

Dreamer

But not only v1

William Ruys

Published as a conference paper at ICLR 2020

DREAM TO CONTROL: LEARNING BEHAVIORS BY LATENT IMAGINATION

Danijar Hafner * **Timothy Lillicrap** **Jimmy Ba** **Mohammad Norouzi**
University of Toronto DeepMind University of Toronto Google Brain
Google Brain

Published as a conference paper at ICLR 2021

MASTERING ATARI WITH DISCRETE WORLD MODELS

Danijar Hafner * **Timothy Lillicrap** **Mohammad Norouzi** **Jimmy Ba**
Google Research DeepMind Google Research University of Toronto

Article

Mastering diverse control tasks through world models

<https://doi.org/10.1038/s41586-025-08744-2> Danijar Hafner¹, Jurgis Pasukonis¹, Jimmy Ba² & Timothy Lillicrap¹

Received: 5 April 2024

Dreamer

But not only v1

William Ruys

v1

DREAM TO CONTROL: LEARNING BEHAVIORS BY LATENT IMAGINATION

Danijar Hafner * Timothy Lillicrap Jimmy Ba Mohammad Norouzi
University of Toronto DeepMind University of Toronto Google Brain
Google Brain

Published as a conference paper at ICLR 2020

v2

MASTERING ATARI WITH DISCRETE WORLD MODELS

Danijar Hafner * Timothy Lillicrap Mohammad Norouzi Jimmy Ba
Google Research DeepMind Google Research University of Toronto

v3

Article

Mastering diverse control tasks through world models

<https://doi.org/10.1038/s41586-025-08744-2> Danijar Hafner^{1✉}, Jurgis Pasukonis¹, Jimmy Ba² & Timothy Lillicrap¹

Received: 5 April 2024

Published as a conference paper at ICLR 2021

Dreamer

But not only v1

William Ruys

Architecture + Idea **v1**

Published as a conference paper at ICLR 2020

DREAM TO CONTROL: LEARNING BEHAVIORS
BY LATENT IMAGINATION

Danijar Hafner * Timothy Lillicrap Jimmy Ba Mohammad Norouzi
University of Toronto DeepMind University of Toronto Google Brain
Google Brain

v2

Published as a conference paper at ICLR 2021

MASTERING ATARI WITH DISCRETE WORLD MODELS

Danijar Hafner * Timothy Lillicrap Mohammad Norouzi Jimmy Ba
Google Research DeepMind Google Research University of Toronto

v3

Article

Mastering diverse control tasks through
world models

<https://doi.org/10.1038/s41586-025-08744-2> Danijar Hafner^{1✉}, Jurgis Pasukonis¹, Jimmy Ba² & Timothy Lillicrap¹

Received: 5 April 2024

Dreamer

But not only v1

William Ruys

Architecture + Idea **v1**

Modify for Atari
Discrete Latent Space
Split KL-objective

v2

DREAM TO CONTROL: LEARNING BEHAVIORS
BY LATENT IMAGINATION

Danijar Hafner * Timothy Lillicrap Jimmy Ba Mohammad Norouzi
University of Toronto DeepMind University of Toronto Google Brain
Google Brain

Published as a conference paper at ICLR 2020

MASTERING ATARI WITH DISCRETE WORLD MODELS

Danijar Hafner * Timothy Lillicrap Mohammad Norouzi Jimmy Ba
Google Research DeepMind Google Research University of Toronto

v3

Article
**Mastering diverse control tasks through
world models**

<https://doi.org/10.1038/s41586-025-08744-2> Danijar Hafner¹, Jurgis Pasukonis¹, Jimmy Ba² & Timothy Lillicrap¹

Received: 5 April 2024

Dreamer

But not only v1

William Ruys

Architecture + Idea **v1**

Modify for Atari
Discrete Latent Space
Split KL-objective

v2

Reduce hyperparameter tuning
Reward Scaling
Distributional RL
Vanishing KL-objective

v3

Published as a conference paper at ICLR 2020

DREAM TO CONTROL: LEARNING BEHAVIORS
BY LATENT IMAGINATION

Danijar Hafner * Timothy Lillicrap Jimmy Ba Mohammad Norouzi
University of Toronto DeepMind University of Toronto Google Brain
Google Brain

Published as a conference paper at ICLR 2021

MASTERING ATARI WITH DISCRETE WORLD MODELS

Danijar Hafner * Timothy Lillicrap Mohammad Norouzi Jimmy Ba
Google Research DeepMind Google Research University of Toronto

Article

Mastering diverse control tasks through world models

<https://doi.org/10.1038/s41586-025-08744-2> Danijar Hafner^{1✉}, Jurgis Pasukonis¹, Jimmy Ba² & Timothy Lillicrap¹

Received: 5 April 2024

Dreamer

But not only v1

William Ruys

Architecture + Idea v1

Published as a conference paper at ICLR 2020

DREAM TO CONTROL: LEARNING BEHAVIORS BY LATENT IMAGINATION

Danijar Hafner * Timothy Lillicrap Jimmy Ba Mohammad Norouzi
University of Toronto DeepMind University of Toronto Google Brain
Google Brain

Modify for Atari
Discrete Latent Space
Split KL-objective

v2

Published as a conference paper at ICLR 2021

MASTERING ATARI WITH DISCRETE WORLD MODELS

Danijar Hafner * Timothy Lillicrap Mohammad Norouzi Jimmy Ba
Google Research DeepMind Google Research University of Toronto

Reduce hyperparameter tuning
Reward Scaling
Distributional RL
Vanishing KL-objective

v3

Article

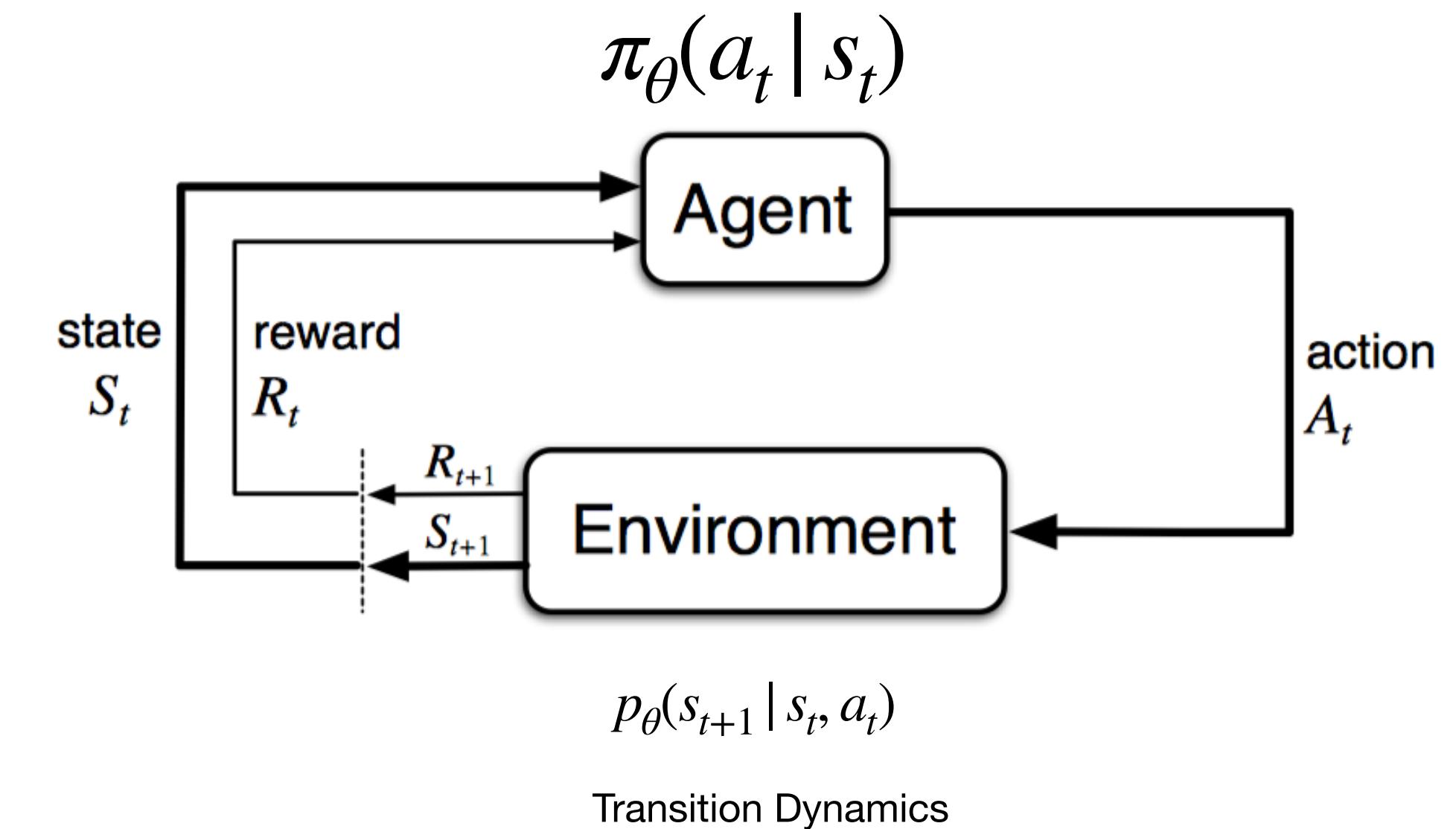
Mastering diverse control tasks through world models

<https://doi.org/10.1038/s41586-025-08744-2> Danijar Hafner^{1✉}, Jurgis Pasukonis¹, Jimmy Ba² & Timothy Lillicrap¹

Received: 5 April 2024

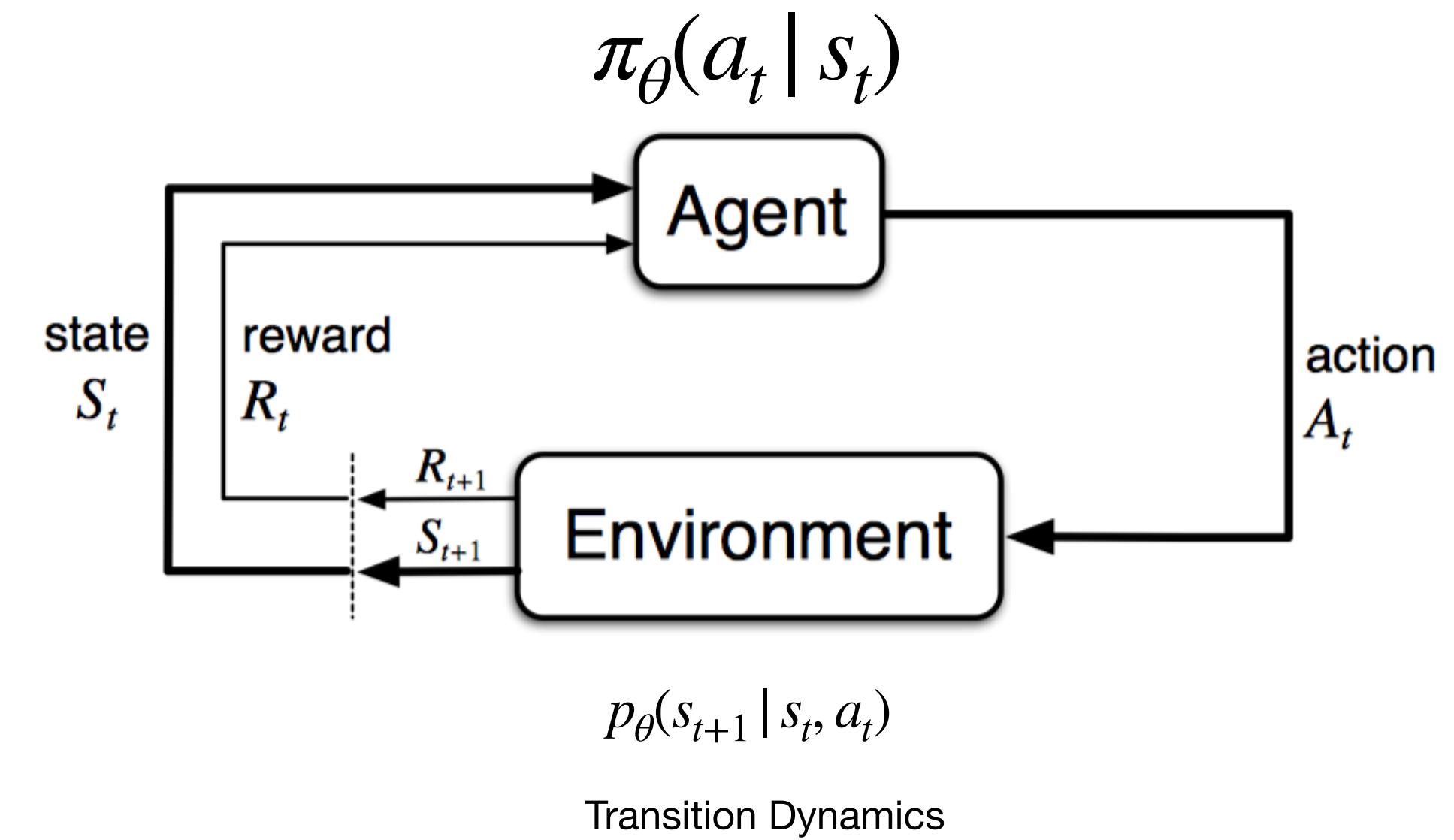
Reinforcement Learning

- We seek to learn a policy $\pi_\theta(a_t | s_t)$ to maximize total expected return
- Return of Trajectory $R(\tau) = \sum_{t=0}^T \gamma^t r(s_t, a_t) \quad (s_t, a_t) \in \tau$
- Maximize: $J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$
- Might only have observations not states
- Rewards might be stochastic



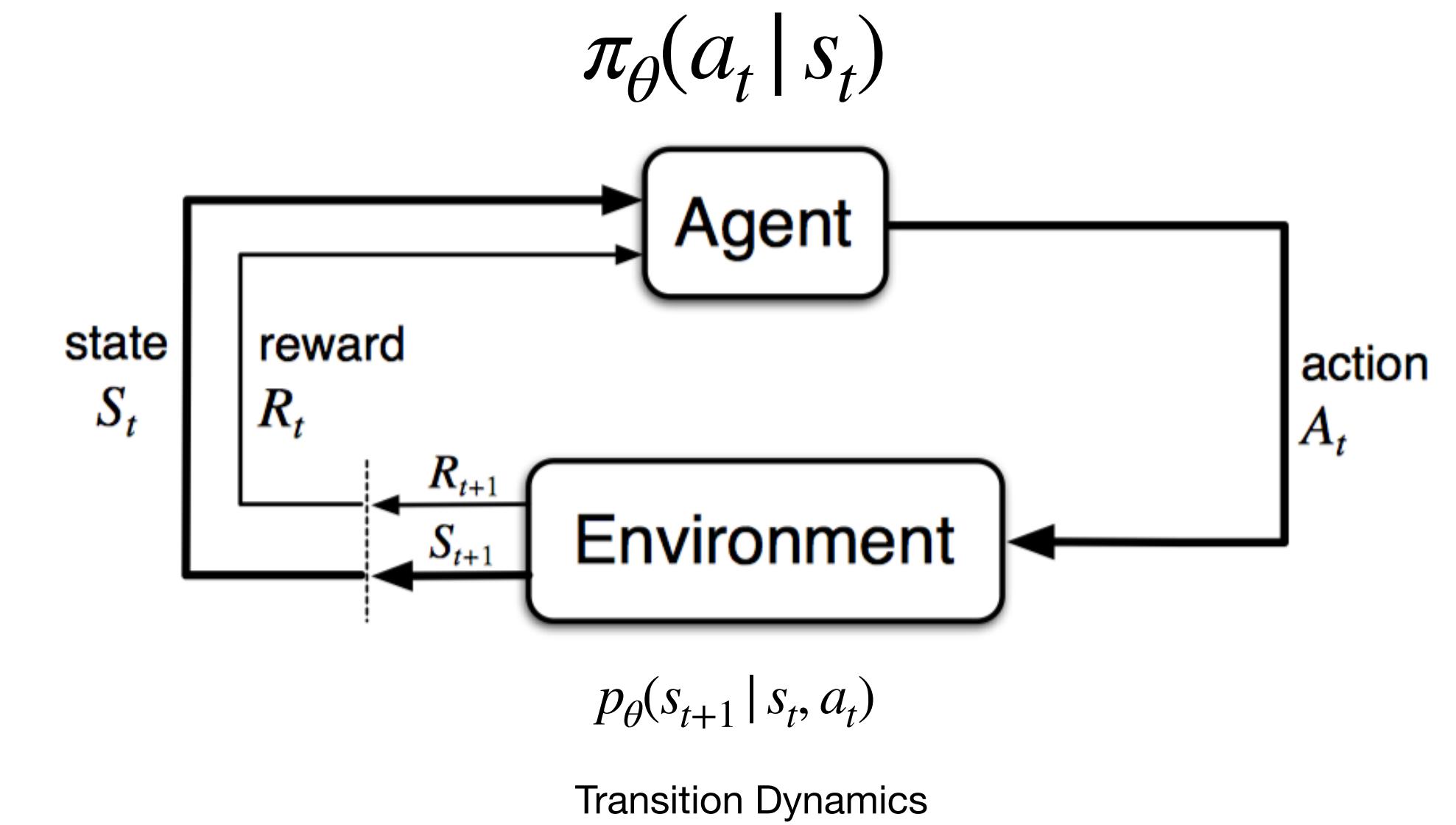
Policy Gradient Methods

- Optimize the objective $J(\theta)$ directly.
 - Gradient ascent on: $\nabla_{\theta}J(\theta) \propto \mathbb{E}_{a_t \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi_{\theta}}(s_t, a_t)]$
- Collect transitions (s_t, a_t, s_{t+1}, r_t) from env
- Increase log-likelihood of high return actions
- Repeat



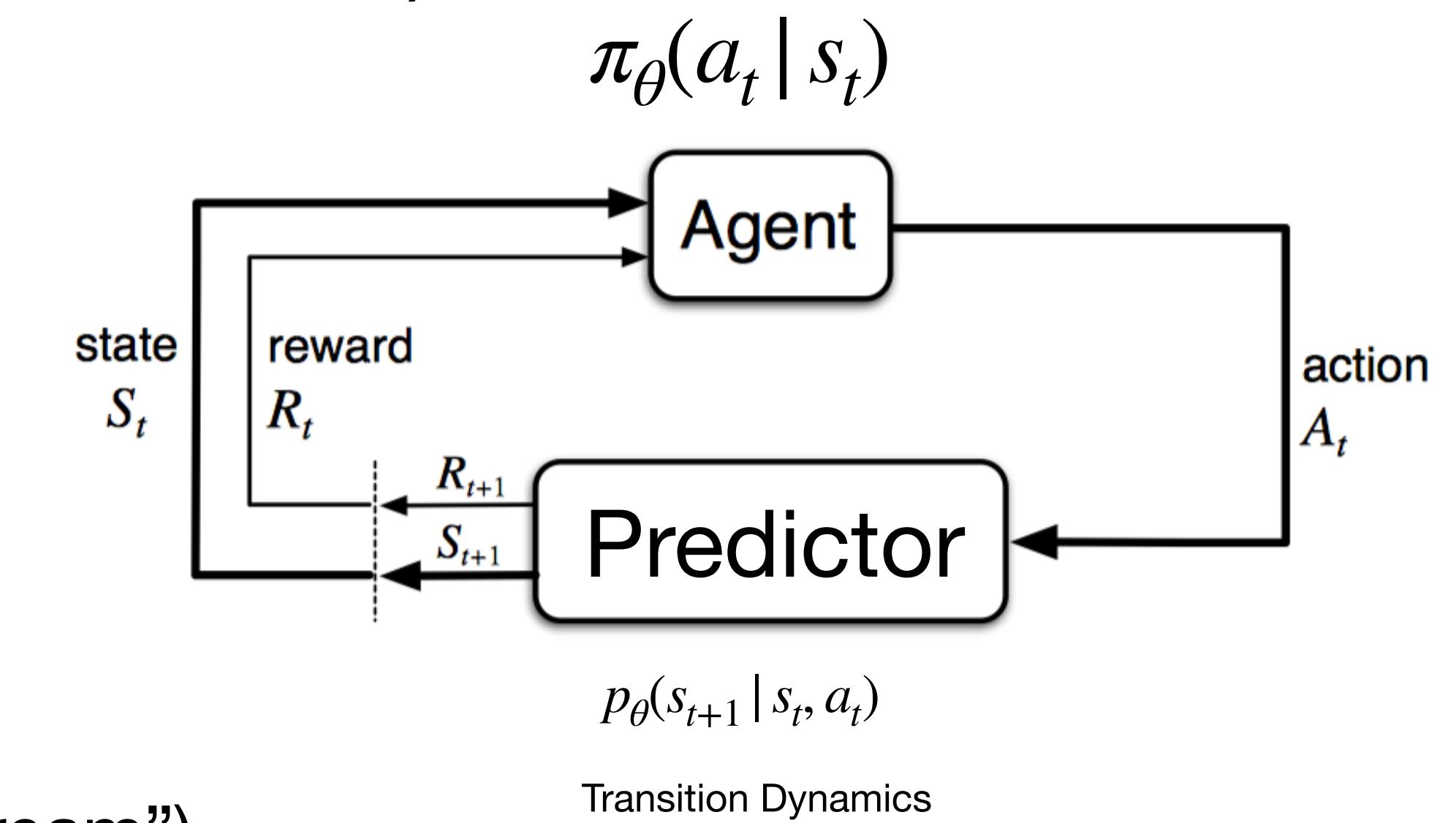
What is Dreamer?

- Policy Gradient methods are *sample inefficient*
 - Need many environment interactions ($\sim 100k - 1B$ transitions)
 - Environment may be ‘slow’ or costly (real time, etc)



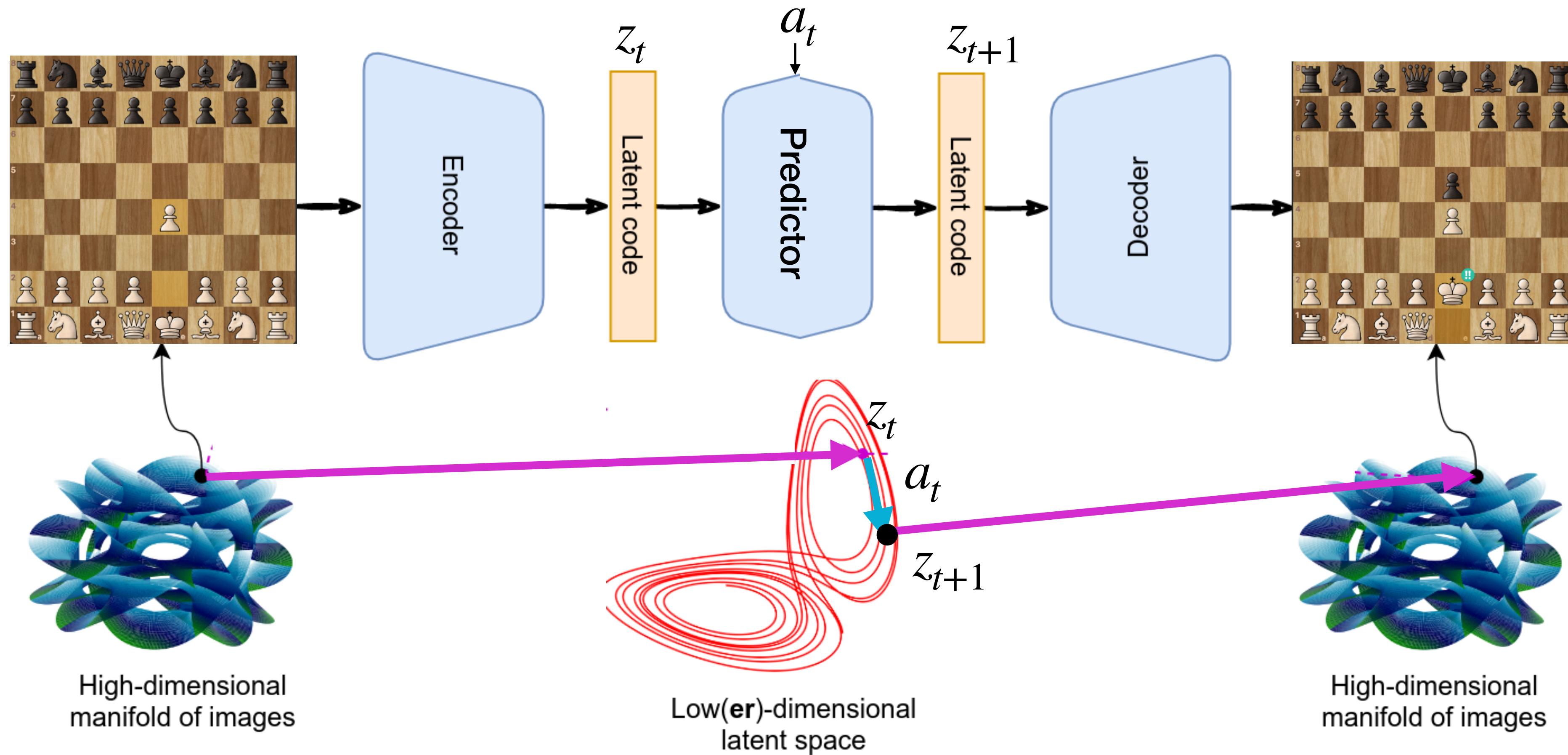
What is Dreamer?

- Policy Gradient methods are *sample inefficient*
 - Need many environment interactions ($\sim 100k - 1B$ transitions)
 - Environment may be ‘slow’ or costly (real time, etc)

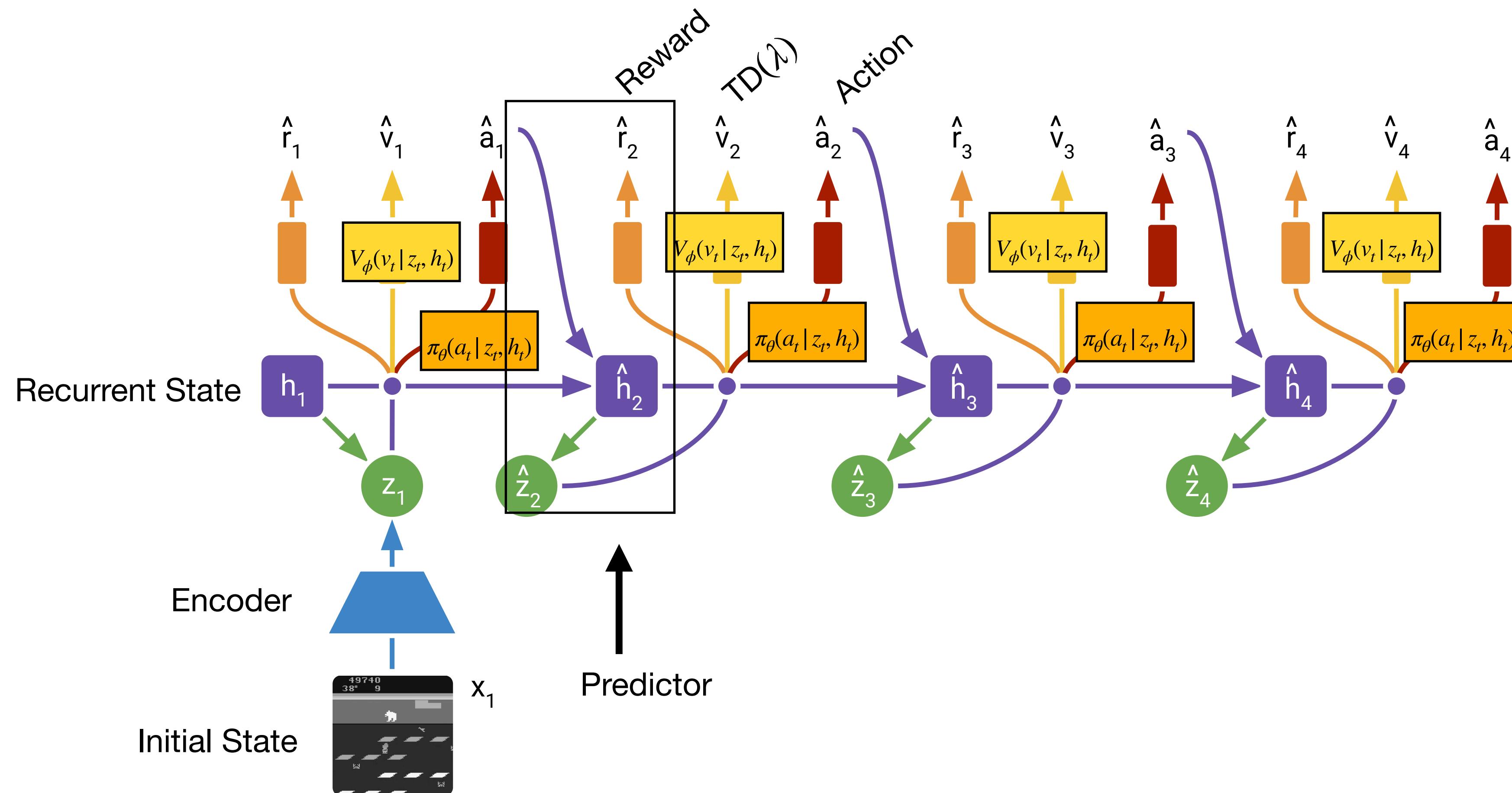


- Replace real interactions with prediction
 - Neural network: model transitions and rewards
 - Network inference faster than real interaction
 - Perform policy gradient with predicted dynamics (“dream”)

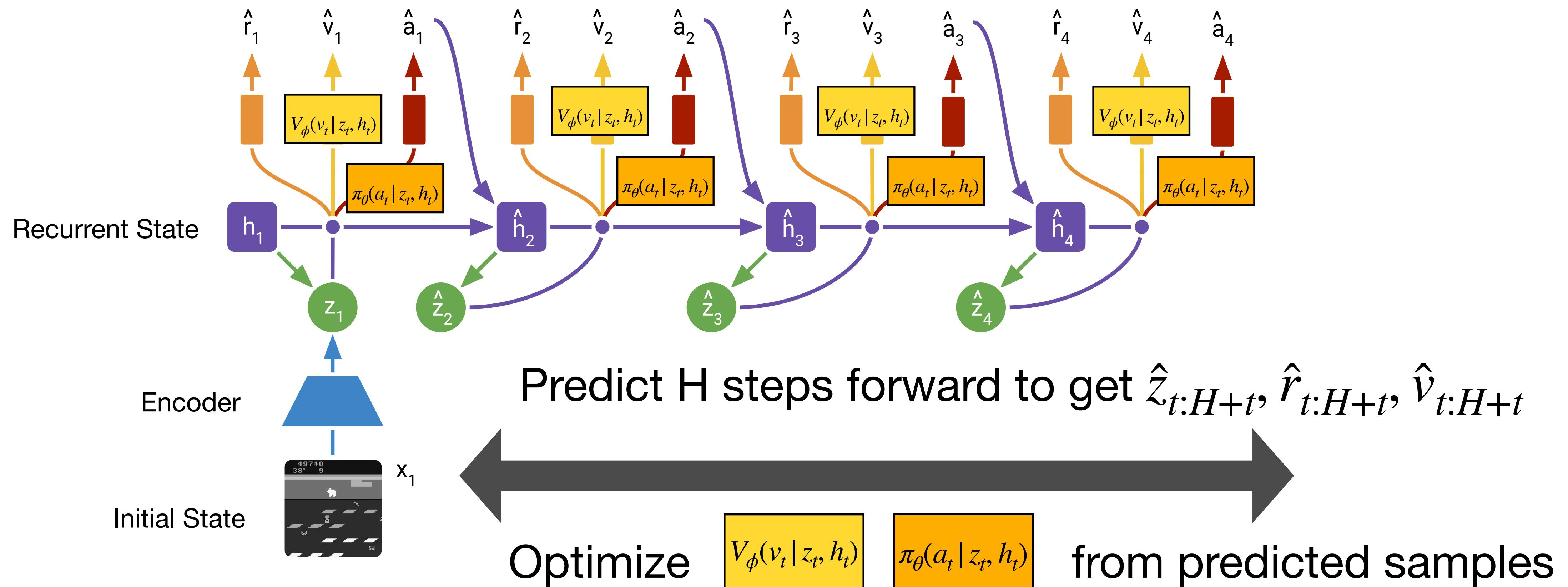
Learn dynamics in lower-dimensional space



Recurrent State Space Model (RSSM)



Recurrent State Space Model (RSSM)

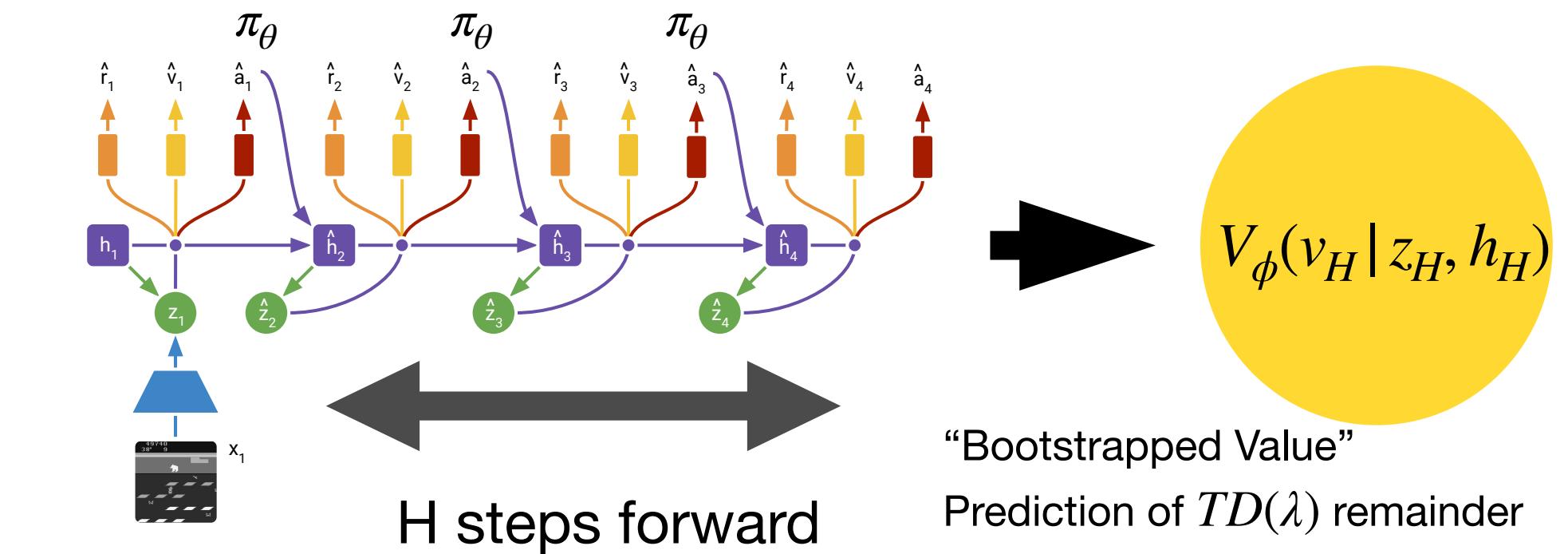


Behavior Learning / Policy Improvement

- Sample a starting state s_t
- ‘Imagine’ trajectories $\{(s_k, a_k, r_k)\}_{k=t}^{t+H}$ from s_t using π_θ
- Compute value estimates:

$$V_{\text{TD}(\lambda)}(s_t) = (1 - \lambda) \sum_{k=1}^H \lambda^{k-1} V(s_t, k)$$

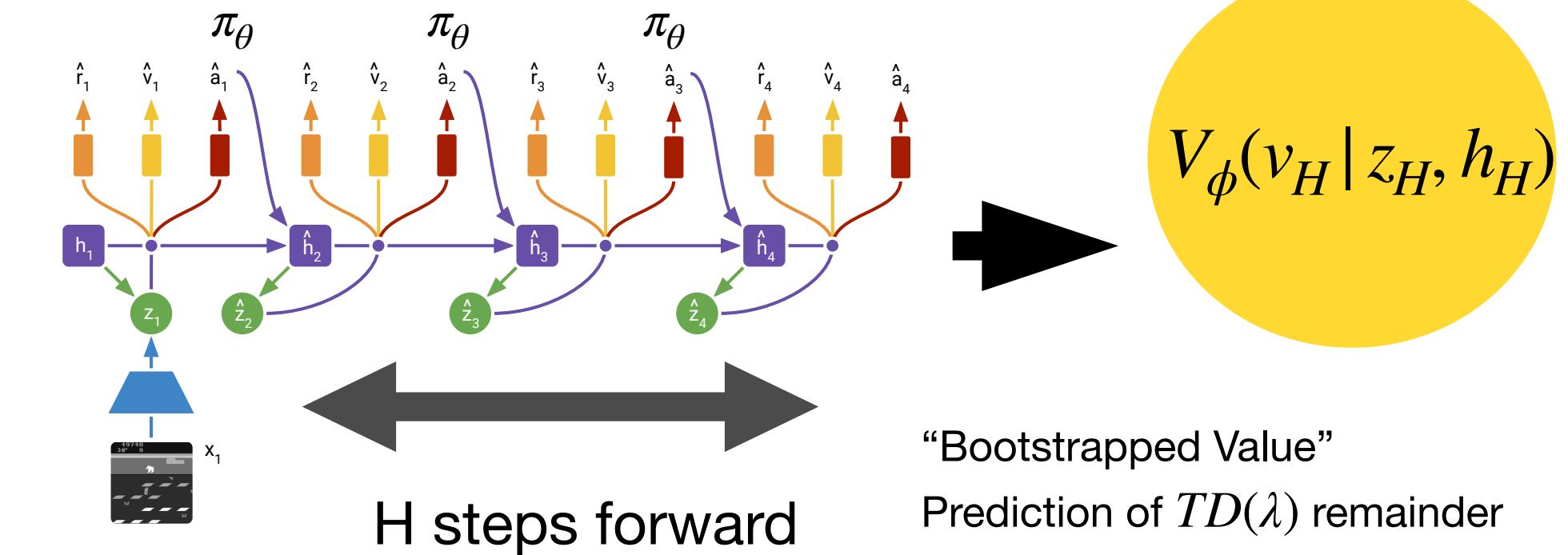
$$V(s_t, k) = \sum_{j=k}^{t+k} \gamma^{j-t} r_j + \gamma^{k-t} V_\phi(v_{t+k} | z_{t+k}, h_{t+k})$$



- Policy Update: $\max_{\theta} \mathbb{E}_{\pi_\theta, v_\phi} \left[\sum_{k=t}^{t+H} V_{\text{TD}(\lambda)}(s_k) \right]$
 $\boxed{\pi_\theta(a_t | z_t, h_t)}$
- Critic / Bootstrap Update: $\min_{\phi} \mathbb{E}_{\pi_\theta, v_\phi} \left[\sum_{k=t}^{t+H} \frac{1}{2} || V_\phi(s_k) - V_{\text{TD}(\lambda)}(s_k) ||^2 \right]$
 $\boxed{V_\phi(v_t | z_t, h_t)}$

Behavior Learning / Policy Improvement

- Sample a starting state s_t
- ‘Imagine’ trajectories $\{(s_k, a_k, r_k)\}_{k=t}^{t+H}$ from s_t using π_θ
- Compute value estimates:

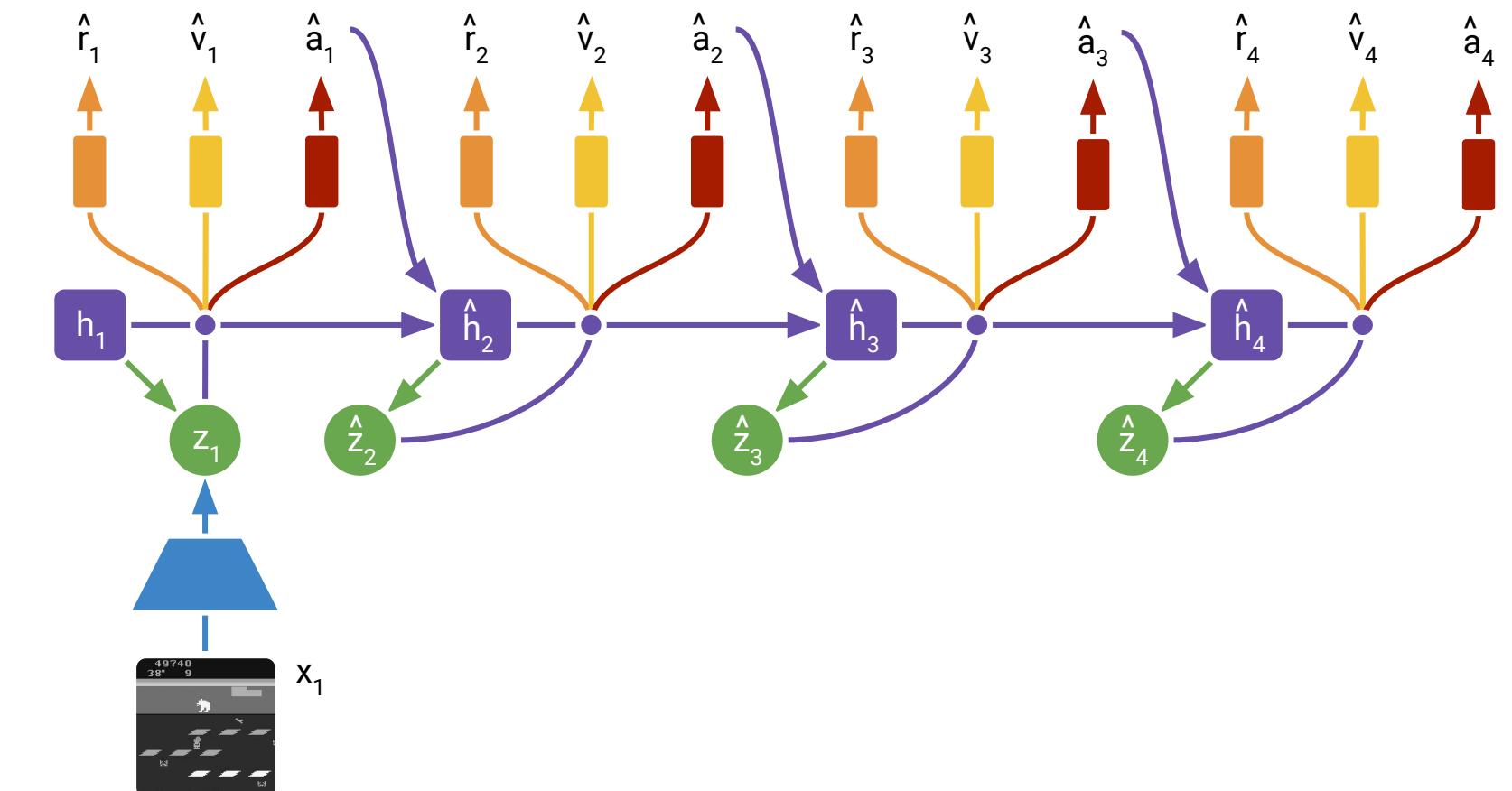


$$V_{TD(\lambda)}(s_t) = (1 - \lambda) \sum_{k=1}^H \lambda^{k-1} V(s_t, k)$$

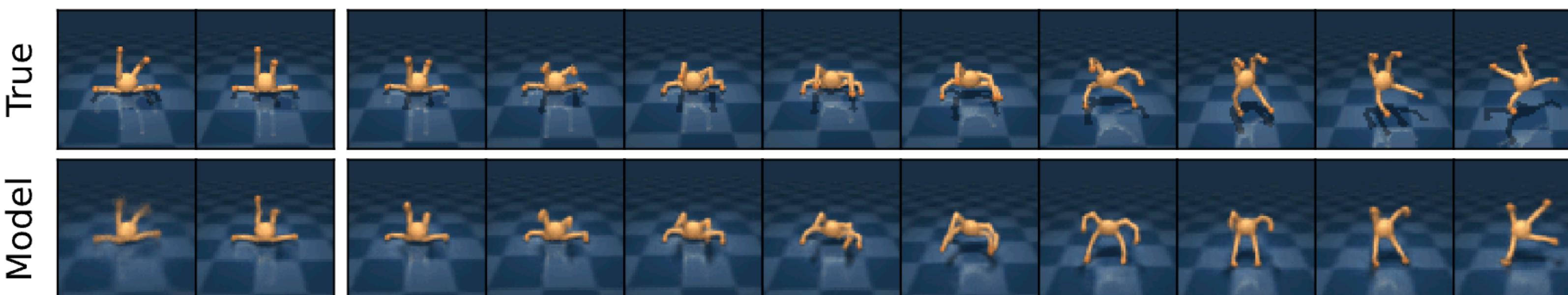
$$V(s_t, k) = \sum_{j=k}^{t+k} \gamma^{j-t} r_j + \gamma^{k-t} V_\phi(v_{t+k} | z_{t+k}, h_{t+k})$$

- Policy Update: $\max_{\theta} \mathbb{E}_{\pi_\theta, v_\phi} \left[\sum_{k=t}^{t+H} V_{TD(\lambda)}(s_k) \right]$
 $\boxed{\pi_\theta(a_t | z_t, h_t)}$
- Critic / Bootstrap Update: $\min_{\phi} \mathbb{E}_{\pi_\theta, v_\phi} \left[\sum_{k=t}^{t+H} \frac{1}{2} || V_\phi(s_k) - V_{TD(\lambda)}(s_k) ||^2 \right]$
 $\boxed{V_\phi(v_t | z_t, h_t)}$
- Use any method we’ve discussed so far.
- Dreamer uses REINFORCE to train policy
 - W/ batch-wise entropy regularization

Need accurate predictor



Context Input Open Loop Prediction



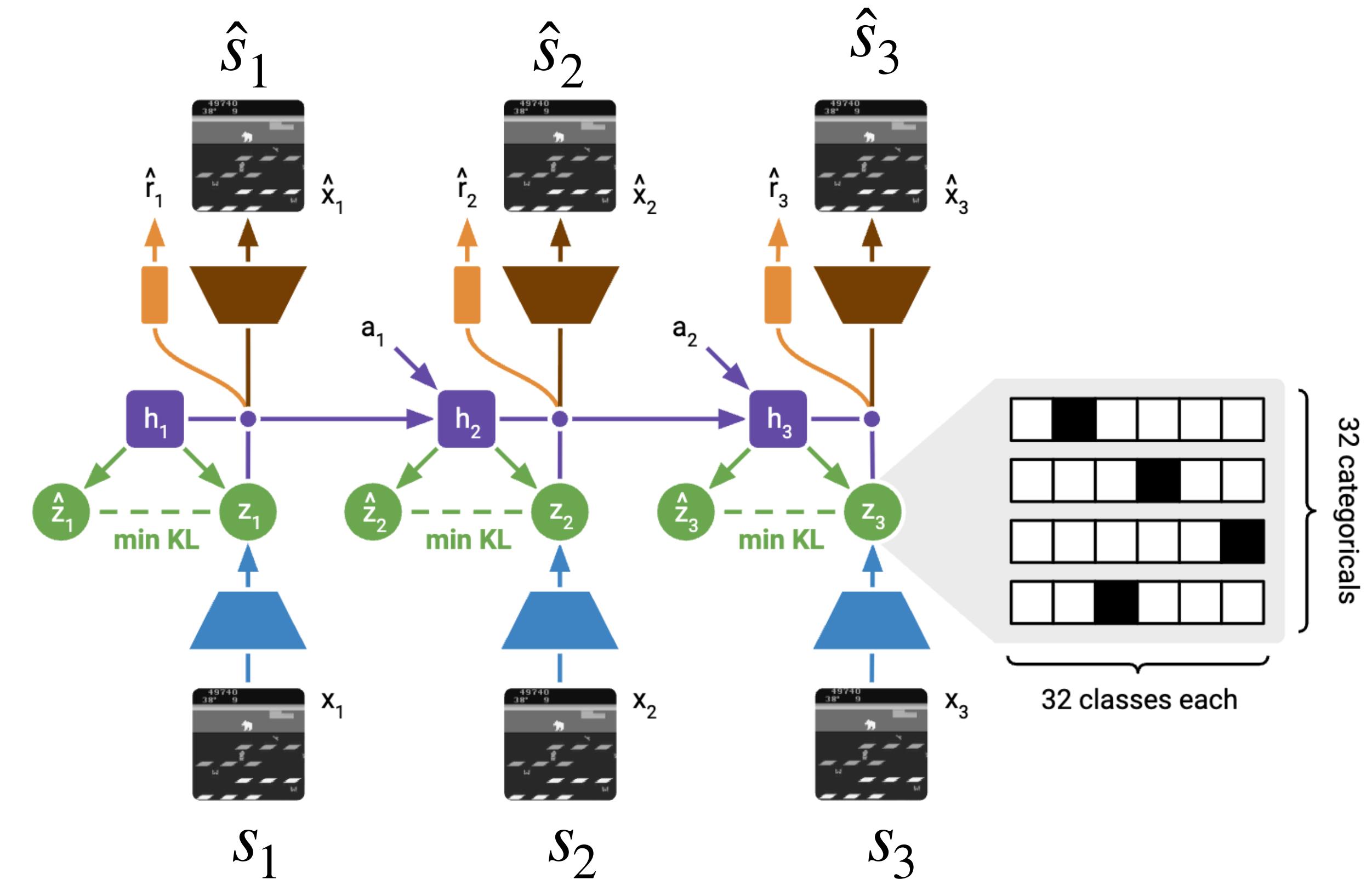
T = 0 5 10 15 20 25 30 35 40 45 50

How do we build a predictor?

World Model Training

- Given representative trajectory $\{(s_k, a_k, r_k)\}_{k=t}^{t+H}$ (assume collected with ϵ -Greedy on current policy)
- Learn a stochastic model to generate it using **only** lower dimensional features

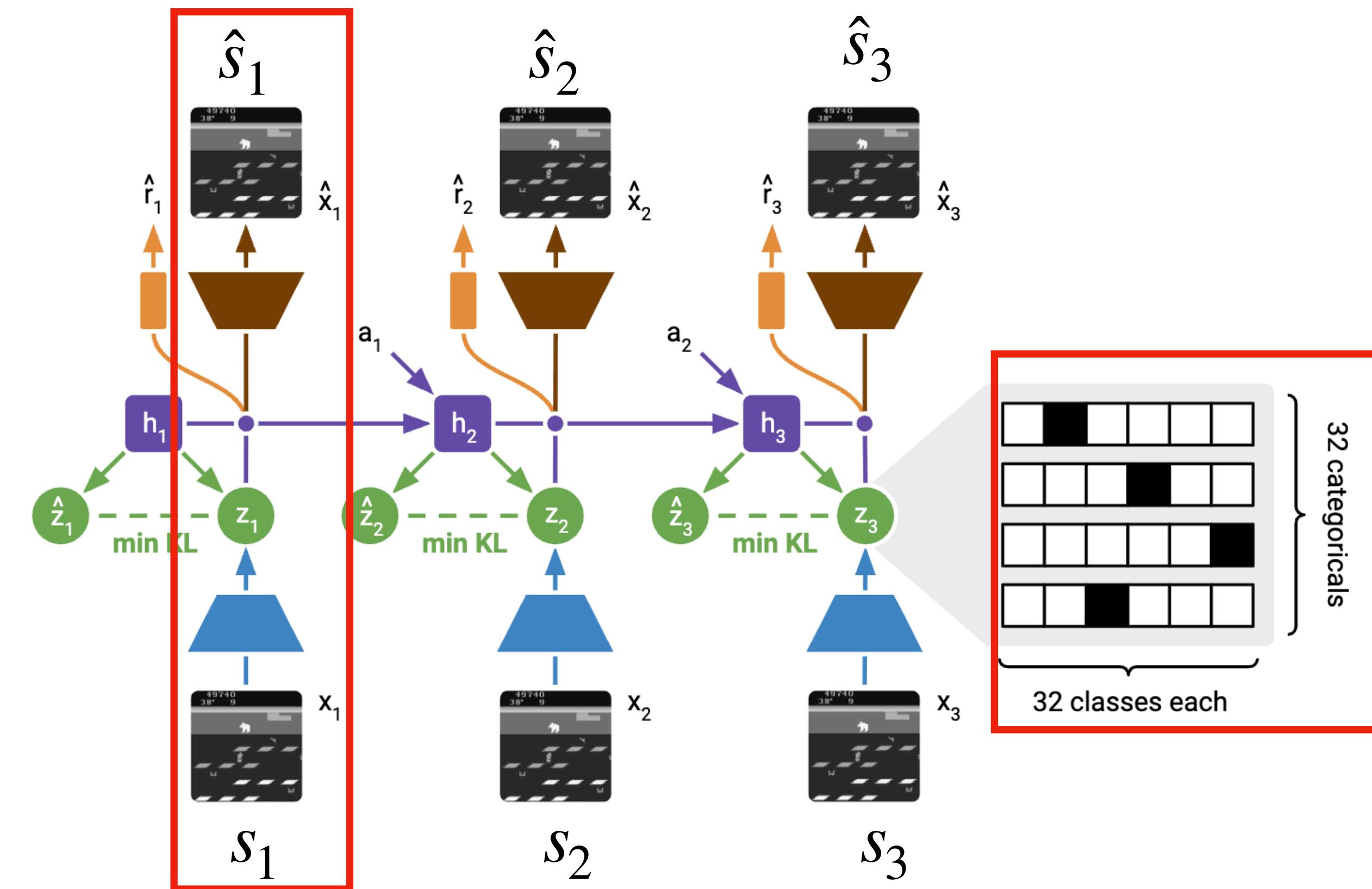
- Multiobjective:
 - Reconstruction from lower dimension:
 - $d(s_t, \hat{s}_t)$ small, learn good state encoder
 - $d(r_t, \hat{r}_t)$ small, predict rewards from state-actions
 - Predict state from history
 - $d(z_t, \hat{z}_t)$ small
 - Predict next encoder without real system observation



Building the RSSM

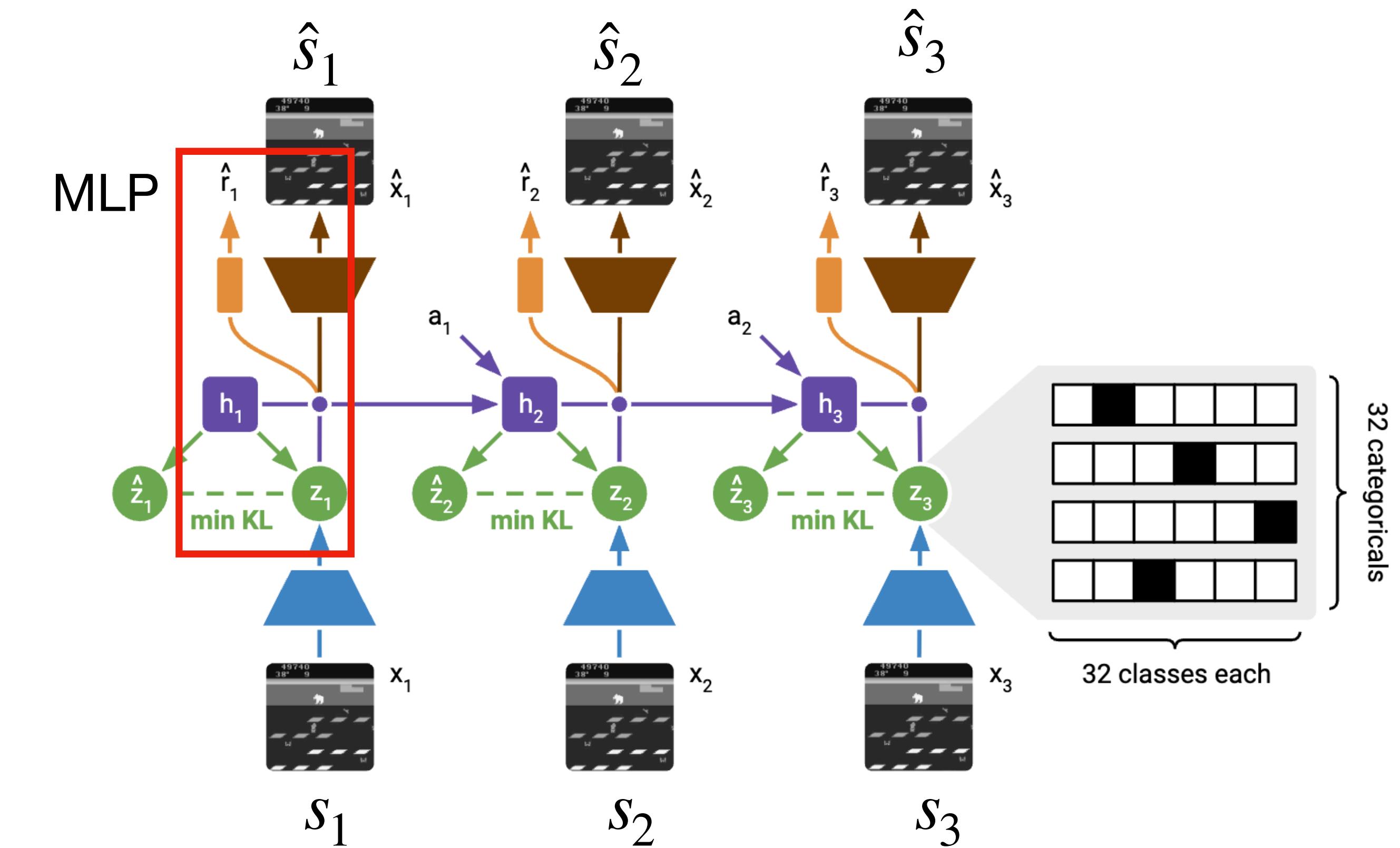
- Enforce z_t is from categorical distribution
- Encoder: $z_t \sim q_\theta(z_t | h_t, s_t)$
- Decoder: $\hat{s}_t \sim p_\theta(\hat{s}_t | h_t, z_t)$

Sequential Discrete Variational Autoencoder (VQ-VAE)



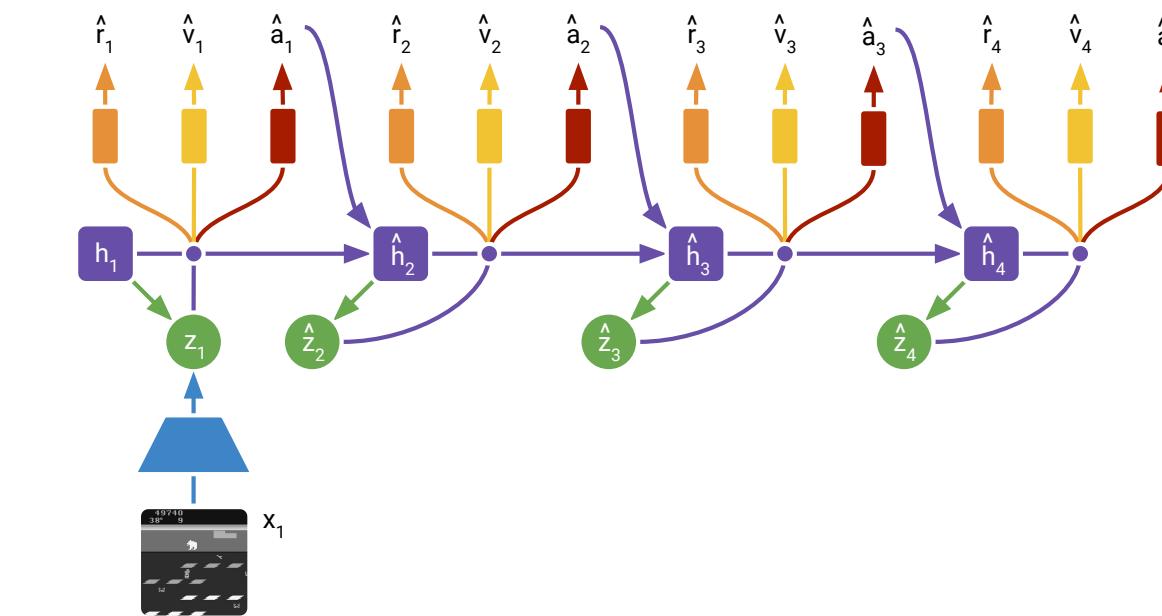
Building the RSSM

- Reward Model: $r_t \sim p_\theta(r_t | h_t, z_t)$

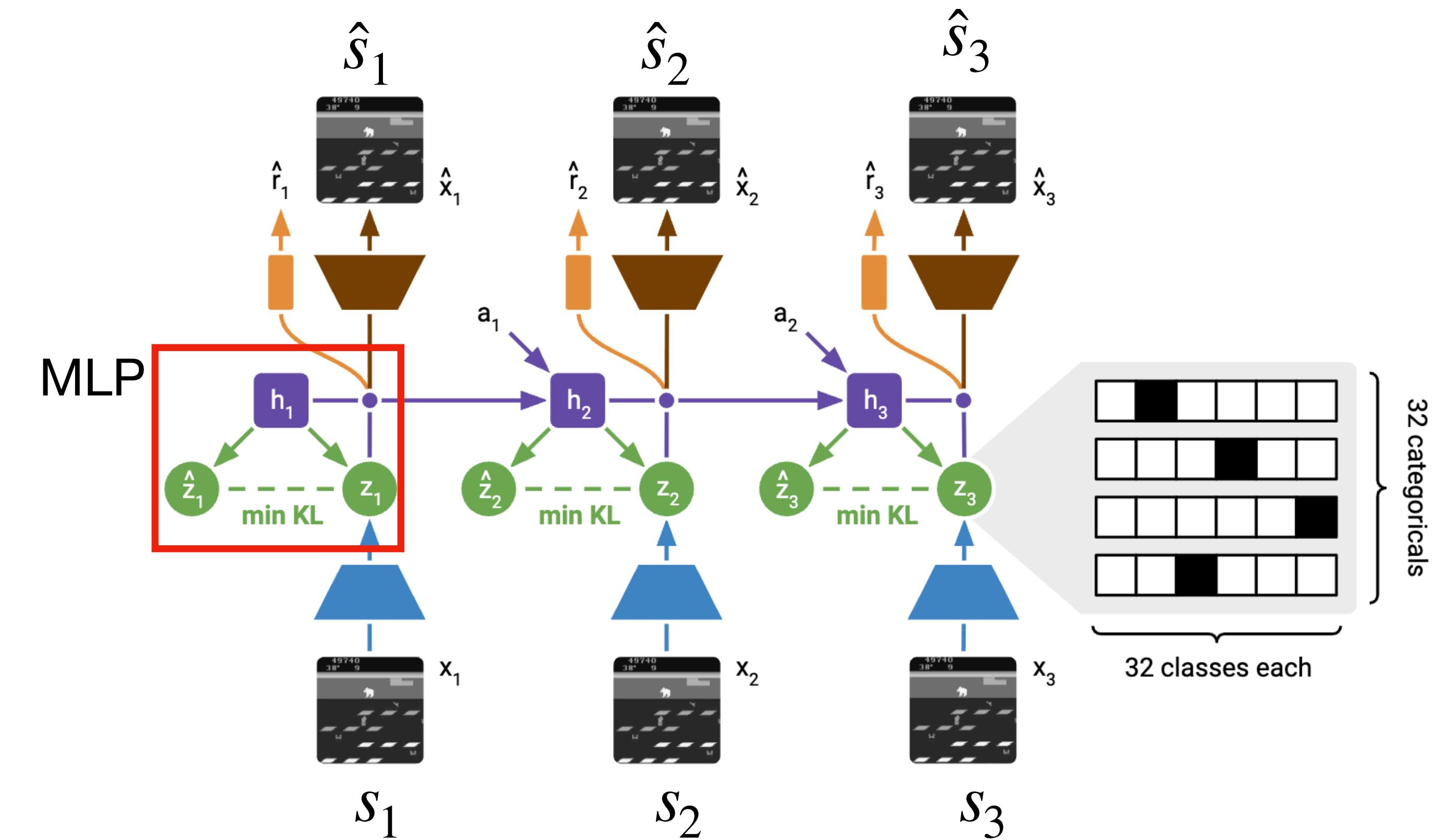


Building the RSSM

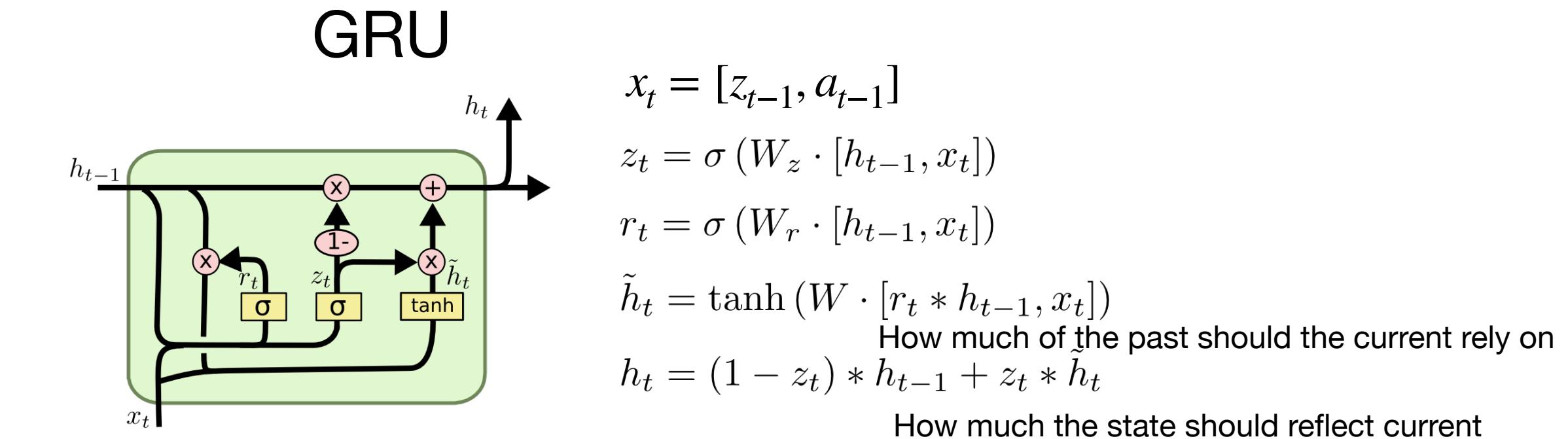
“Replacement for encoder during imagined rollouts”



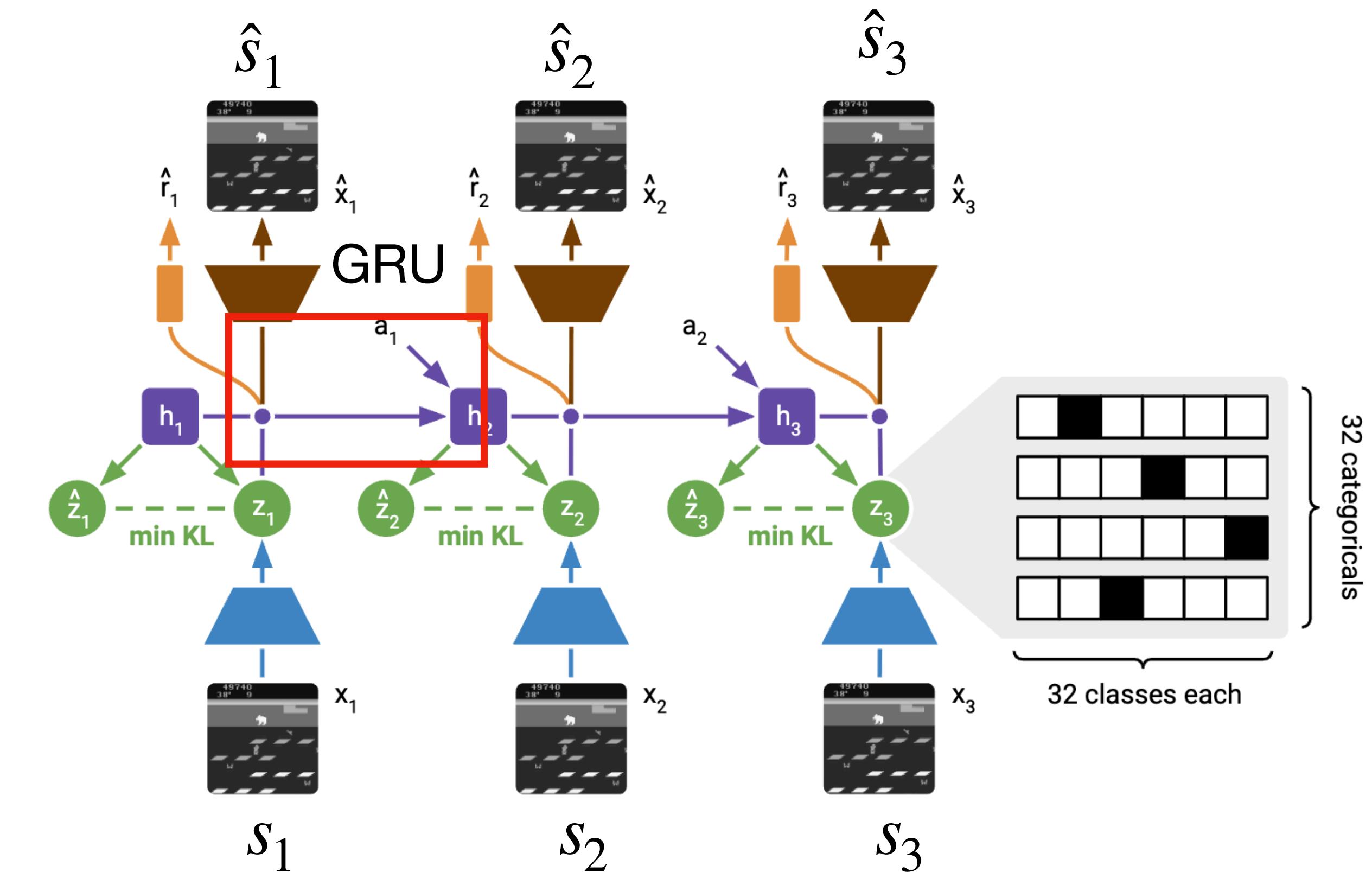
- Transition Predictor: $\hat{z}_t \sim p_\theta(\hat{z}_t | h_t)$



Building the RSSM



- Recurrent Model: $h_t = f_\theta(h_{t-1}, z_{t-1}, a_{t-1})$



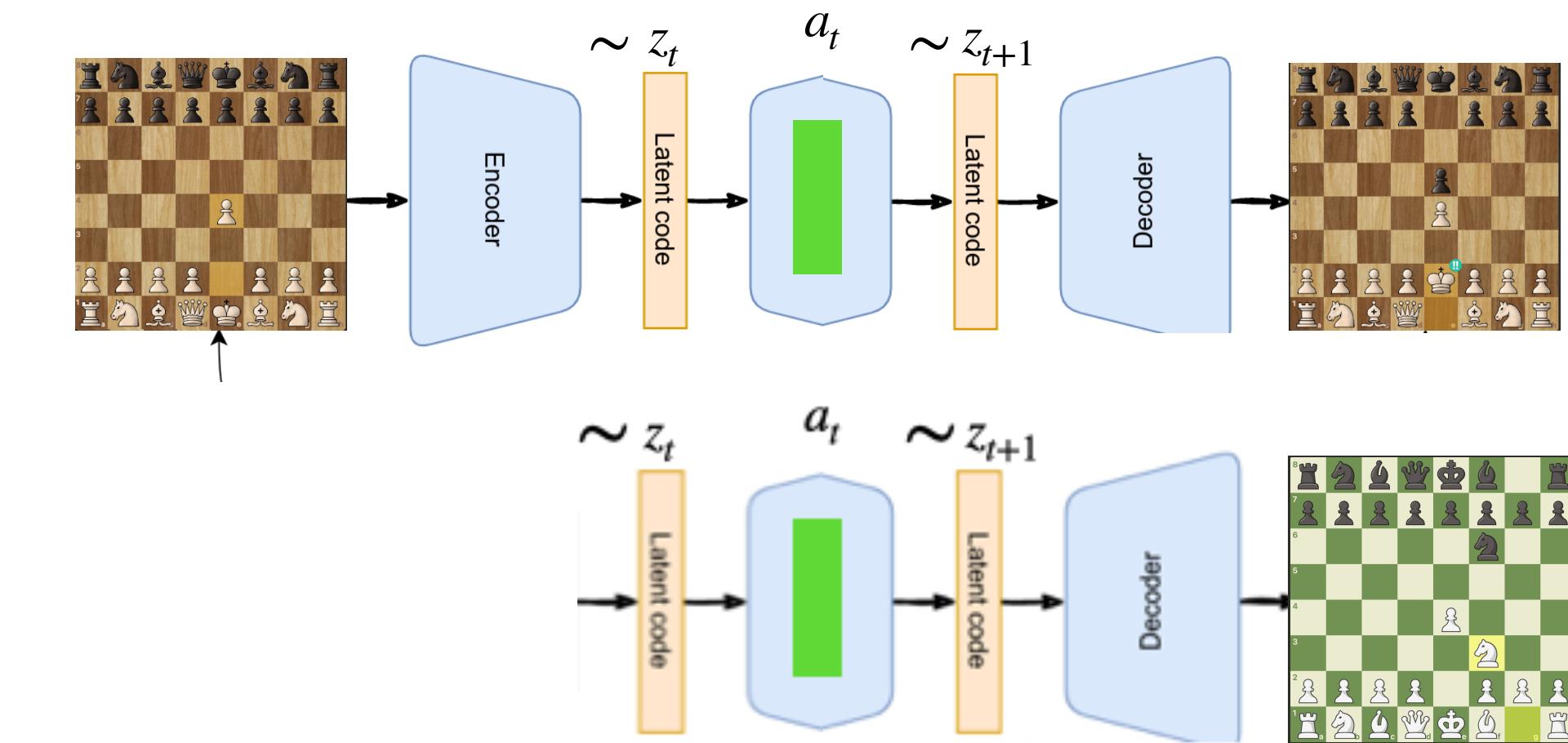
Building the generative RSSM

- Enforce z_t is from categorical distribution
- Encoder: $z_t \sim q_\theta(z_t | h_t, s_t)$
- Decoder: $\hat{s}_t \sim p_\theta(\hat{s}_t | h_t, z_t)$
- Reward Model: $r_t \sim p_\theta(r_t | h_t, z_t)$
- Transition Predictor: $\hat{z}_t \sim p_\theta(\hat{z}_t | h_t)$

Why?

- Capture uncertain dynamics (don't know opponents move or inherent randomness)
- Capture uncertainty in current state from partial observations
- Pragmatically: Variational Autoencoders are effective representative models
- See: MuZero for a deterministic latent space model

These are probabilistic models
Randomness from sampling latent z_t, \hat{z}_t



World Model Training

- Pull trajectories from buffer (of interactions with real environment)
- All components share parameters
- Minimize:
$$\mathcal{L}(\phi) \doteq \mathbb{E}_{q_\phi(z_{1:T}|a_{1:T}, s_{1:T})} \left[\sum_{t=1}^T \underbrace{-\beta_1 \ln p_\phi(s_t | h_t, z_t)}_{\text{state log loss}} + \underbrace{-\beta_2 \ln p_\phi(r_t | h_t, z_t)}_{\text{reward log loss}} + \underbrace{\beta_3 \text{KL} \left[q_\phi(z_t | h_t, s_t) \| p_\phi(\hat{z}_t | h_t) \right]}_{\text{KL loss}} \right].$$

KL Term does not work well as written

- Minimize: $\mathcal{L}(\phi) \doteq \mathbb{E}_{q_\phi(z_{1:T} | a_{1:T}, s_{1:T})} \left[\sum_{t=1}^T \underbrace{-\beta_1 \ln p_\phi(s_t | h_t, z_t)}_{\text{state log loss}} + \underbrace{-\beta_2 \ln p_\phi(r_t | h_t, z_t)}_{\text{reward log loss}} - \beta_3 \underbrace{\text{KL} \left[q_\phi(z_t | h_t, s_t) \| p_\phi(\hat{z}_t | h_t) \right]}_{\text{KL loss}} \right].$

- Problems from multi-objective optimization:
 - An accurate $p_\phi(\hat{z}_t | h_t)$ is more important than an accurate $q_\phi(z_t | h_t, s_t)$
 - Solution: Split the KL term with separate scaling coefficients ($B_{3_1} > B_{3_2}, B_{3_1} + B_{3_2} = 1,$)
 - $\text{KL} \left[q_\phi(z_t | h_t, s_t) \| p_\phi(\hat{z}_t | h_t) \right] = B_{3_1} \text{KL} \left[\text{sg}[q_\phi(z_t | h_t, s_t)] \| p_\phi(\hat{z}_t | h_t) \right] + B_{3_2} \text{KL} \left[q_\phi(z_t | h_t, s_t) \| \text{sg}[p_\phi(\hat{z}_t | h_t)] \right]$
 - “sg” is the stop gradient operator. “detach” in PyTorch.

KL Term does not work well as written

- Minimize: $\mathcal{L}(\phi) \doteq \mathbb{E}_{q_\phi(z_{1:T} | a_{1:T}, s_{1:T})} \left[\sum_{t=1}^T \underbrace{-\beta_1 \ln p_\phi(s_t | h_t, z_t)}_{\text{state log loss}} + \underbrace{-\beta_2 \ln p_\phi(r_t | h_t, z_t)}_{\text{reward log loss}} - \beta_3 \underbrace{\text{KL} \left[q_\phi(z_t | h_t, s_t) \| p_\phi(\hat{z}_t | h_t) \right]}_{\text{KL loss}} \right].$
- Problems from multi-objective optimization:
 - Degenerate case (z_t easy to predict) not informative.
 - Focus should be on state and reward recovery once predicted embedding is “good enough”
 - Solution: clip minimum, allow “free unoptimized bits of difference between them”
 - $B_{3_1} \min[c, \text{KL} \left[\text{sg}[q_\phi(z_t | h_t, s_t)] \| p_\phi(\hat{z}_t | h_t) \right]] + B_{3_2} \min[c, \text{KL} \left[q_\phi(z_t | h_t, s_t) \| \text{sg}[p_\phi(\hat{z}_t | h_t)] \right]]$

Putting it all together

- Collect transitions from real env with epsilon-greedy
- Train predictive generative world model on real transitions
- Train policy and value models on generated transitions
- Repeat

Results (DreamerV3)

- Surprisingly sample efficient (in terms of interactions with real environment)

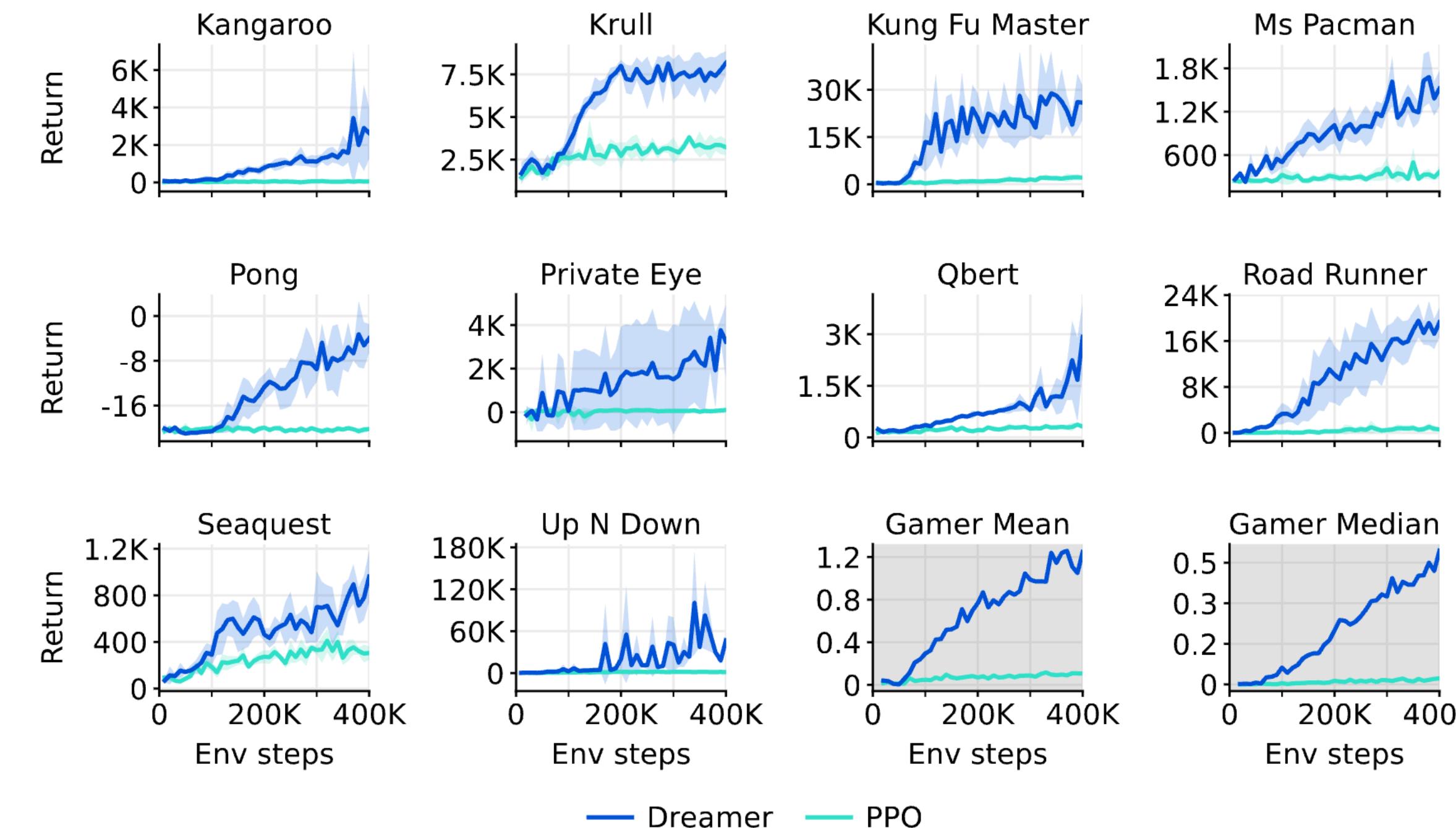


Figure 13: Atari100k learning curves.

Results (DreamerV3)

- Surprisingly sample efficient (in terms of interactions with real environment)

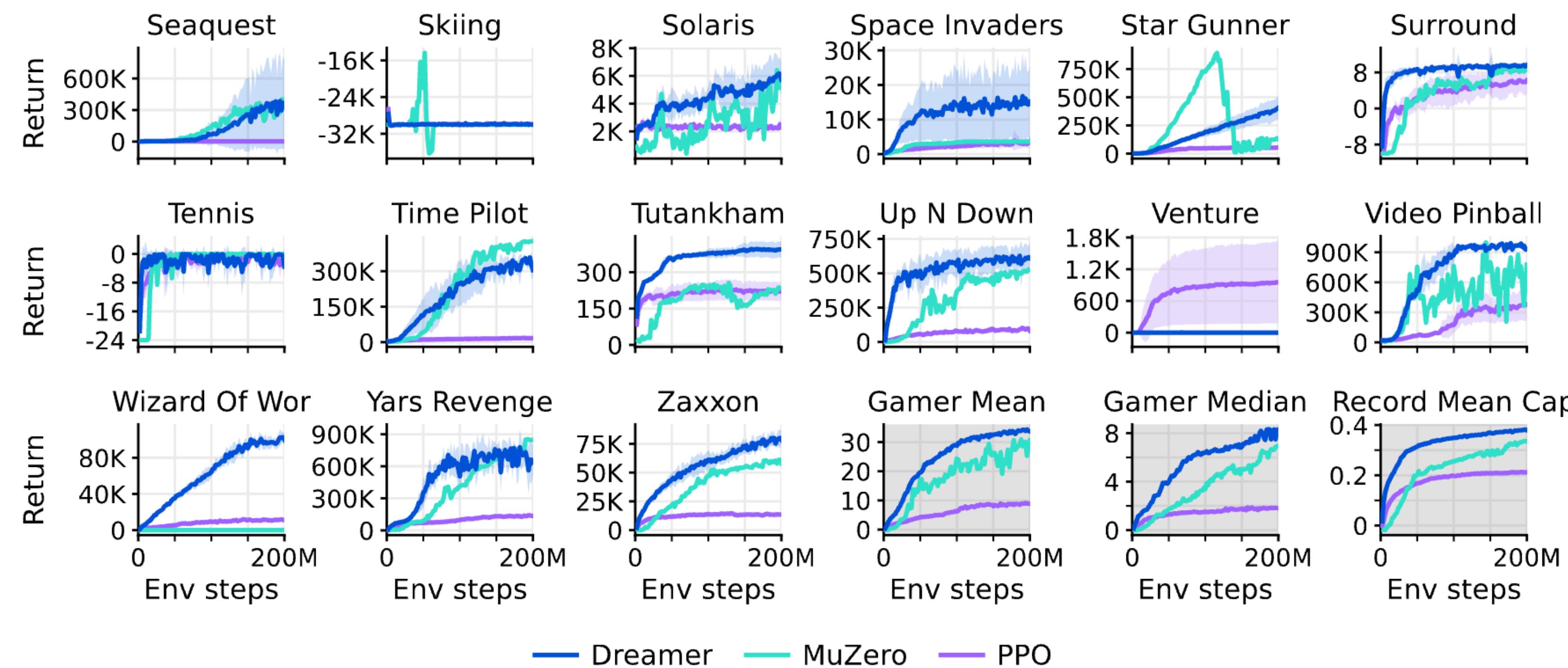


Figure 10: Atari learning curves.

Results (DreamerV3)

each room has a different layout, treasure, and enemies, and some rooms have special features, such as moving walls

- Surprisingly sample efficient (in terms of interactions with real environment)

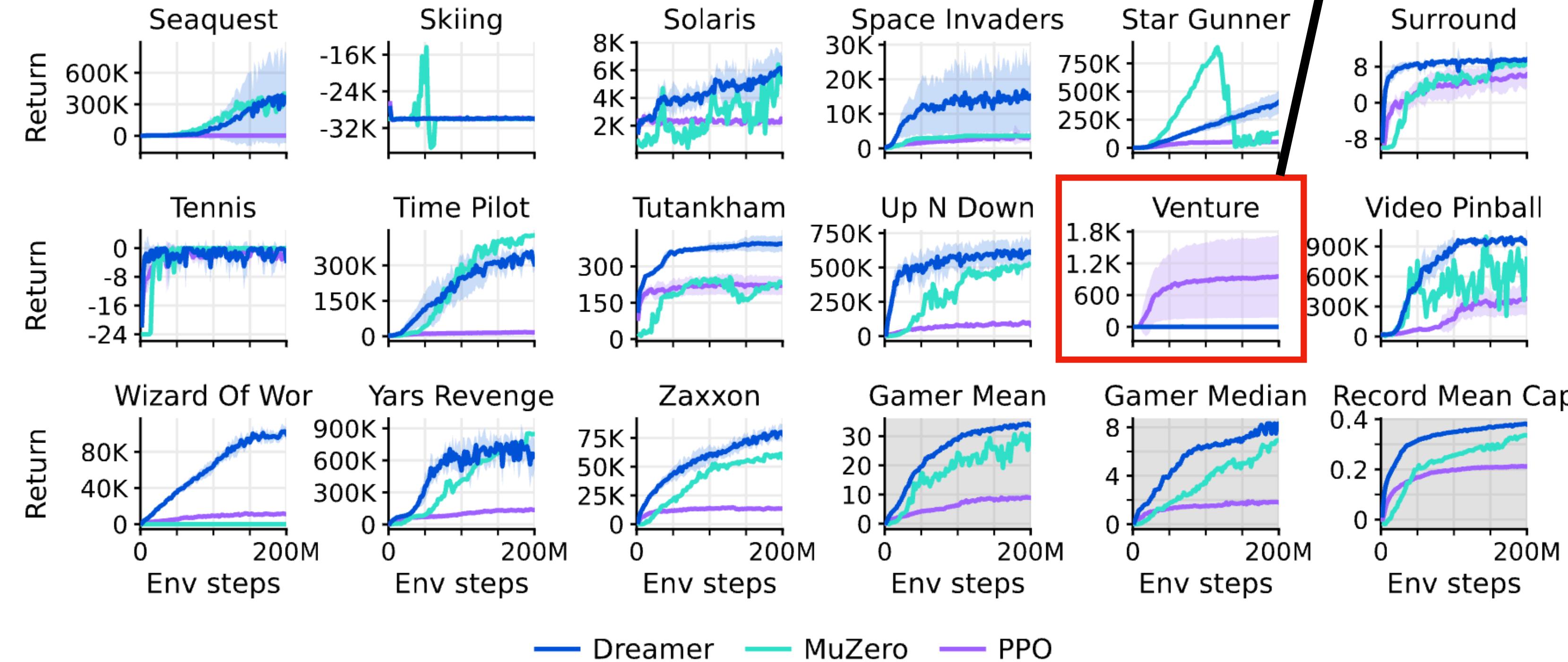
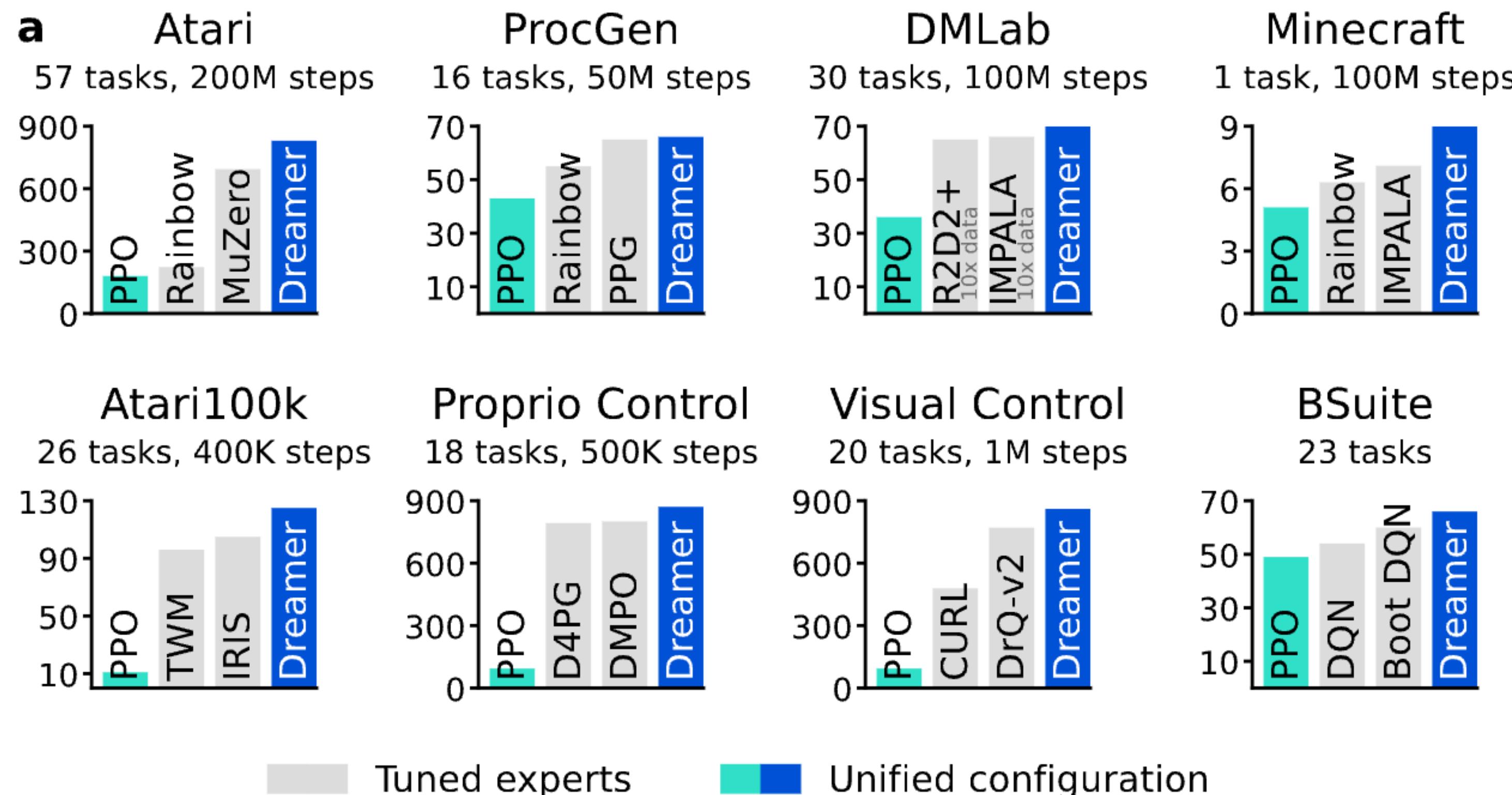


Figure 10: Atari learning curves.

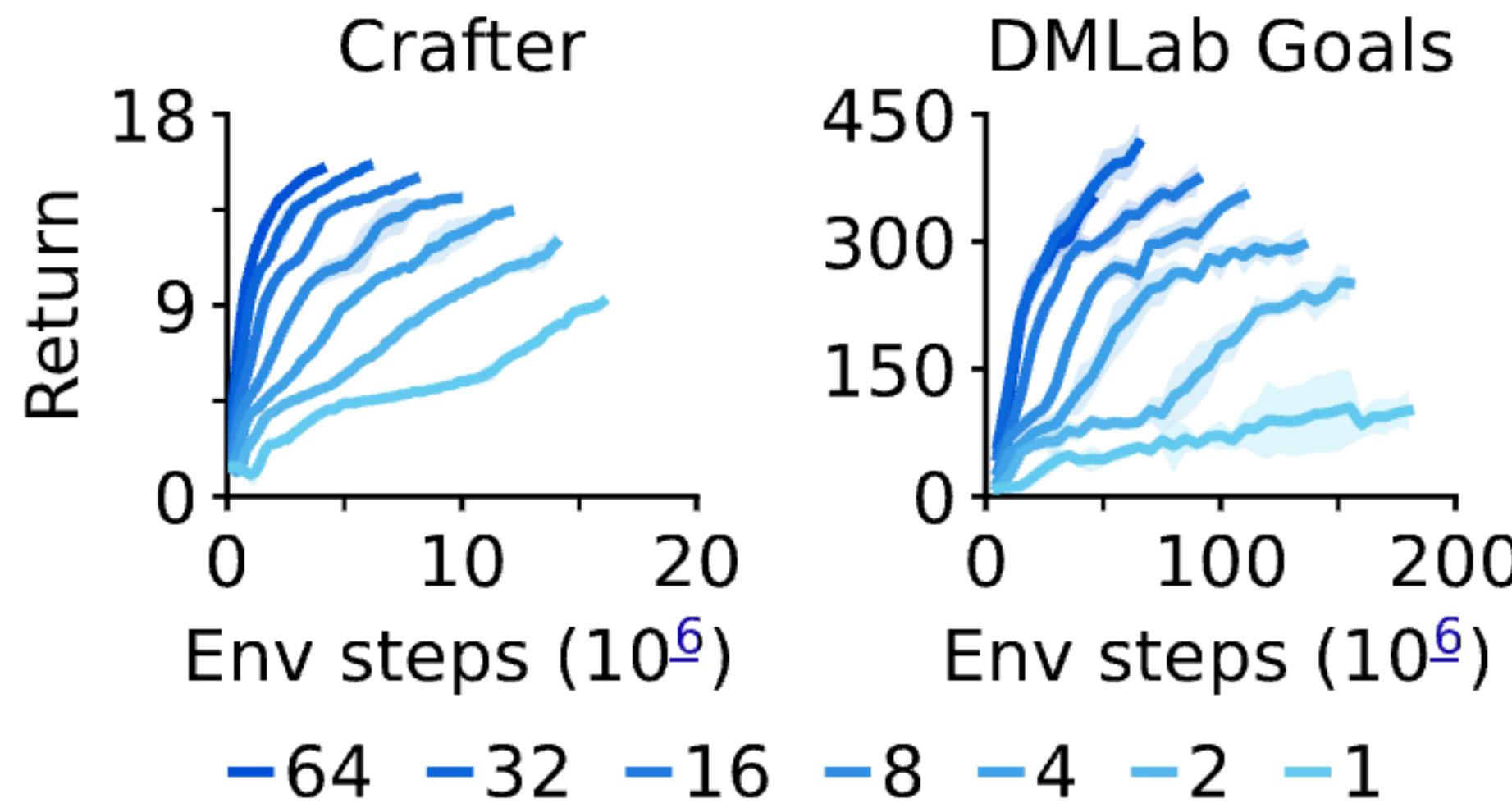
Results (DreamerV3)

- Claimed effective across a range of tasks without hyperparameter search

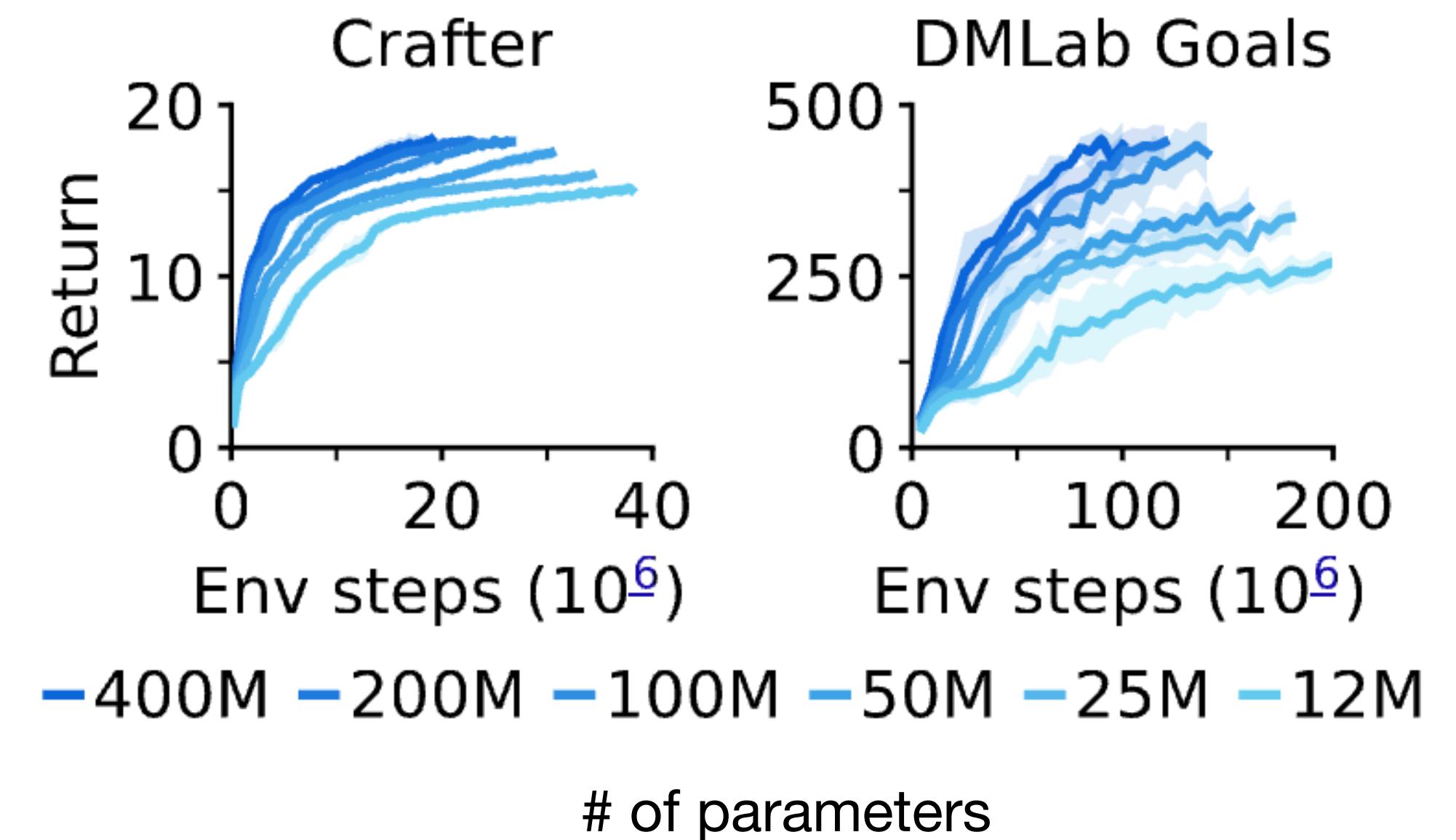


Results (DreamerV3)

d Replay scaling



c Model size scaling



- Imagination vs Buffer ratio
 - X-axis: Real steps to train prediction
 - Lines: Ratio of imagined steps to train policy per real step to train predictor generative model
 - Without a reference A2C or PPO on this plot it is hard to compare sample efficiency to not using a generative model

- Models are **LARGE**
- But performance increases as parameters increase for same setup

Crafter is a synthetic 2D Minecraft
DMLab is 3D Maze navigation