

# MuDreamer: Learning Control-Focused World Models without Pixel Reconstruction

Withdraw From ICLR 2024

November 10, 2025

# Problem Setup in DreamerV3

## **DreamerV3 trains latents by reconstructing pixels.**

- ▶ Pixel reconstruction shapes the representation toward *all* visual details (incl. backgrounds).
- ▶ Under strong visual distractors, this can dilute task-relevant features.
- ▶ Reconstruction is compute-heavy and not strictly necessary for control.

**Goal:** Focus the learned latent state (  $\mathbf{s}_t = \{h_t, z_t\}$  ) on *control-relevant* information (rewards, continuation, value, action effects), not textures.

# MuDreamer's Core Idea

**Drop pixel reconstruction as a shaping signal.** Learn a latent world model by predicting *task-relevant* quantities only:

- ▶ Reward ( $\hat{r}_t$ ) and Continue ( $\hat{c}_t$ )
- ▶ **Value** (distributional, discretized ( $\lambda$ )-return)
- ▶ **Past Action** (inverse dynamics)

Optional: keep a decoder for visualization **with stop-grad**, so it does not influence latents.

**Why:** Dense supervision each step, but aligned with *control*.

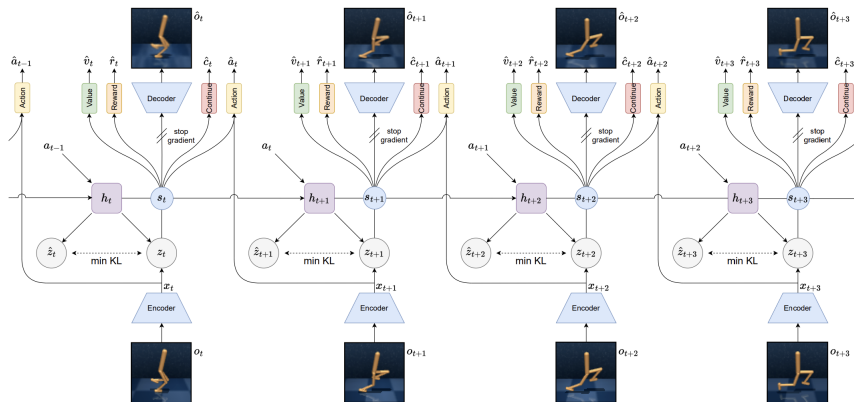
# Architecture at a Glance

## Per-timestep (strictly sequential)

1. Update memory:  $(h_t = f_\phi(h_{t-1}, z_{t-1}, a_{t-1}))$
2. Prior over state:  $(\hat{z}_t \sim p_\phi(z \mid h_t))$
3. Encode image:  $(x_t = \text{enc}_\phi(o_t))$
4. Posterior:  $(z_t \sim q_\phi(z \mid h_t, x_t))$
5. Bundle:  $(s_t = \{h_t, z_t\})$
6. Heads:  $(\hat{r}_t, \hat{c}_t, \hat{v}_t)$  from  $(s_t)$ ;  $(\hat{a}_{t-1})$  from  $((x_t, s_{t-1}))$ ; optional decoder  $(\hat{o}_t)$  from  $(\text{sg}(s_t))$ .

**Normalization:** Use BatchNorm in the representation net to avoid collapse.

# Architecture Diagram



# Loss Stack (World Model)

Total model loss (sequence length ( $T$ )):

$$\mathcal{L}_{\text{model}} = \sum_{t=1}^T \left( \beta_{\text{pred}} \mathcal{L}_{\text{pred},t} + \beta_{\text{dyn}} \mathcal{L}_{\text{dyn},t} + \beta_{\text{rep}} \mathcal{L}_{\text{rep},t} \right).$$

**Prediction loss ( $\mathcal{L}_{\text{pred}}$ )** trains:

- ▶ Reward:  $(-\log p_{\phi}(r_t | s_t))$
- ▶ Continue:  $(-\log p_{\phi}(c_t | s_t))$
- ▶ **Value:**  $(-\log p_{\phi}(R_t^{\lambda} | s_t))$  (two-hot symlog discretization)
- ▶ **Past Action:**  $(-\log p_{\phi}(a_{t-1} | x_t, s_{t-1}))$
- ▶ Optional reconstruction (aux):  $(-\log p_{\phi}(o_t | \text{sg}(s_t)))$

# Loss Stack (KL Terms)

## Dynamics KL (train prior to match posterior)

$$\mathcal{L}_{\text{dyn},t} = \max(1, \text{KL}[\text{sg}q_{\phi}(z_t \mid h_t, x_t) \mid p_{\phi}(z \mid h_t)]).$$

## Representation KL (train posterior to be predictable)

$$\mathcal{L}_{\text{rep},t} = \max(1, \text{KL}[q_{\phi}(z_t \mid h_t, x_t) \mid \text{sg}p_{\phi}(z \mid h_t)]).$$

**Typical weights:** ( $\beta_{\text{pred}} = 1.0, \beta_{\text{dyn}} = 0.95, \beta_{\text{rep}} = 0.05$ ).

**Note:** KLs are capped below at 1 for stability.

# Value Head Target (on Real Sequences)

Use a slow **EMA teacher** ( $v_{\phi'}$ ) to define  $(\lambda)$ -return targets:

$$R_t^\lambda = r_{t+1} + \gamma c_{t+1} ((1 - \lambda) v_{\phi'}(s_{t+1}) + \lambda R_{t+1}^\lambda), \quad R_T^\lambda = v_{\phi'}(s_T).$$

- ▶ Transform with symlog, discretize with two-hot bins; train the student value head ( $v_\phi$ ) with cross-entropy.
- ▶ Purpose: shape **representations** (not used for behavior learning on imagined rollouts).



# Behavior Learning (Dreamer-style)

## Imagination in latent space (unchanged from V3):

- ▶ Start from a posterior state, roll out ( $H = 15$ ) steps with the **prior** and the **actor**.
- ▶ Use model heads ( $\hat{r}, \hat{c}$ ) on imagined states to construct discretized ( $\lambda$ )-return targets.
- ▶ **Critic** ( $v_\psi$ ) learns on imagined trajectories, with EMA regularizer toward ( $v_{\psi'}$ ).
- ▶ **Actor** maximizes normalized advantages plus entropy; backprop-through-model (continuous) or REINFORCE (discrete).

**Key separation:** behavior uses ( $v_\psi$ ), not the world-model ( $v_\phi$ ).

# Why Past-Action & Value Heads?

- ▶ **Value head (on real data):** forces latents to encode *long-horizon, return-relevant* features; improves stability.
- ▶ **Past-action head:** inverse dynamics from  $((x_t, s_{t-1})) \rightarrow (a_{t-1})$  gives dense, control-aligned supervision each step, especially when rewards are sparse.
- ▶ Together they replace pixel reconstruction as a representation-learning signal, while keeping behavior learning unchanged.

## Results: DeepMind Visual Control Suite (1M steps)

- ▶ Competitive or better than DreamerV3 on several tasks (e.g., Cheetah Run, Quadruped Walk, Reacher Hard).
- ▶ Overall mean/median scores comparable or improved with faster convergence in many cases.
- ▶ Long-horizon latent predictions remain accurate without reconstruction gradients.

# DeepMind Visual Control Suite Results

Task	Random	TPC <sup>‡</sup>	DreamerPro <sup>‡</sup>	DreamerV3 <sup>‡</sup>	MuDreamer
Acrobot Swingup	0.3	5.1	13.1	9.1	<b>41.9</b>
Cartpole Balance	329.3	792.9	870.1	198.7	<b>974.8</b>
Cartpole Balance Sparse	53.9	26.9	198.4	18.4	<b>898.7</b>
Cartpole Swingup	67.4	574.8	689.2	145.7	<b>794.4</b>
Cartpole Swingup Sparse	0.0	0.2	<b>17.8</b>	0.3	0.0
Cheetah Run	6.7	440.8	<b>380.7</b>	94.3	318.1
Cup Catch	31.5	451.5	437.5	27.9	<b>904.5</b>
Finger Spin	0.9	696.8	<b>724.2</b>	96.5	644.2
Finger Turn Easy	48.8	<b>479.5</b>	232.4	197.8	229.4
Finger Turn Hard	35.0	198.3	<b>228.3</b>	39.8	226.7
Hopper Hop	0.0	0.2	<b>1.4</b>	0.6	0.2
Hopper Stand	1.9	14.5	<b>296.5</b>	3.0	5.4
Pendulum Swingup	2.0	<b>778.7</b>	777.6	8.0	606.8
Quadruped Run	8.8	162.9	470.8	108.9	<b>735.0</b>
Quadruped Walk	110.0	681.4	784.5	61.2	<b>872.8</b>
Reacher Easy	52.6	642.4	692.7	154.2	<b>914.4</b>
Reacher Hard	7.4	7.0	9.4	10.6	<b>13.5</b>
Walker Run	25.9	137.9	402.9	78.7	<b>432.8</b>
Walker Stand	139.4	935.4	940.6	254.4	<b>966.7</b>
Walker Walk	36.8	428.3	736.1	164.7	<b>759.0</b>
Mean	47.9	372.8	445.2	83.6	<b>517.0</b>
Median	28.7	409.8	416.2	57.1	<b>620.0</b>

## Results: Natural Background (Robustness to Visual Distractors)

- ▶ Natural-video backgrounds make pixel reconstruction overly track textures.
- ▶ MuDreamer's aux decoder reconstructions visually **filter irrelevant background**, evidence the latent ignores distractors.
- ▶ Reported mean/median (example): Mean  $\approx 517$ , Median  $\approx 620$  (vs. Dreamer baselines under same setting).

# Results: Atari

Game	Random	Human	Lookahead search		No lookahead search			
			MuZero	EffZero	SimPLe	IRIS	DreamerV3	MuDreamer
Alien	228	7128	530	1140	617	420	<b>959</b>	951
Amidar	6	1720	39	102	74	143	139	<b>153</b>
Assault	222	742	500	1407	527	<b>1524</b>	706	891
Asterix	210	8503	1734	16844	1128	854	932	<b>1411</b>
Bank Heist	14	753	193	362	34	53	<b>649</b>	156
Battle Zone	2360	37188	7688	17938	4031	<b>13074</b>	12250	12080
Boxing	0	12	15	44	8	70	78	<b>96</b>
Breakout	2	30	48	406	16	<b>84</b>	31	34
Chopper Com.	811	7388	1350	1794	979	<b>1565</b>	420	808
Crazy Climber	10780	35829	56937	80125	62584	59324	<b>97190</b>	96128
Demon Attack	152	1971	3527	13298	208	<b>2034</b>	303	553
Freeway	0	30	22	22	17	<b>31</b>	0	5
Frostbite	65	4335	255	314	237	259	909	<b>1652</b>
Gopher	258	2412	1256	3518	597	2236	<b>3730</b>	1500
Hero	1027	30826	3095	8530	2657	7037	<b>11161</b>	8272
James Bond	29	303	88	459	100	<b>463</b>	445	409
Kangaroo	52	3035	63	962	51	838	4098	<b>4380</b>
Krull	1598	2666	4891	6047	2205	6616	7782	<b>9644</b>
Kung Fu Mas.	258	22736	18813	31112	14862	21760	21420	<b>26832</b>
Ms Pacman	307	6952	1266	1387	1480	999	1327	<b>2311</b>
Pong	-21	15	-7	21	13	15	<b>18</b>	<b>18</b>
Private Eye	25	69571	56	100	35	100	882	<b>1042</b>
Qbert	164	13455	3952	15458	1289	746	3405	<b>4061</b>
Road Runner	12	7845	2500	18512	5641	9615	<b>15565</b>	8460
Seaquest	68	42055	208	1020	<b>683</b>	661	618	428
Up N Down	533	11693	2897	16096	3350	3546	7600	<b>26494</b>
#Superhuman	0	N/A	5	13	1	10	9	<b>11</b>
Human Mean	0%	100%	56%	190%	33%	105%	112%	<b>126%</b>
Human Median	0%	100%	23%	116%	13%	29%	<b>49%</b>	43%

# Ablations

## Action & Value Heads

- ▶ Remove **Action** → performance drops (notably on sparse-reward tasks like Hopper Hop, Cartpole Swingup Sparse).
- ▶ Remove **Value** → less stable training.
- ▶ Remove **Both** → larger drop.

**Conclusion:** Both heads help the model learn dynamics-aware, control-relevant latents.

## BatchNorm & KL Balancing

- ▶ **BatchNorm in representation** prevents collapse without pixel reconstruction; improves stability and speed.
- ▶ **KL balance** ( $(\beta_{\text{rep}})$  vs.  $(\beta_{\text{dyn}})$ ) matters: too little rep-KL destabilizes; too much slows learning.
- ▶ Reported curves (mean score across 20 DMC-Visual tasks) highlight sensitivity.

## My Skepticism (Failure Modes)

- ▶ **Action–state ambiguity:** inverse dynamics may be ill-posed; weak learning signal.
- ▶ **Partial observability/occlusion:**  $(x_t)$  may not reflect  $(a_{t-1})$ 's effect.
- ▶ **Shortcut risks:** leaks if using  $(s_t)$  as input; overfit to trivial cues without strong data augs.
- ▶ **Non-controllable but relevant factors:** inverse dynamics won't emphasize them; value head helps but not always enough.



# What I'd Test or Swap In

- ▶ **Stress tests:** sparse reward + distractors; measure linear probes for velocities/contacts; MI ( $I(a_{t-1}; x_t \mid s_{t-1})$ ); rollout error.
- ▶ **Alternatives:** action-conditioned contrastive (InfoNCE), bisimulation-style regularization, forward-consistency KL/JSD, controllability Jacobian regularizers.

**Takeaway:** Treat past-action & value heads as *toggleable biases toward controllability and long-horizon relevance*; keep them if ablations pay off.