

Florian Vahl

08.12.2022

---

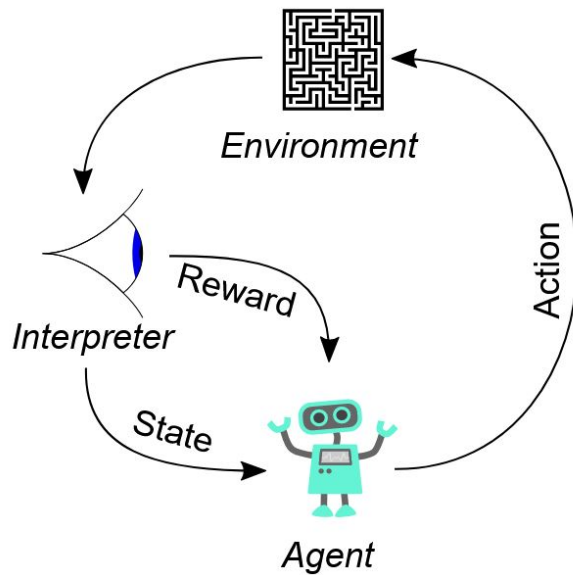
# Reinforcement learning with Dreamer / Dreaming V1 / V2

# Outline

- Fundamentals
  - Markov Decision Process
  - Reinforcement Learning
- Approaches
  - Dreamer V1
  - Dreamer V2
  - Dreaming V1
  - Dreaming V2
- Results
- Discussion

# Fundamentals

# Markov Decision Process



# Markov Decision Process

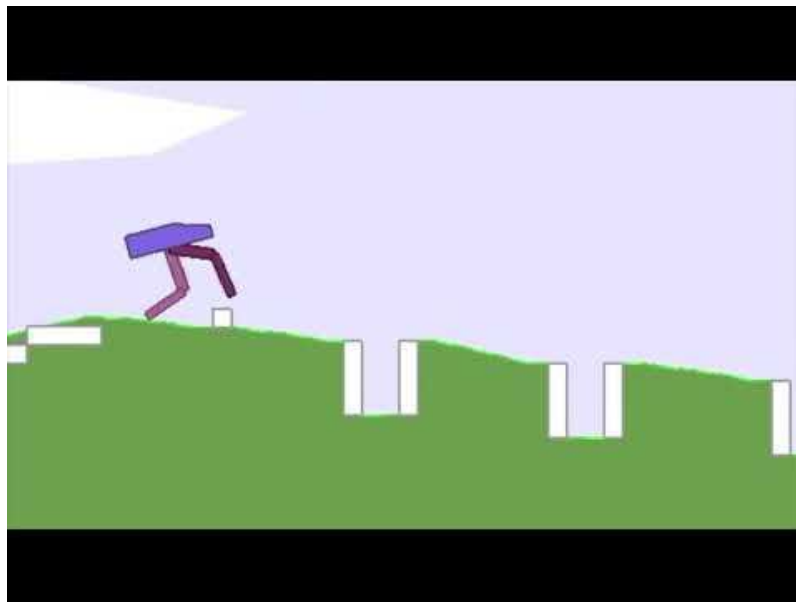
An MDP is a tuple  $\langle S, A, P, R \rangle$  where:

- $S$  is a set of states.
- $A$  is a set of actions.
- $P : S \times A \times S \rightarrow [0, 1]$  is the state transition probability function.
- $R : S \times A \rightarrow R$  is the reward function.

# Reinforcement Learning

- Type of machine learning that involves agents learning based on a reward
- No ground truth labels
- Optimizes policy to maximize the reward signal
- Applied in many domains including robotic control and game playing
- Typically formulated using the Markov Decision Process

# Reinforcement Learning - Environments



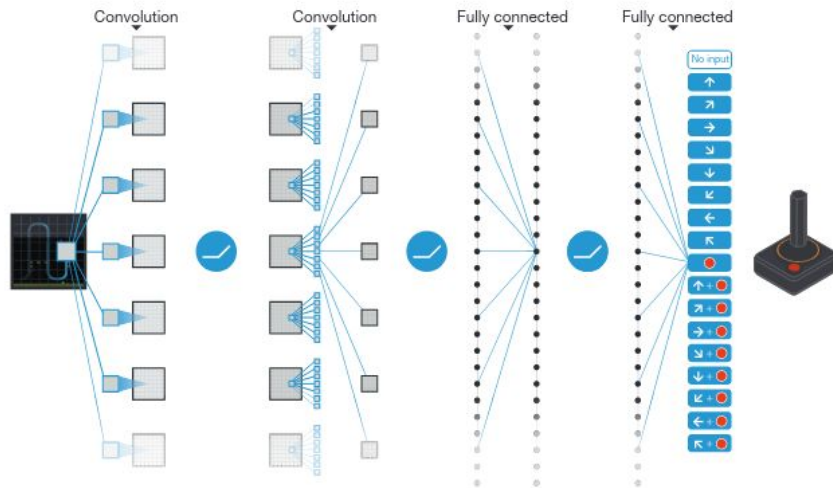
<https://www.youtube.com/watch?v=SyonPJc8hMw>

# Reinforcement Learning - Environments





# Reinforcement Learning - Deep Q-Learning (Example)

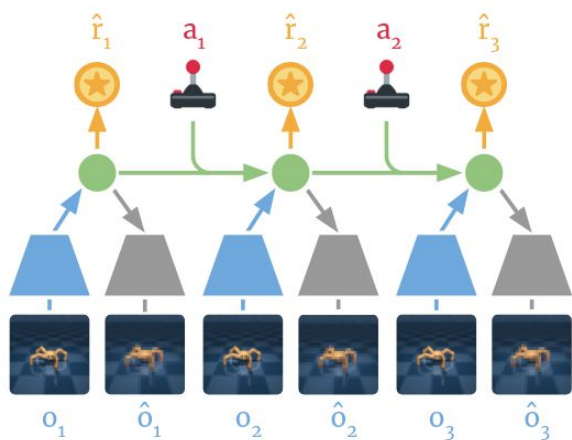


# Approaches

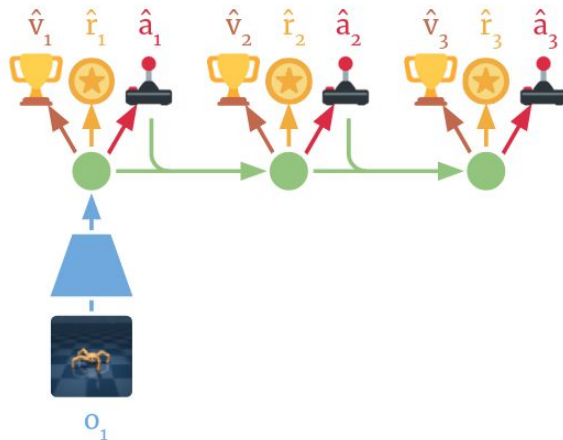
## Dreamer V1

- Paper: “Dream to Control: Learning Behaviors by Latent Imagination” by Hafner et al. from Google Brain at ICLR 2020
- Model based reinforcement learning algorithm
- Learns world model
- Solves long-horizon tasks purely from image base observations
- Imagines novel trajectories in the world models latent space
- Policy is trained using the imagined trajectories

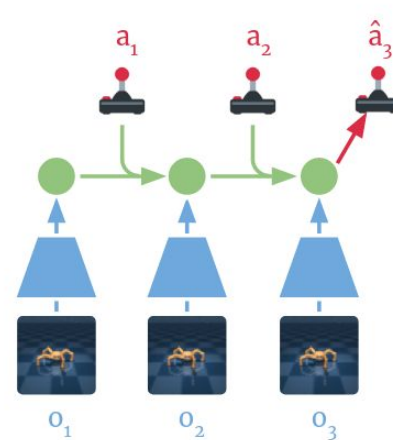
# Dreamer V1 – Architecture



(a) Learn dynamics from experience



(b) Learn behavior in imagination



(c) Act in the environment

Dream to Control: Learning Behaviors by Latent Imagination (2020)

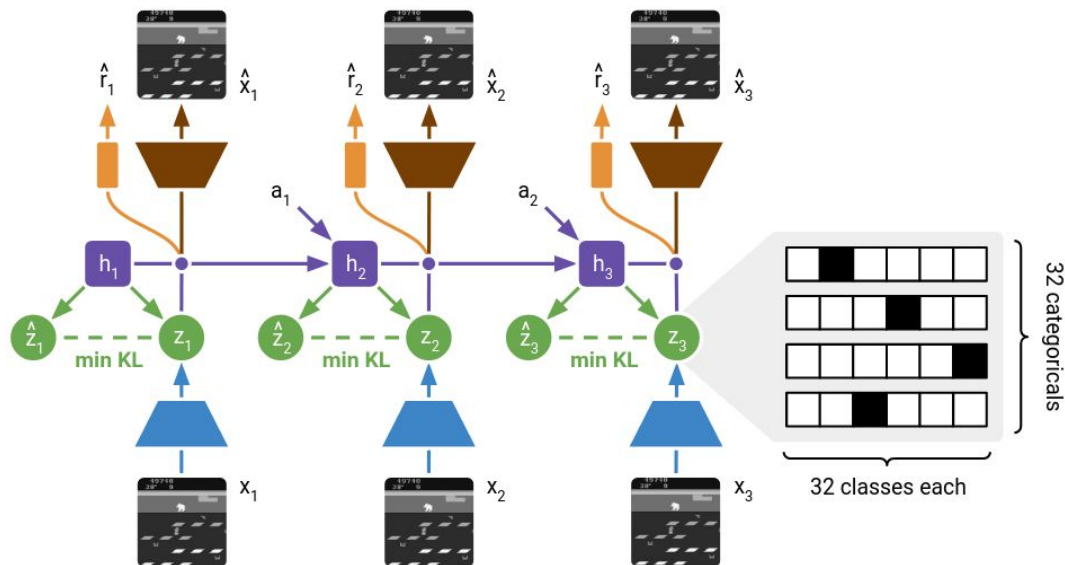
## Dreamer V1 – World Model

- Representation model  $p_{\theta}(s_t \mid s_{t-1}, a_{t-1}, o_t)$
- Observation model  $q_{\theta}(o_t \mid s_t)$
- Reward model  $q_{\theta}(r_t \mid s_t)$
- Transition model  $q_{\theta}(s_t \mid s_{t-1}, a_{t-1})$ .
  - Predicts next latent state
  - Outputs tanh-transformed Gaussian distribution
- Frozen while learning behaviors

## Dreamer V2

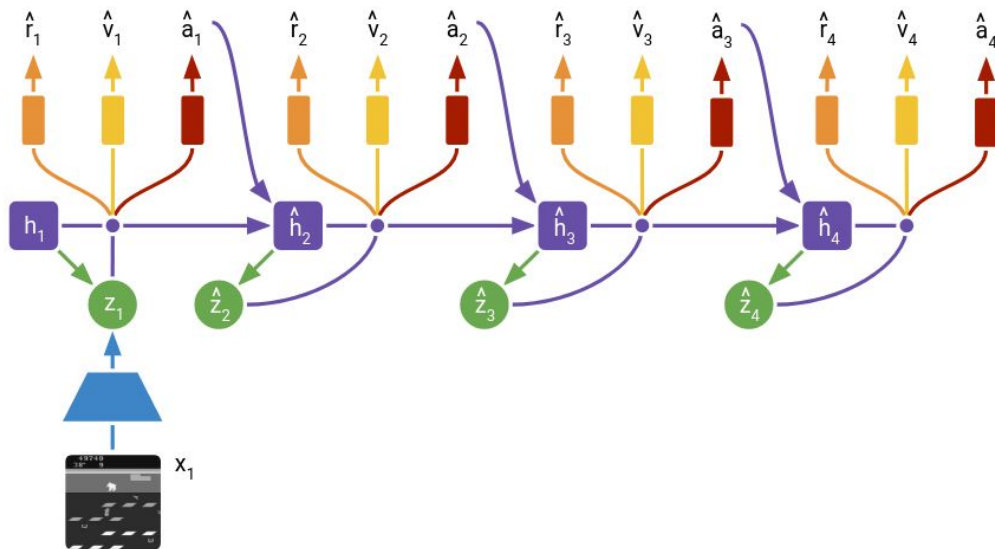
- Paper: “Mastering Atari with Discrete World Models” by Hafner et al. from Google Brain at ICLR 2021
- Similar to Dreamer V1
- Transition model uses categorical variables instead of Gaussians
- Performed well in image based Atari tasks

# Dreamer V2 – Architecture



Mastering Atari with Discrete World Models (2022)

# Dreamer V2 – Latent Imagination



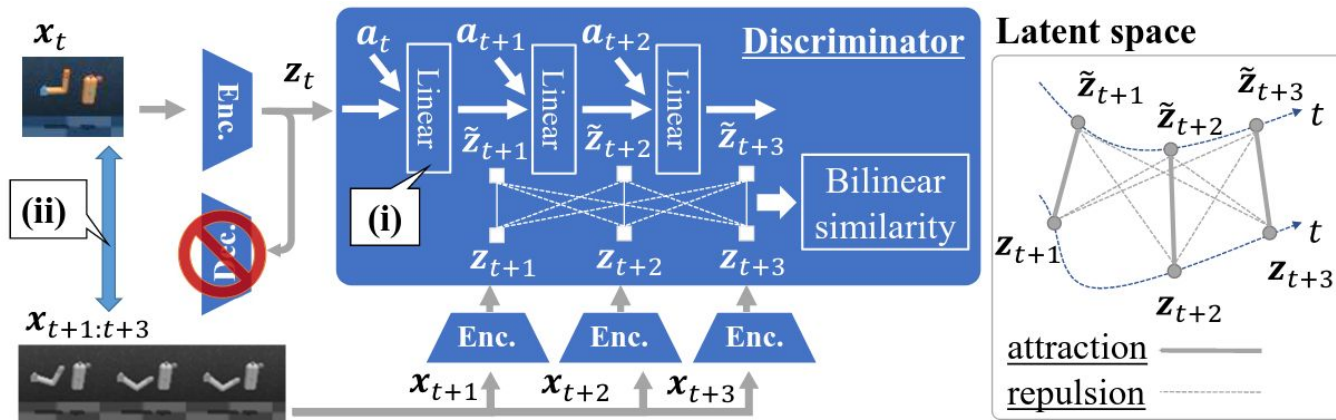
Mastering Atari with Discrete World Models (2022)



## Dreaming V1

- Paper: “Dreaming: Model-based Reinforcement Learning by Latent Imagination without Reconstruction” by Okada et al. from Panasonic Corporation in ICRA 2021
- Similar to Dreamer V1
- World model is trained via contrastive learning
- No generative decoder / observation model
- Tries to avoid object vanishing

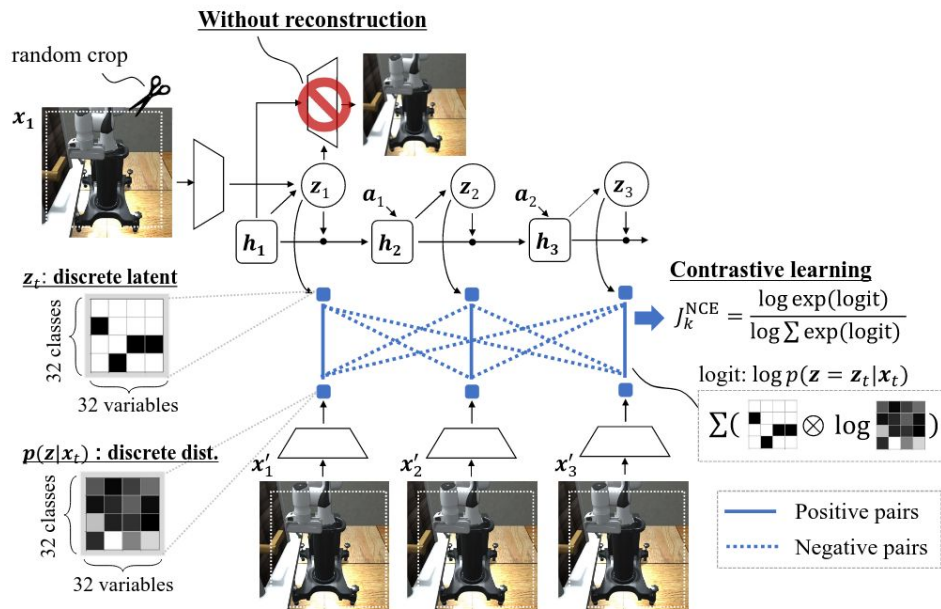
# Dreaming V1 – Architecture



## Dreaming V2

- Paper: “DreamingV2: Reinforcement Learning with Discrete World Models without Reconstruction” by Okada et al. from Panasonic Corporation on arXiv (2022)
- Applies the concept concepts introduced in Dreamer V2 to Dreaming V1
- Also decoder free

# Dreaming V2 – Architecture

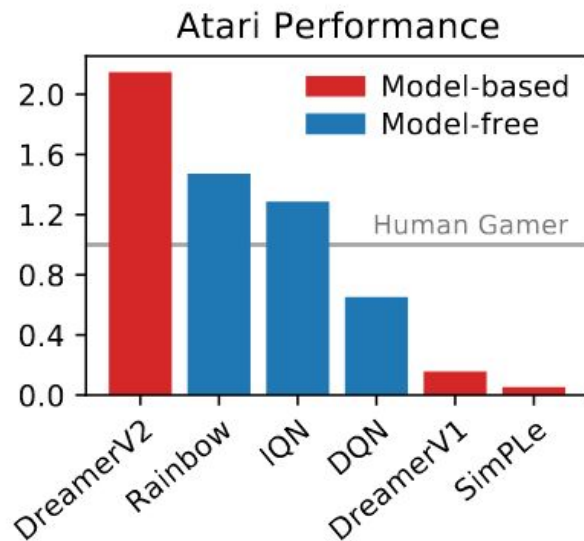


# Results

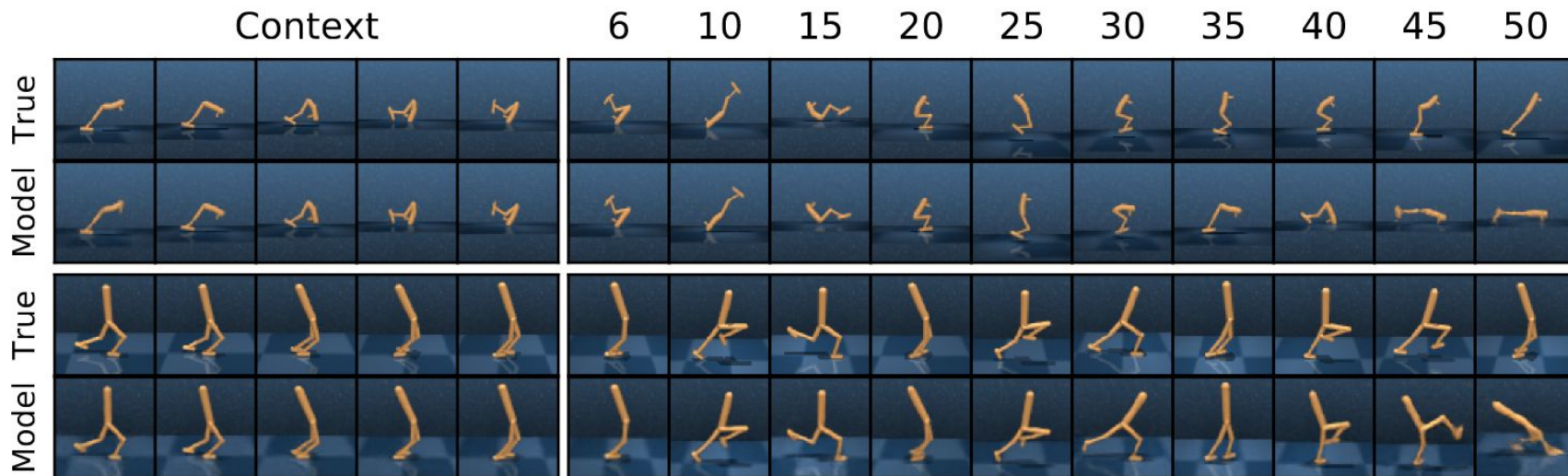
# Results

	DreamingV2 (ours)	DreamerV2 [6]	Dreaming [7]	Dreamer [5]
Discrete latent	✓	✓		
Reconstruction free	✓		✓	
<b>(A) 3D Rotot-arm tasks</b>				
UR5-reach	<u>776</u> ±194	<u>704</u> ±222	<u>752</u> ±1178	<u>701</u> ±223
Reach-duplo	<u>199</u> ±43	149±62	145±61	5±11
Lift	<u>327</u> ±150	165±126	174±107	134±46
Door	<u>383</u> ±143	190±126	319±173	154±32
PegInHole	<u>436</u> ±26	<u>376</u> ±59	353±50	354±47
<b>(B) 2D Manipulation tasks</b>				
Reacher-hard	<u>598</u> ±447	175±340	<u>743</u> ±346	247±392
Finger-turn-hard	484±434	<u>600</u> ±417	<u>858</u> ±210	533±426
Reacher-easy	<u>924</u> ±210	923±215	<u>947</u> ±100	658±429
Finger-turn-easy	434±469	498±469	<u>842</u> ±286	<u>665</u> ±430
<b>(C) Difficult pole-swingup tasks</b>				
Acrobot-swingup	<u>470</u> ±129	309±131	359±111	<u>382</u> ±147
Cartpole-two-poles	<u>308</u> ±55	248±103	<u>273</u> ±53	256±65
<b>(D) 3D locomotion robot tasks</b>				
Quadruped-walk	<u>492</u> ±127	350±89	<u>379</u> ±189	242±120
Quadruped-run	<u>385</u> ±91	<u>352</u> ±68	339±128	269±114
<b>(E) 2D locomotion tasks</b>				
Cheetah-run	768±24	<u>811</u> ±75	542±132	<u>776</u> ±120
Walker-walk	857±115	<u>951</u> ±28	518±76	<u>906</u> ±70

# Results



# Results

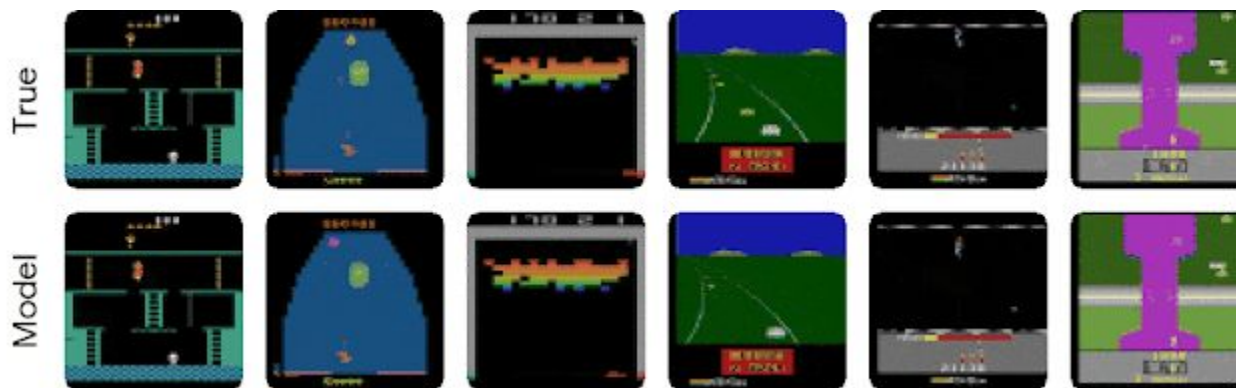




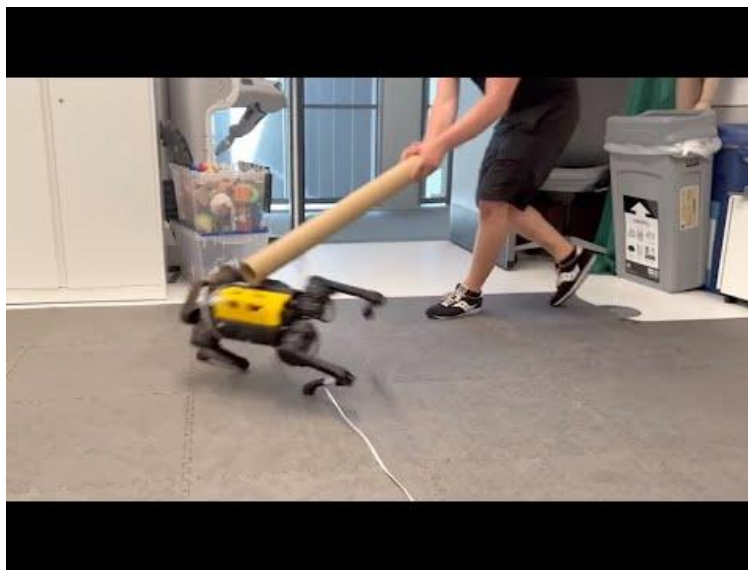
## Results – Video



# Results – Video



## Results – Real World Reinforcement Learning



## Discussion

- Sampling efficiency can be improved using latent imagination
- Object vanishing is still a problem even in Dreaming
- Atari performance only for “single GPU” setups
- Categorical stochastic representations improve performance
- Good fit for situations where exploration is costly (e.g. real-world reinforcement learning)

**Questions?**