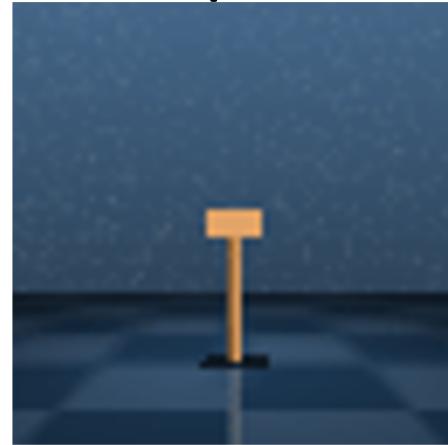
# Planning to Explore



# Planning to Explore

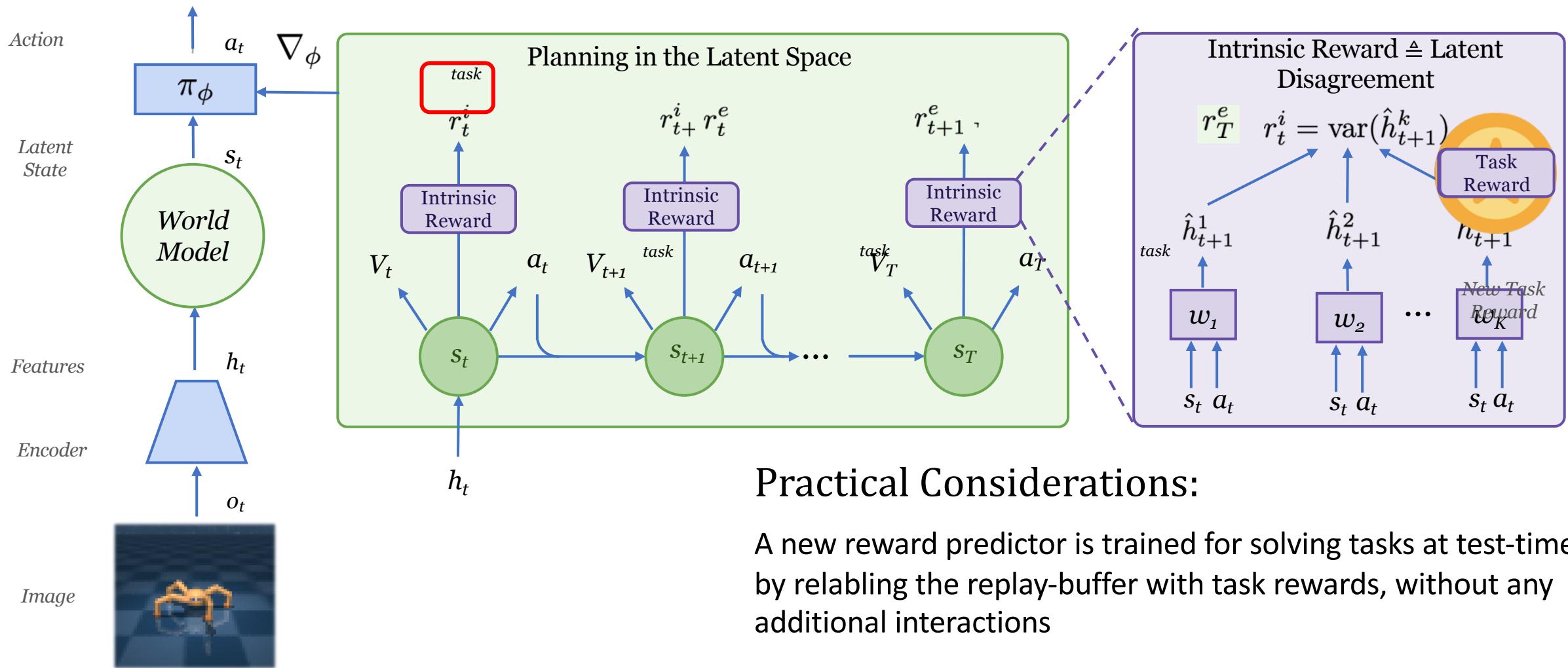


# Self-Supervised Exploration Results



How do we go from exploration to solving tasks?

# Exploration → Tasks



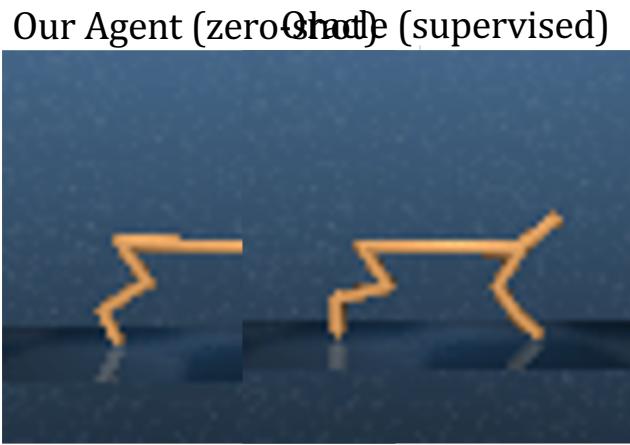
# Experiments Outline

1. Solving a new-task in zero-shot
2. What if we add 20 supervised episodes?
3. Multi-task performance

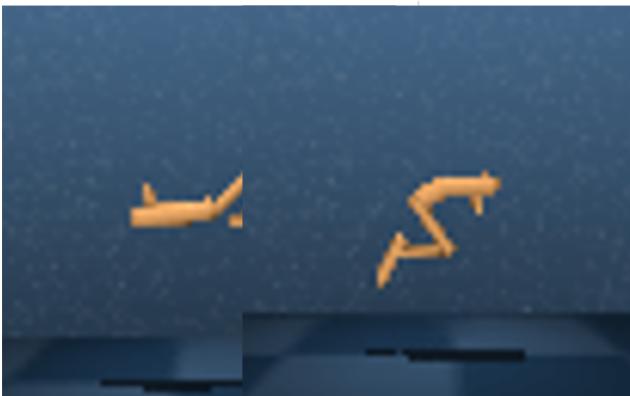
# Experiments Outline

1. Solving a new-task in zero-shot
2. What if we add 20 supervised episodes?
3. Multi-task performance

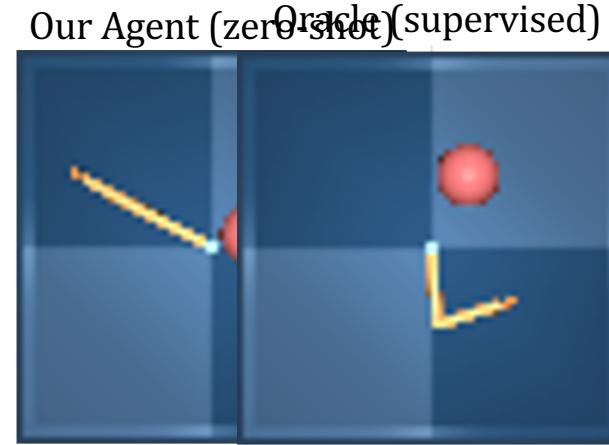
# Zero-Shot Reinforcement Learning



Cheetah Run



Hopper Hop

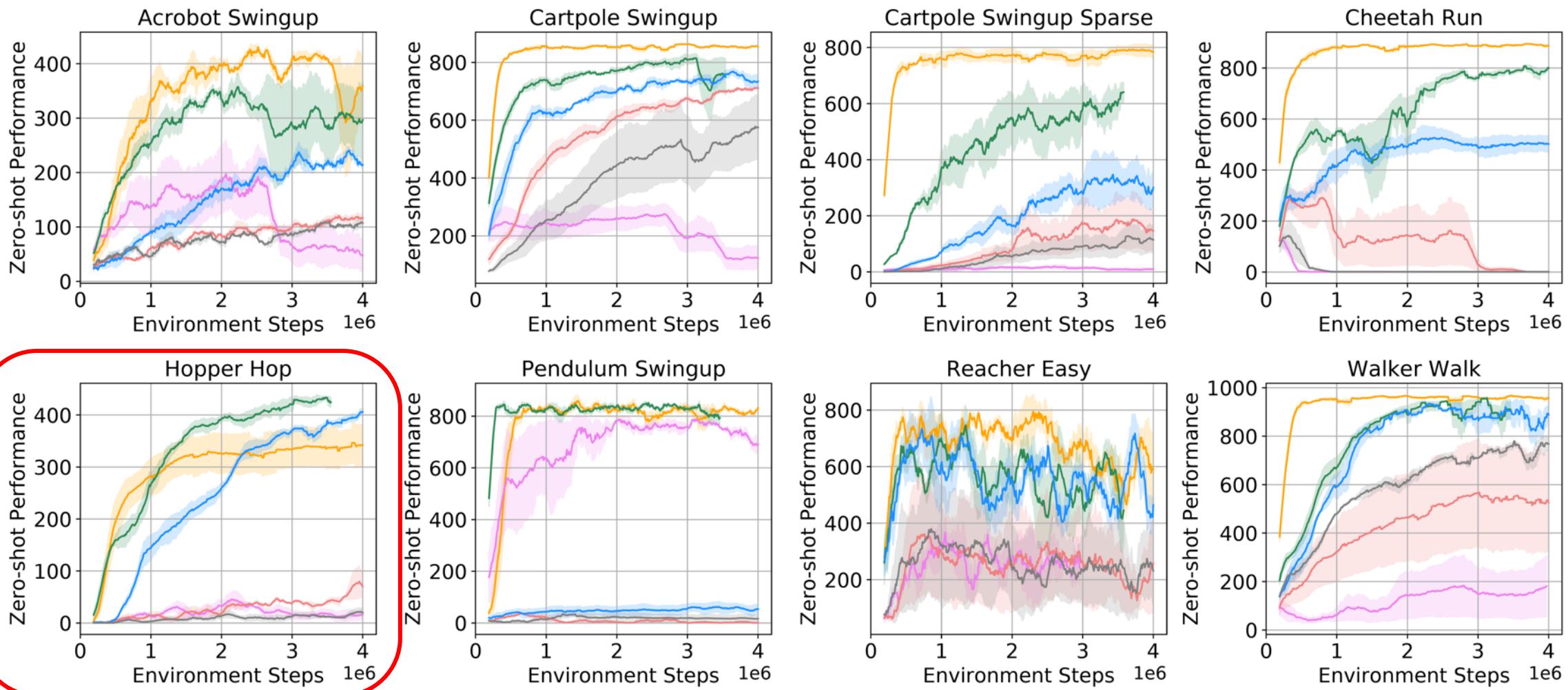


Reacher Easy



Walker Walk

# Zero-Shot Reinforcement Learning



# Experiments Outline

1. Solving a new-task in zero-shot
2. What if we add 20 supervised episodes?
3. Multi-task performance

# Few-Shot Adaptation

Our Agent (few-~~Oracle~~) (supervised)

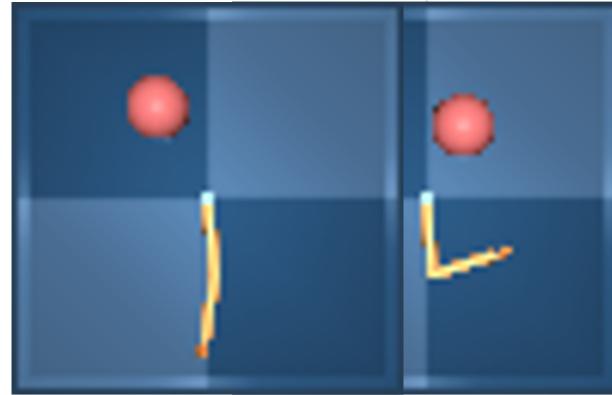


Cheetah Run



Hopper Hop

Our Agent (few-~~Oracle~~) (supervised)



Reacher Easy

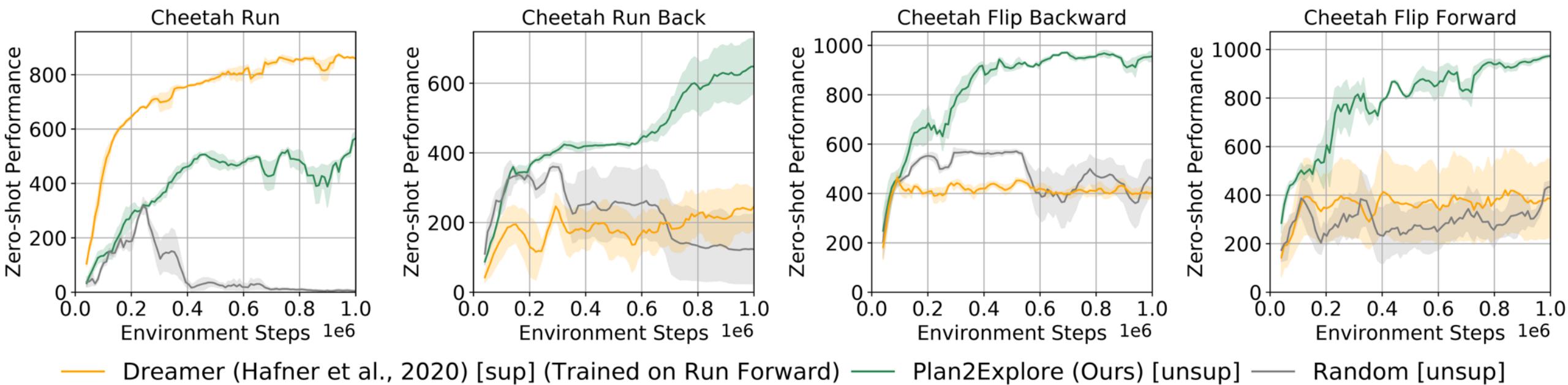


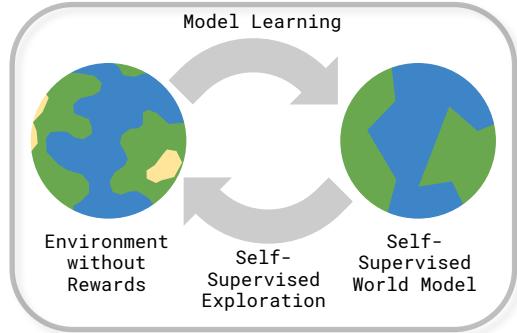
Walker Walk

# Experiments Outline

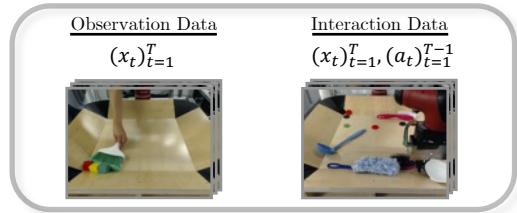
1. Solving a new-task in zero-shot
2. What if we add 20 supervised episodes?
3. Multi-task performance

# Can one model be used for multiple tasks?

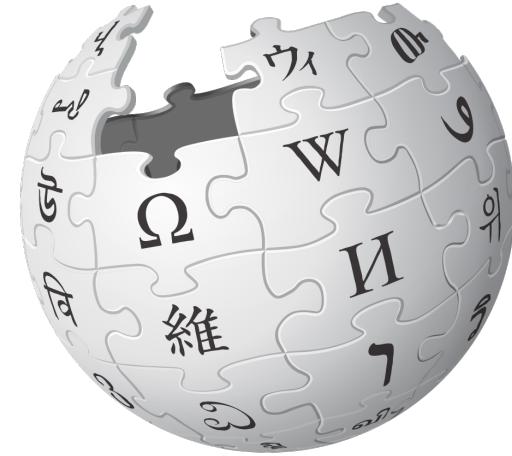




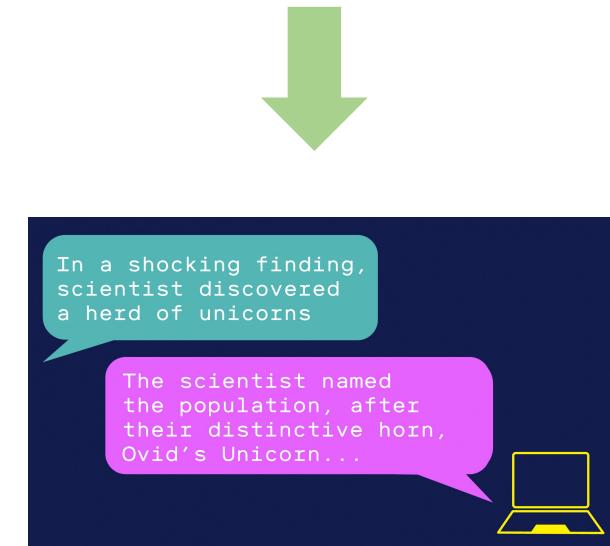
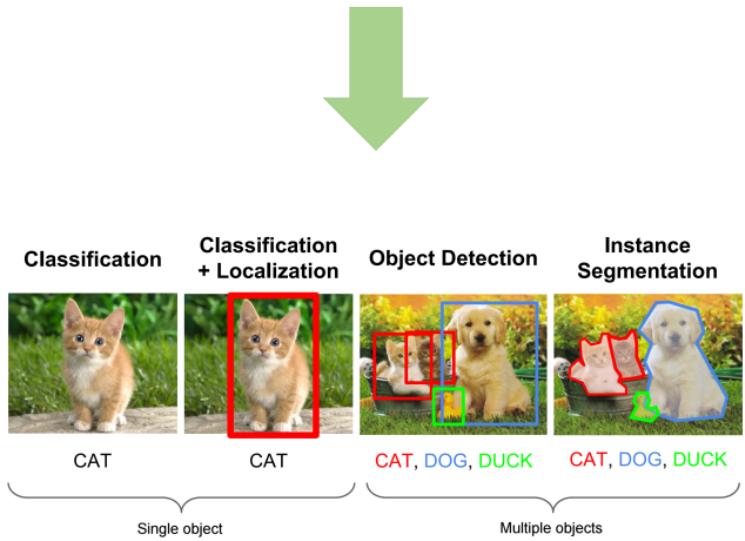
## Self-supervised data collection



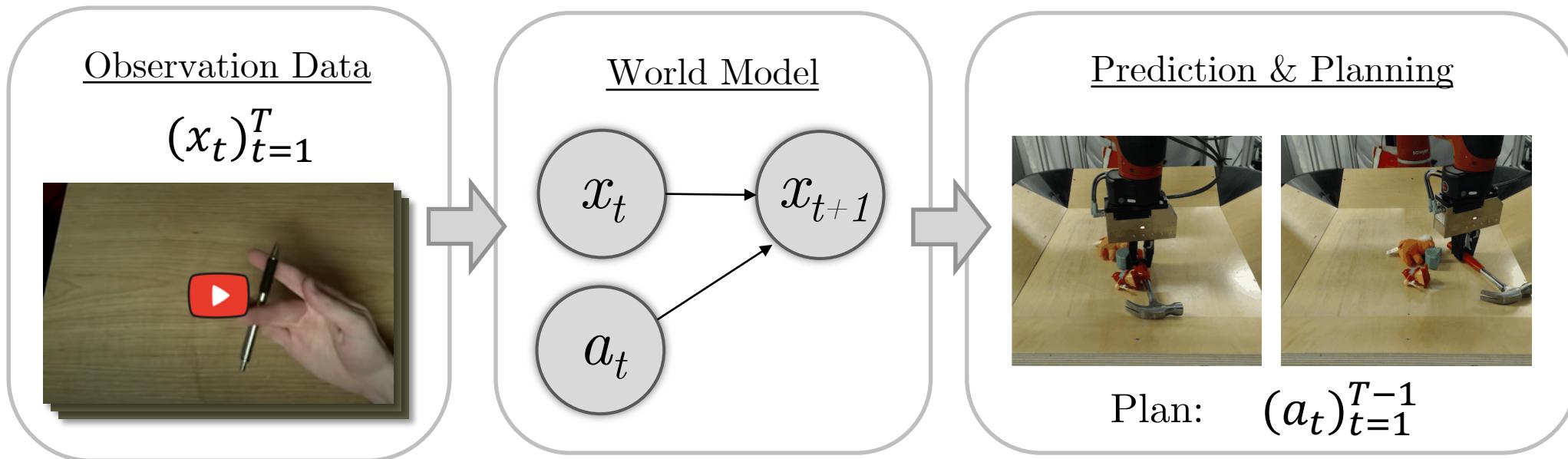
Leveraging human-collected data



# WIKIPEDIA

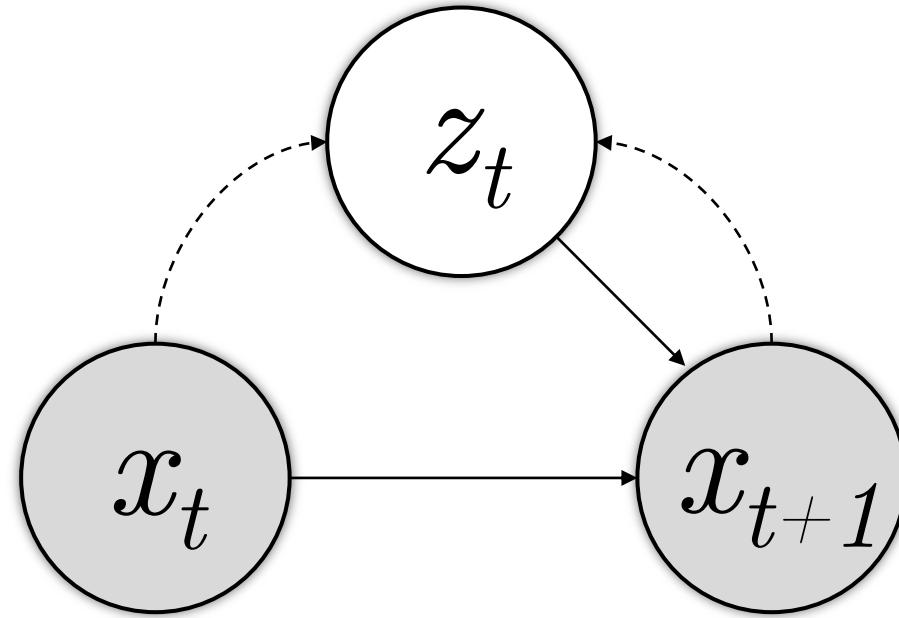


# Learning from Observation

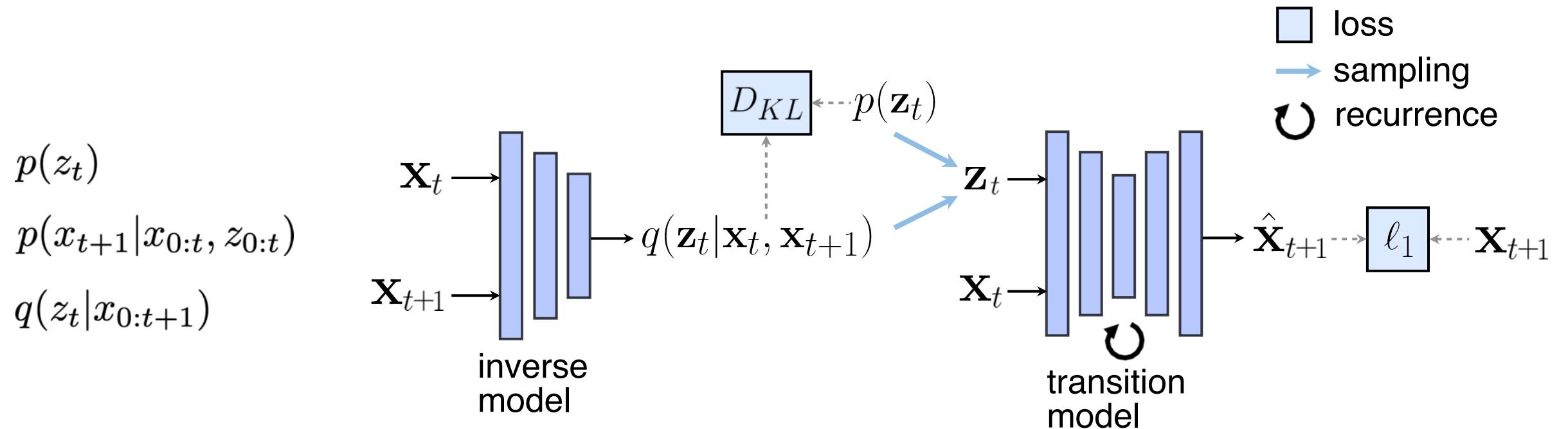


# Learning Action Representations

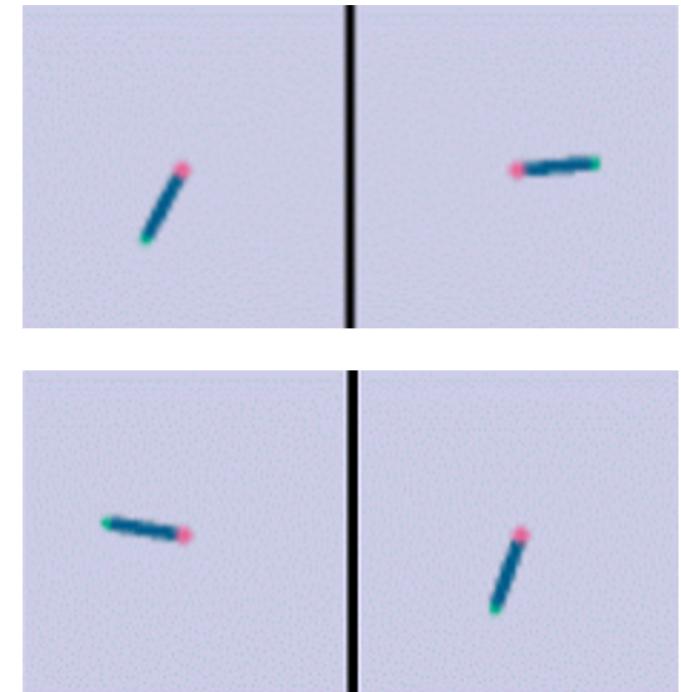
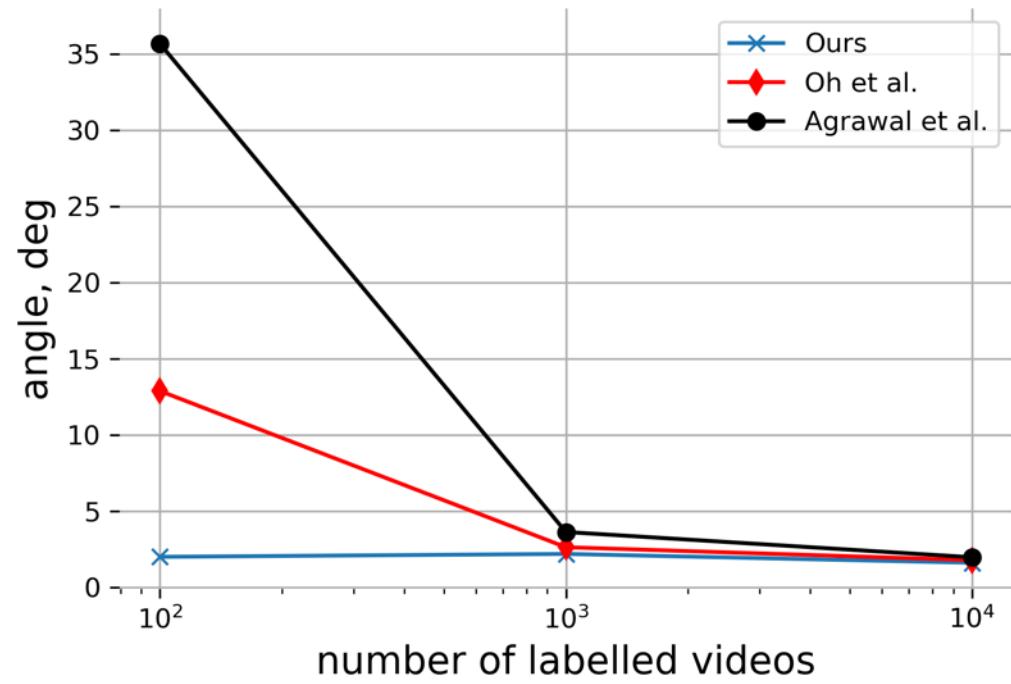
$p(z_t)$   
 $p(x_{t+1}|x_{0:t}, z_{0:t})$   
 $q(z_t|x_{0:t+1})$



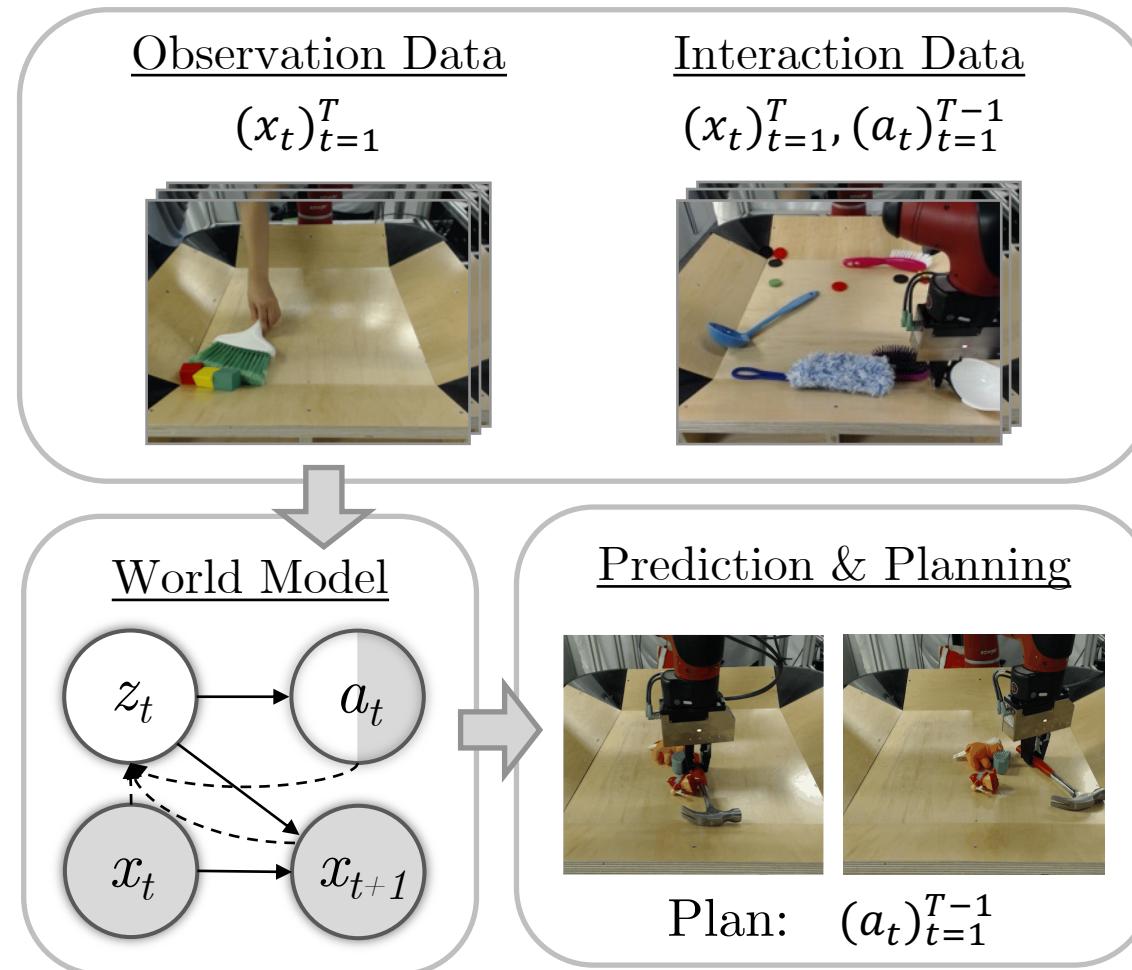
# Learning Action Representations



# Learning Action Representations

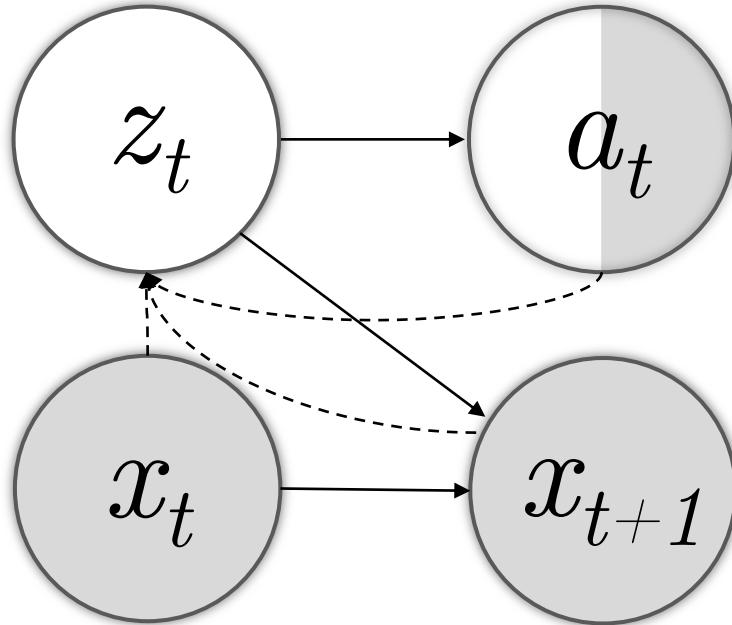


# Learning From Observation and Interaction

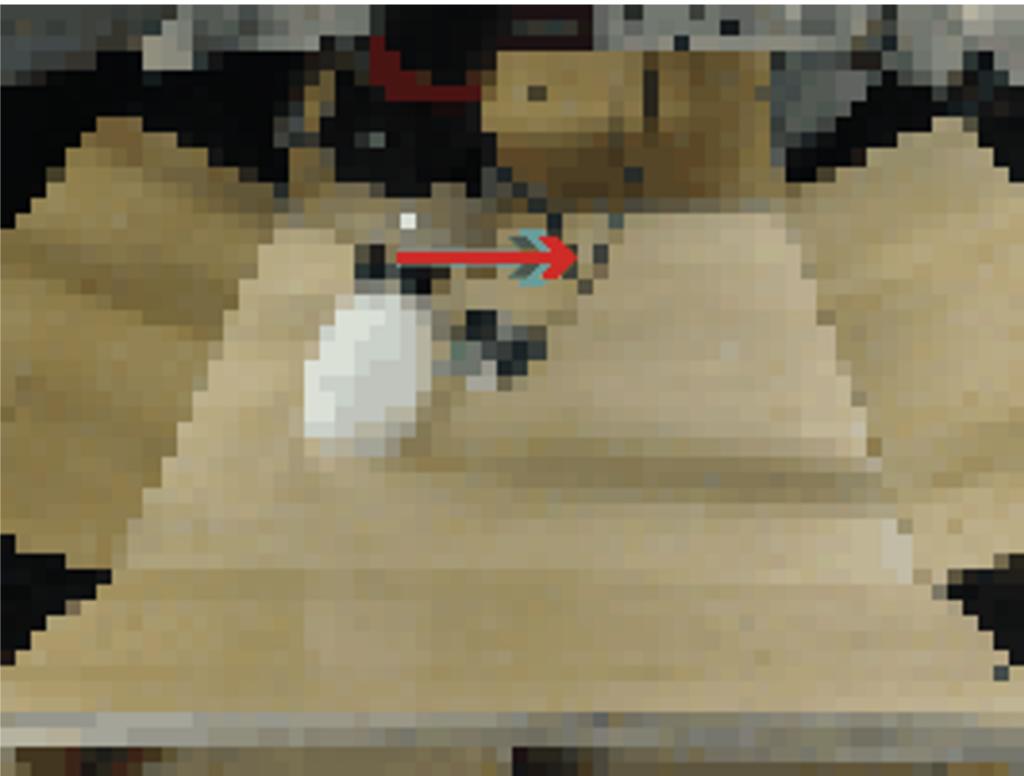


# Learning From Observation and Interaction

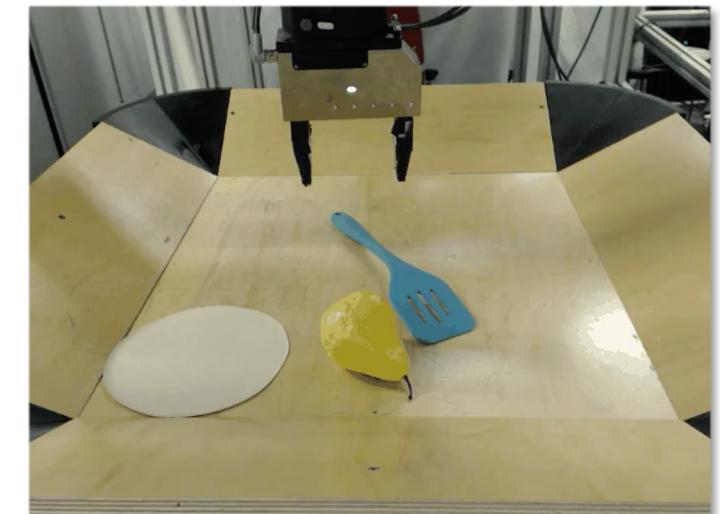
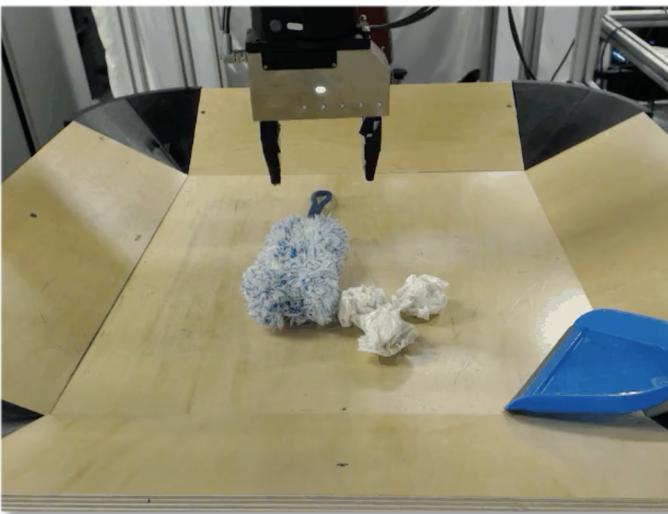
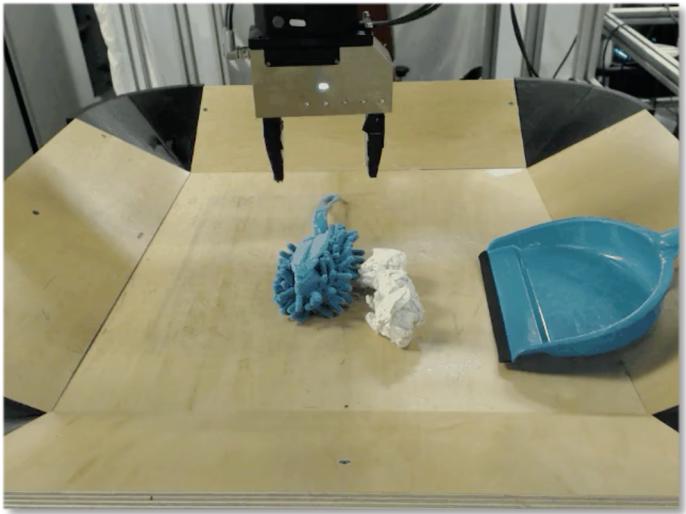
$p(z_t)$   
 $p(x_{t+1}|x_{0:t}, z_{0:t})$   
 $p(a_t|z_t)$   
 $q(z_t|x_{0:t+1})$   
 $q(z_t|a_t)$



# Action Prediction

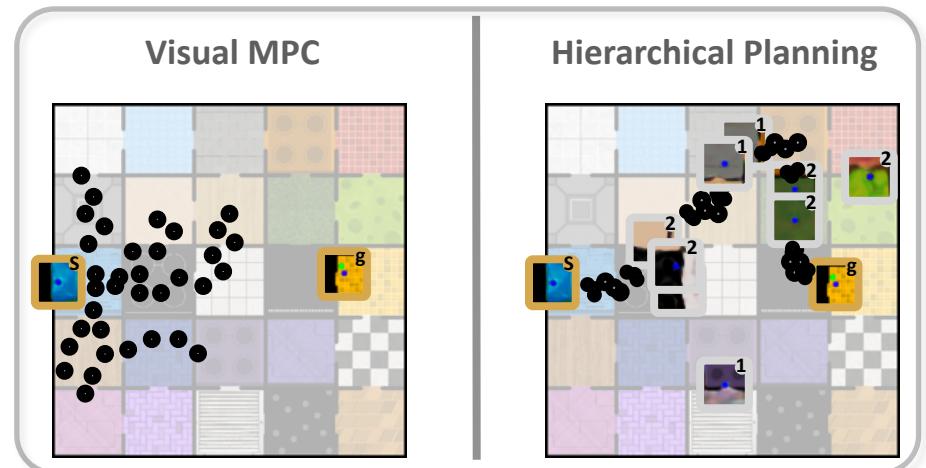
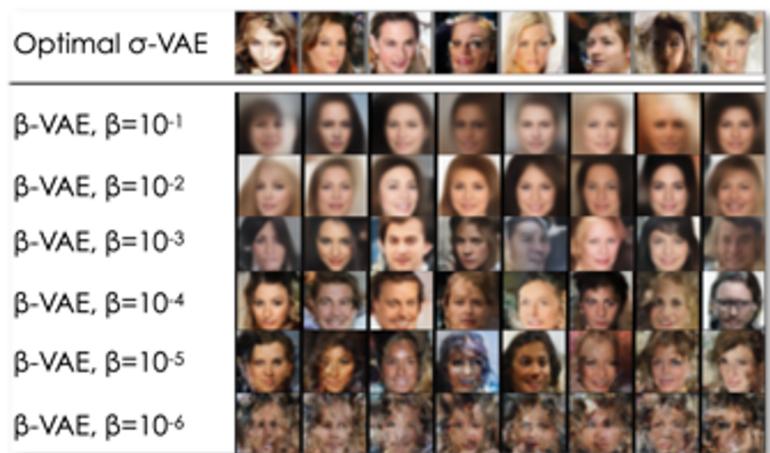
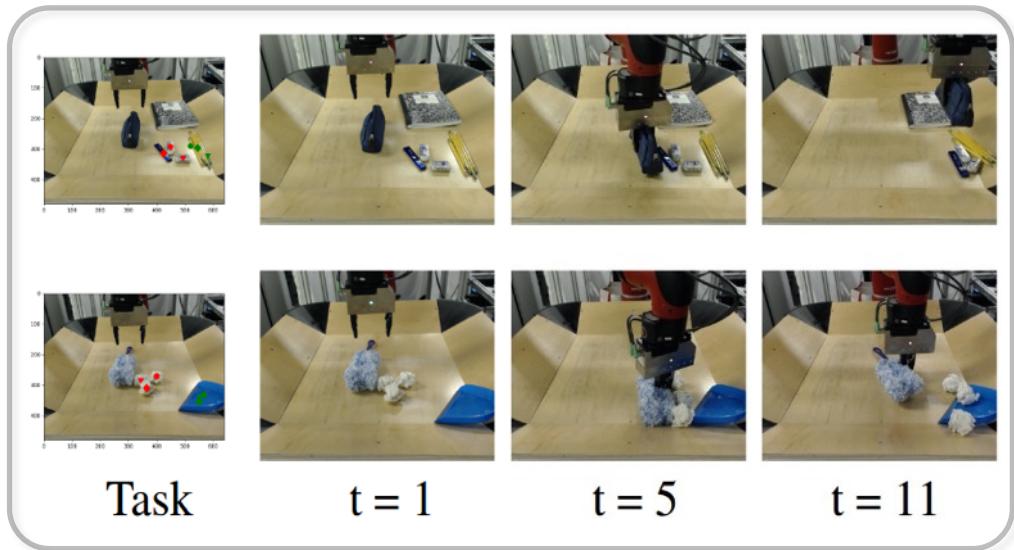
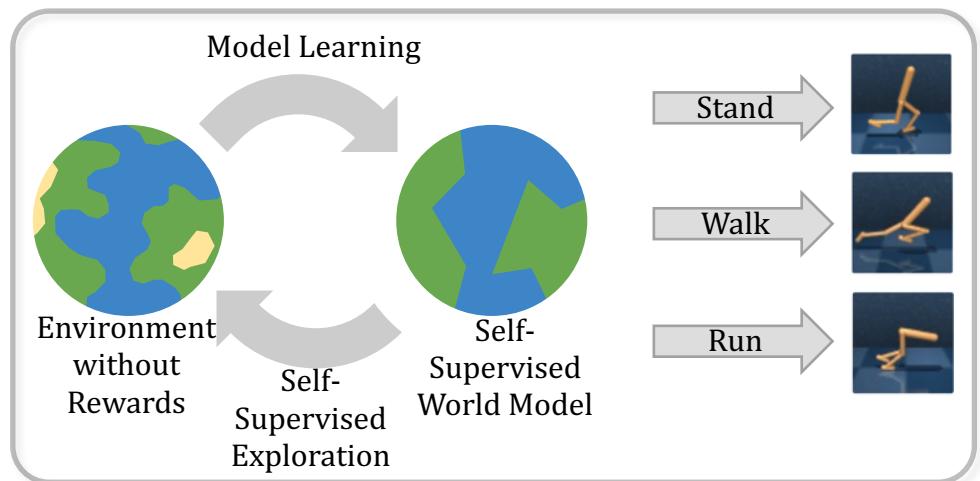


# Robotic Planning



Schmeckpeper, Han, Daniilidis, Rybkin. Visual Planning with Semi-Supervised Action Representations. ICML workshops, 2019

Schmeckpeper, Xie, Rybkin, Tian, Daniilidis, Levine, Finn. Learning Predictive Models from Observation and Interaction. ECCV, 2020.



# Thanks to my amazing collaborators!



Karl  
Pertsch



Ramanan  
Sekar



Danijar  
Hafner



Pieter  
Abbeel



Deepak  
Pathak



Frederik  
Ebert



Dinesh  
Jayaraman



Chelsea  
Finn



Sergey  
Levine



Karl  
Schmeckpeper



Stephen  
Tian



Annie  
Xie



Joseph  
Lim



Jingyun  
Yang



Shenghao  
Zhou



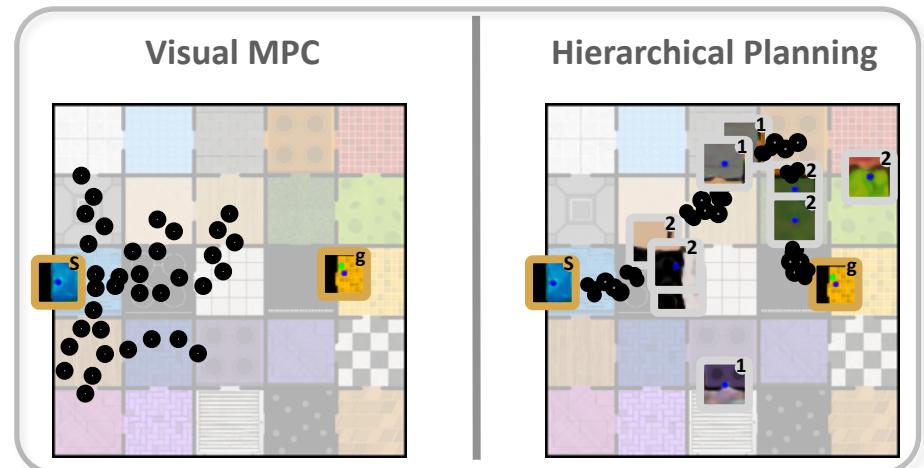
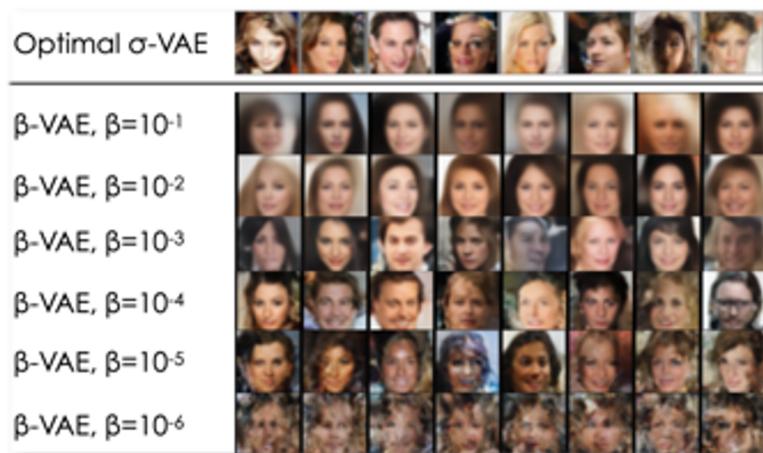
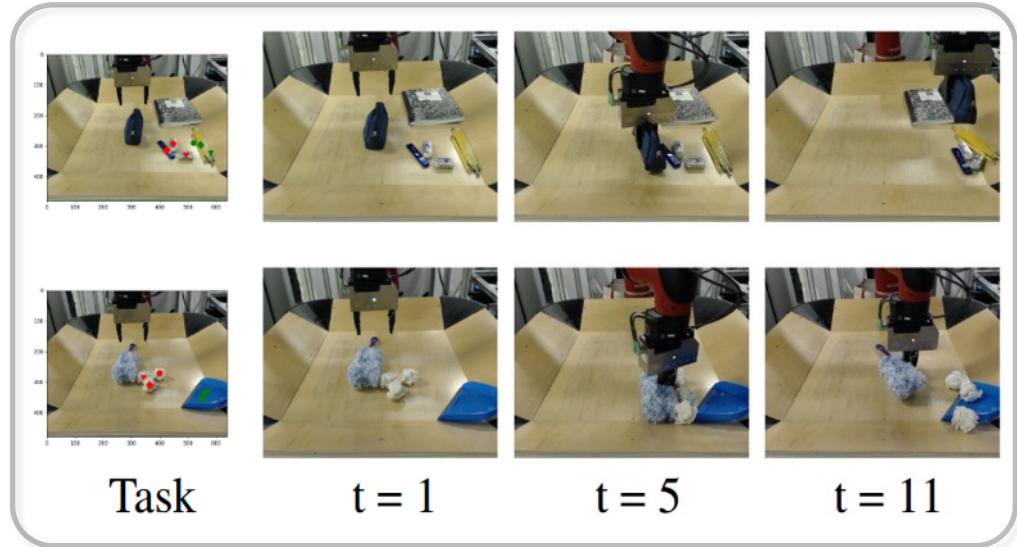
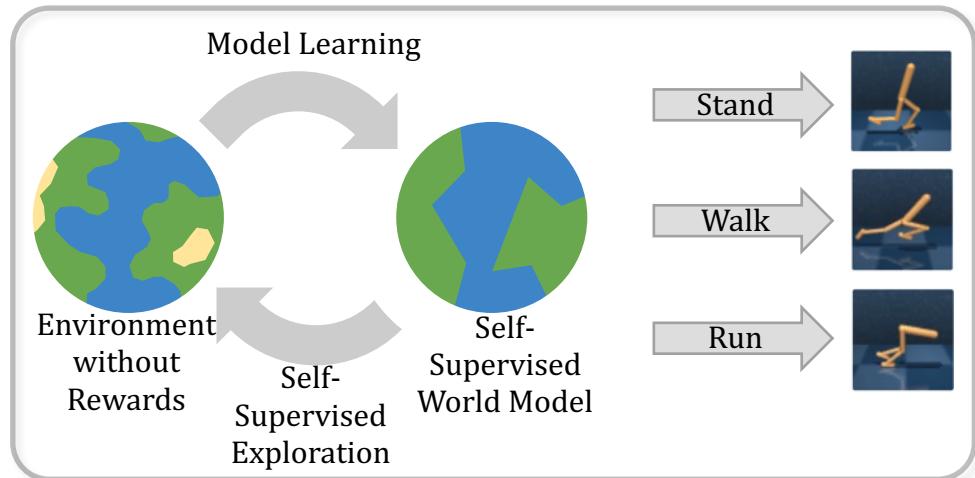
Kosta  
Derpanis



Andrew  
Jaegle



Kostas  
Daniilidis



# Questions?

Oleh Rybkin  
[oleh@seas.upenn.edu](mailto:oleh@seas.upenn.edu)