



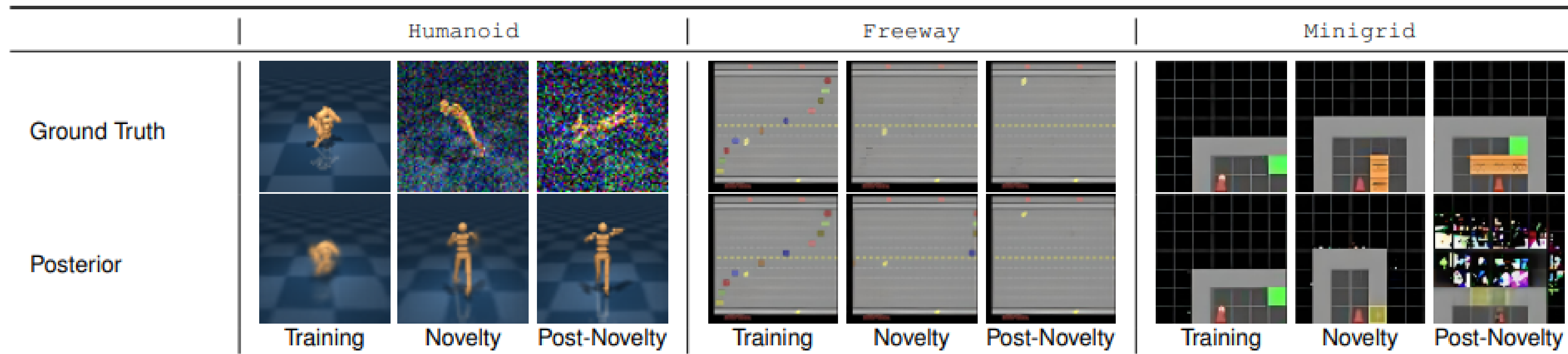
Novelty Detection in Reinforcement Learning with World Models

Geigh Zollicoffer, Kenneth Eaton*, Jonathan C Balloch, Julia Kim, Wei Zhou, Robert Wright*, Mark Riedl

Georgia Institute of Technology, Georgia Tech Research Institute*

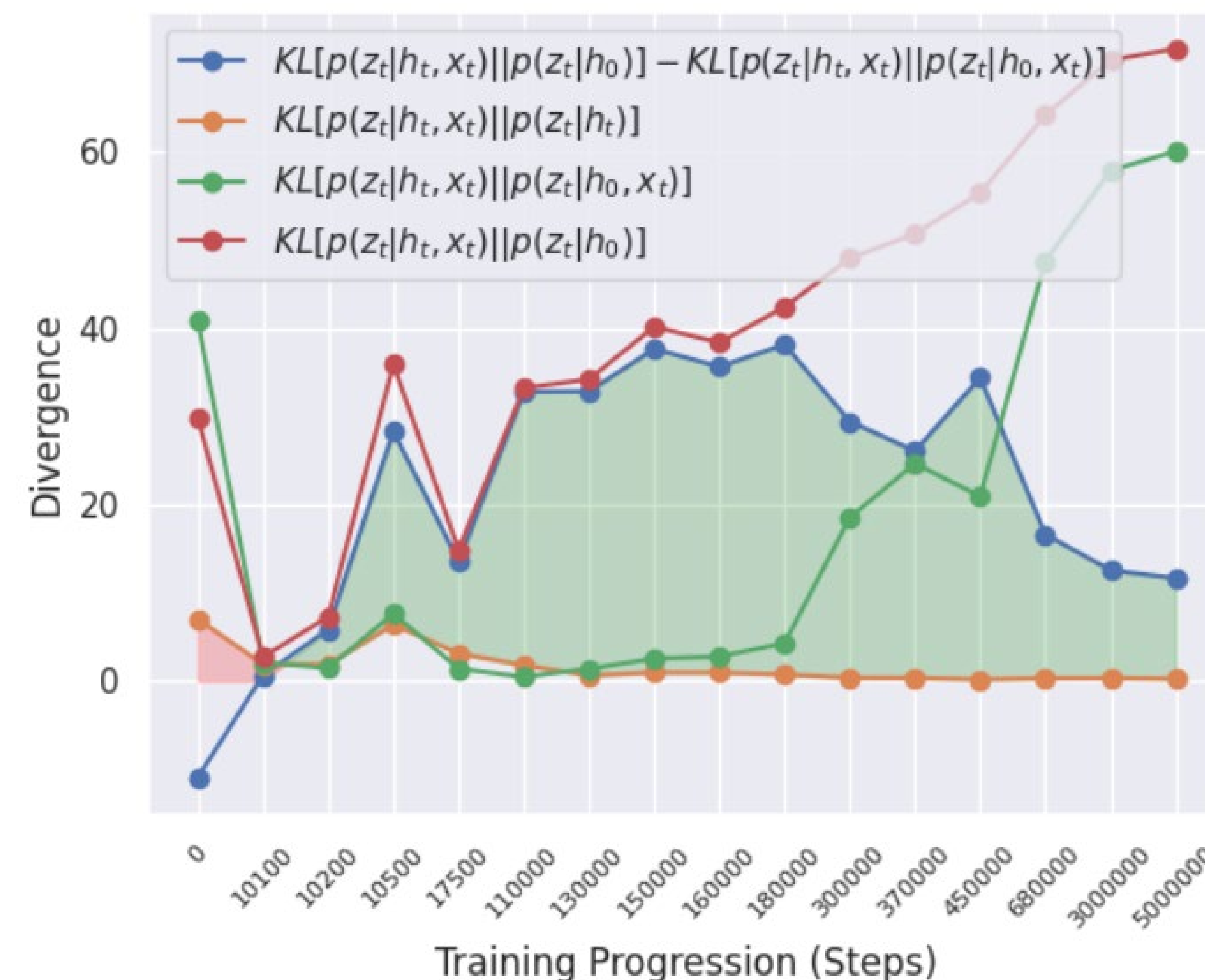
1. Problem Statement

Reinforcement learning (RL) using world models has found significant recent successes. However, when a sudden change to world mechanics or properties occurs then agent performance and reliability can dramatically decline.



Our technique calculates a novelty threshold bound without additional hyperparameters by considering how much the actual world observation deviates from the distribution of world observations that the agent predicts it will encounter.

2. Latent Based Detection



World Model Architecture:

- Recurrent model: $h_t = f_\phi(h_{t-1}, z_{t-1}, a_{t-1})$
- Representation model: $p_\phi(z_t|h_t, x_t)$
- Transition prediction model: $p_\phi(\hat{z}_t|h_t)$
- Image prediction model: $p_\phi(\hat{x}_t|h_t, z_t)$
- Reward prediction model: $p_\phi(\hat{r}_t|h_t, z_t)$
- Discount prediction model: $p_\phi(\hat{\gamma}_t|h_t, z_t)$

$$\mathcal{L}(\phi) = \mathbb{E}_{p_\phi(z|a,x)} \left[\sum_t^T -\ln p_\phi(x_t|h_t, z_t) - \ln p_\phi(r_t|h_t, z_t) - \ln p_\phi(\gamma_t|h_t, z_t) + \beta KL[p_\phi(z_t|h_t, x_t)||p_\phi(z_t|h_0)] \right]$$

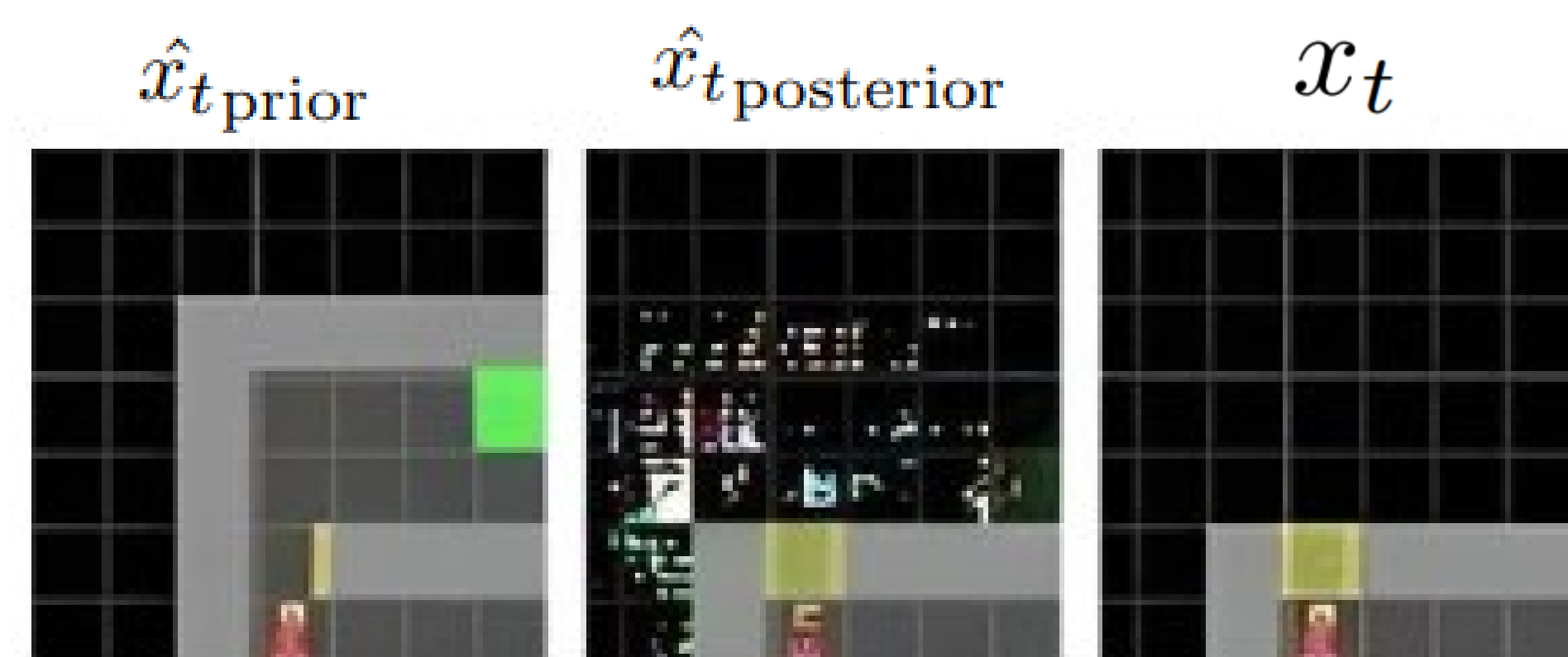
1. Measure the cross-entropy score comparison and derive in terms of KL to use world model components:

$$H(p_\phi(z_t|h_t, x_t), p_\phi(z_t|h_0)) - H(p_\phi(z_t|h_t, x_t), p_\phi(z_t|h_0, x_t)) = (KL[p_\phi(z_t|h_t, x_t)||p_\phi(z_t|h_0)] + H(p_\phi(z_t|h_t, x_t))) - KL[p_\phi(z_t|h_t, x_t)||p_\phi(z_t|h_0, x_t)] - H(p_\phi(z_t|h_t, x_t)) = KL[p_\phi(z_t|h_t, x_t)||p_\phi(z_t|h_0)] - KL[p_\phi(z_t|h_t, x_t)||p_\phi(z_t|h_0, x_t)]$$

2. Derive a bound using the minimization objective of the world model as an anomaly score:

$$KL[p_\phi(z_t|h_t, x_t)||p_\phi(z_t|h_t)] \leq KL[p_\phi(z_t|h_t, x_t)||p_\phi(z_t|h_0)] - KL[p_\phi(z_t|h_t, x_t)||p_\phi(z_t|h_0, x_t)]$$

3. Observation based Detection



Derive a bound from difference in pixel values compared to tuned hyperparameter λ :

$$\frac{\sum_i^N |\hat{x}_{t,prior}^i - \hat{x}_{t,posterior}^i|}{N} \leq \lambda$$

Note, the prior and posterior images are generated without providing the ground truth, x_t , as input

4. Experimental Evaluation

Task: Detect a novel transition as soon as the agent experiences it, while minimizing false positives.

Techniques: KL Bound (Latent Based), RIQN (Ensemble Based), PP-MARE (Observation Based).

Results: Agents can utilize world models to detect novel transitions at a more effective rate than previous methods.

Method	False Positive Rate			
	Boxing	Kangaroo	Freeway	SeaQuest
KL Bound	$\leq 10^{-2}$	$\leq 10^{-2}$	$\leq 10^{-2}$	$\leq 10^{-2}$
PP-Mare	.04	$\leq 10^{-2}$	$\leq 10^{-2}$	$\leq 10^{-2}$
RIQN	.17	.39	.24	.43

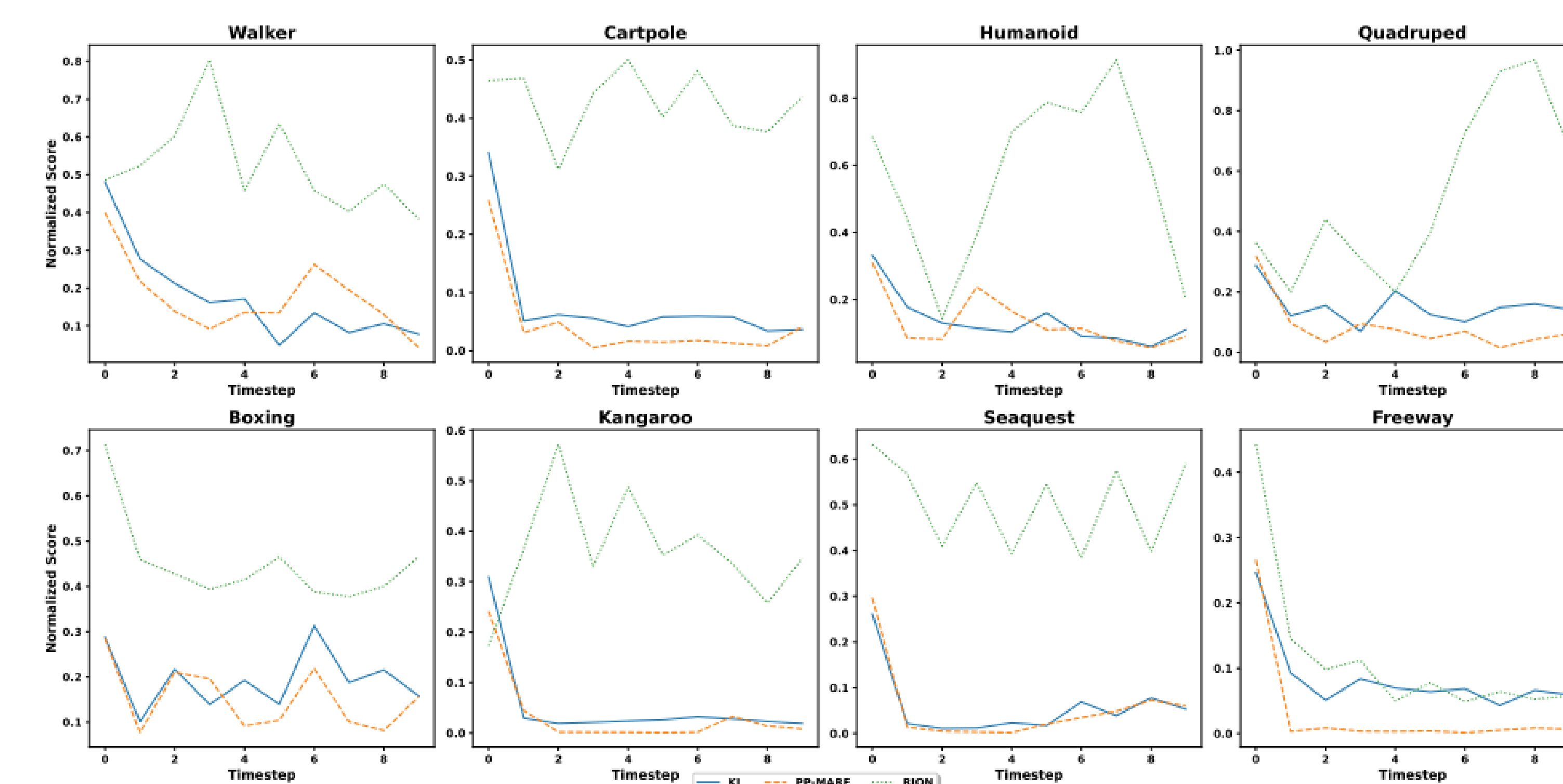
Method	Average False Positive Rate		Inference Run-time Speedup	
	DMC	Atari	DMC	Atari
RIQN	.275	.308	$\times 1$	$\times 1$
PP-Mare	$\leq 10^{-2}$	$\leq 10^{-2}$	$\times 1.16 \cdot 10^2$	$\times 4.45 \cdot 10^1$
KL bound	$\leq 10^{-2}$	$\leq 10^{-2}$	$\times 1.34 \cdot 10^3$	$\times 5.12 \cdot 10^2$

World Model based detection improves the speed and minimizes the average false positive rate.

	ADE↓	AUC↑	ADE↓	AUC↑	ADE↓	AUC↑
Boxing						
OneArm			BodySwitch		Doppelganger	
KL Bound	52.6	.708	$\leq 10^{-2}$	$\geq .99$	$\leq 10^{-2}$	$\geq .99$
PP-Mare (2)	22.5	.605	6.30	.915	5.4	.862
RIQN ($10^{-5}, 10^{-7}$)	347.5	.505	509.3	.380	103.1	.401
Kangaroo						
Floorswap			Difficulty+		DisableMonkey	
KL Bound	9.9	.787	$\leq 10^{-2}$	$\geq .99$.960	$\geq .99$
PP-Mare (.5)	85.2	.281	42.5	$\geq .99$	41.3	.937
RIQN ($10^{-2}, 10^{-2}$)	166.3	.541	94.2	$\geq .99$	93.1	.710
Freeway						
InvisibleCars			ColorCars		FrozenCars	
KL Bound	$\leq 10^{-2}$	$\geq .99$	$\leq 10^{-2}$	$\geq .99$	$\leq 10^{-2}$.985
PP-Mare (.5)	2.33	$\geq .99$	1.84	$\geq .99$	2.60	.931
RIQN ($10^{-7}, 10^{-9}$)	2.87	.938	2.62	$\geq .99$.502	.980
SeaQuest						
DisableEnemy			Gravity		UnlimitedOxygen	
KL Bound	.202	.962	45.8	.949	$\leq 10^{-2}$	$\geq .99$
PP-Mare (7)	111.6	.678	24.1	.882	3.73	.938
RIQN ($10^{-2}, 10^{-3}$)	45.3	.272	157.3	.390	16.0	.701

Left: Evaluation scores of AUC and Average Delay Error (ADE) on Minigrid and Atari environments

Right: Evaluation scores of AUC and Average Delay Error (ADE) on DeepMind Control environments.



The normalized anomaly score trend of each detection method as the episode progresses in the nominal environment.

