

# STORM: Efficient Stochastic Transformer based World Models for Reinforcement Learning

March 12, 2025

## Reference

Weipu Zhang, Gang Wang, Jian Sun, Yetian Yuan, and Gao Huang. Storm: Efficient stochastic transformer based world models for reinforcement learning, 2023. URL <https://arxiv.org/abs/2310.09615>.

# Introduction

- ▶ Previous world model approaches:
  - ▶ SimPLe: Uses LSTM
  - ▶ DreamerV3: Employs GRU
  - ▶ RNNs excel at sequence modeling but lack parallel computing capabilities
- ▶ Transformer architecture advantages:
  - ▶ Superior performance over RNNs
  - ▶ Better handles long-term dependencies
  - ▶ Enables efficient parallel computing
- ▶ STORM (Stochastic Transformer-based wORLD Model):
  - ▶ Uses categorical VAE as image encoder
  - ▶ Incorporates Transformer for sequence modeling
  - ▶ Achieves 126.7% mean human normalized score on Atari 100k
  - ▶ Training requires only 4.3 hours on RTX 3090

# Comparison with Recent Approaches

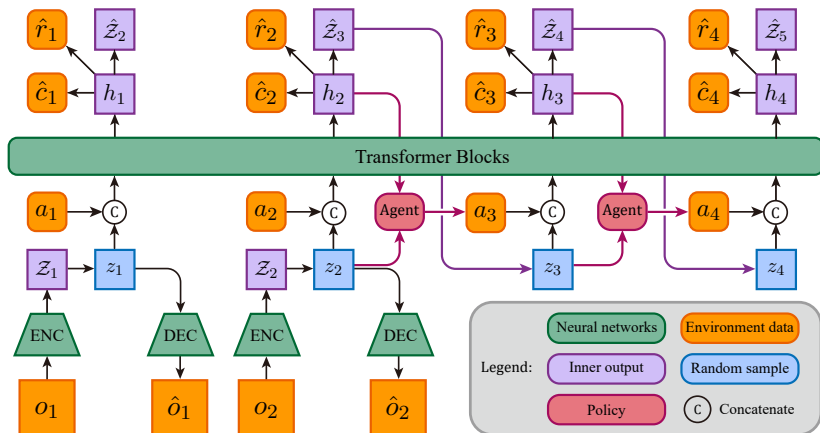
Comparison between STORM and recent approaches. "Tokens" refers to the input tokens introduced to the sequence model during a single timestep. "Historical information" indicates whether the VAE reconstruction process incorporates historical data, such as the hidden states of an RNN.

Attributes	SimPLe	TWM	IRIS	DreamerV3	STORM
Sequence model	LSTM	Transformer-XL	Transformer	GRU	Transformer
Tokens	Latent	Latent, action, reward	Latent( $4 \times 4$ )	Latent	Latent
Latent representation	Binary-VAE	Categorical-VAE	VQ-VAE	Categorical-VAE	Categorical-VAE
Historical information	Yes	No	Yes	Yes	No
Agent state	Reconstructed image	Latent	Reconstructed image	Latent, hidden	Latent, hidden
Agent training	PPO	As DreamerV2	As DreamerV2	DreamerV3	As DreamerV3

## Detailed Comparison with Prior Work

- ▶ SimPLe and Dreamer rely on RNN-based models, whereas STORM employs a GPT-like Transformer as the sequence model.
- ▶ In contrast to IRIS that employs multiple tokens, STORM utilizes a single stochastic latent variable to represent an image.
- ▶ STORM follows a vanilla Transformer structure, while TWM adopts a Transformer-XL structure.
- ▶ In the sequence model of STORM, an observation and an action are fused into a single token, whereas TWM treats observation, action, and reward as three separate tokens of equal importance.
- ▶ Unlike Dreamer and TransDreamer, which incorporate hidden states, STORM reconstructs the original image without utilizing this information.

# STORM (Zhang et al. [2023]) Architecture Overview



# Data Structure and Training Process

- ▶ At each timestep  $t$ , a data point consists of:
  - ▶ Observation  $o_t$
  - ▶ Action  $a_t$
  - ▶ Reward  $r_t$
  - ▶ Continuation flag  $c_t$  (indicates if episode is ongoing)
- ▶ Replay buffer:
  - ▶ Maintains first-in-first-out queue structure
  - ▶ Enables sampling of consecutive trajectories
- ▶ Training process:
  - S1 Execute policy to gather environment data into replay buffer
  - S2 Update world model using sampled trajectories
  - S3 Improve policy using imagined experiences from world model

# Model Structure: VAE Formulation

- ▶ Convert image observations to latent space:

$$\text{Encoder: } z_t \sim q_\phi(z_t|o_t) = \mathcal{Z}_t$$

$$\text{Decoder: } \hat{o}_t = p_\phi(z_t)$$

- ▶ Latent structure  $\mathcal{Z}_t$ :
  - ▶ Stochastic categorical distribution
  - ▶ 32 categories  $\times$  32 classes
  - ▶ Sampled variable  $z_t \sim \mathcal{Z}_t$  represents  $o_t$
- ▶ Implementation details:
  - ▶ CNN-based encoder  $q_\phi$  and decoder  $p_\phi$
  - ▶ Straight-through gradient estimation for backpropagation



# Model Structure: Sequence Model

Action mixer:  $\mathbf{e}_t = m_\phi(\mathbf{z}_t, \mathbf{a}_t)$

Sequence model:  $\mathbf{h}_{1:T} = f_\phi(\mathbf{e}_{1:T})$

Dynamics predictor:  $\hat{\mathbf{z}}_{t+1} = g_\phi^D(\hat{\mathbf{z}}_{t+1} | \mathbf{h}_t)$

Reward predictor:  $\hat{r}_t = g_\phi^R(\mathbf{h}_t)$

Continuation predictor:  $\hat{c}_t = g_\phi^C(\mathbf{h}_t)$

- ▶ Model components and their roles:
  - ▶  $m_\phi$ : Combines latent states and actions into unified tokens through MLP-based fusion
  - ▶  $f_\phi$ : Processes temporal dependencies using a GPT-style Transformer with causal masking
  - ▶  $g_\phi^D$ : Predicts the next latent distribution to capture environment dynamics
  - ▶  $g_\phi^R$ : Estimates immediate rewards based on current hidden states
  - ▶  $g_\phi^C$ : Determines episode termination probability from current context
- ▶ Key feature: Each token  $\mathbf{e}_t$  can only attend to its causal history  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_t$

# Model Structure: Loss Functions

- ▶ Total loss function with hyperparameters  $\beta_1 = 0.5$ ,  $\beta_2 = 0.1$ :

$$\mathcal{L}(\phi) = \frac{1}{BT} \sum_{n=1}^B \sum_{t=1}^T \left[ \mathcal{L}_t^{\text{rec}}(\phi) + \mathcal{L}_t^{\text{rew}}(\phi) + \mathcal{L}_t^{\text{con}}(\phi) + \beta_1 \mathcal{L}_t^{\text{dyn}}(\phi) + \beta_2 \mathcal{L}_t^{\text{rep}}(\phi) \right]$$

- ▶ Environment prediction losses:

- ▶  $\mathcal{L}_t^{\text{rec}}(\phi) = \|\hat{o}_t - o_t\|_2$
- ▶  $\mathcal{L}_t^{\text{rew}}(\phi) = \mathcal{L}^{\text{sym}}(\hat{r}_t, r_t)$
- ▶  $\mathcal{L}_t^{\text{con}}(\phi) = c_t \log \hat{c}_t + (1 - c_t) \log(1 - \hat{c}_t)$

- ▶ Latent dynamics losses (KL-divergence based):

- ▶  $\mathcal{L}_t^{\text{dyn}}(\phi) = \max(1, \text{KL}[q_\phi(z_{t+1}|o_{t+1}) \parallel g_\phi^D(\hat{z}_{t+1}|h_t)])$
- ▶  $\mathcal{L}_t^{\text{rep}}(\phi) = \max(1, \text{KL}[q_\phi(z_{t+1}|o_{t+1}) \parallel \text{sg}(g_\phi^D(\hat{z}_{t+1}|h_t))])$
- ▶ Both losses are KL divergences clamped to minimum value of 1 with selective gradient stopping

# Agent Learning

- ▶ During inference:
  - ▶ Sample  $z_t$  from prior  $\hat{\mathcal{Z}}_t$  rather than posterior  $\mathcal{Z}_t$
  - ▶ Use KV cache in Transformer for faster inference
- ▶ Agent state and policy:

State:  $s_t = [z_t, h_t]$

Critic:  $V_\psi(s_t) \approx \mathbb{E}_{\pi_\theta, p_\phi} \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} \right]$

Actor:  $a_t \sim \pi_\theta(a_t | s_t)$

# Agent Learning: Actor-Critic Algorithm

- Actor-critic loss functions (from DreamerV3):

$$\mathcal{L}(\theta) = \frac{1}{BL} \sum_{n=1}^B \sum_{t=1}^L \left[ -\text{sg} \left( \frac{G_t^\lambda - V_\psi(s_t)}{\max(1, S)} \right) \ln \pi_\theta(a_t | s_t) - \eta H(\pi_\theta(a_t | s_t)) \right] \quad (1a)$$

$$\mathcal{L}(\psi) = \frac{1}{BL} \sum_{n=1}^B \sum_{t=1}^L \left[ \left( V_\psi(s_t) - \text{sg}(G_t^\lambda) \right)^2 + \left( V_\psi(s_t) - \text{sg}(V_{\psi^{\text{EMA}}}(s_t)) \right)^2 \right] \quad (1b)$$

- $\lambda$ -return calculation:

$$G_t^\lambda \doteq r_t + \gamma c_t \left[ (1 - \lambda) V_\psi(s_{t+1}) + \lambda G_{t+1}^\lambda \right] \quad (2a)$$

$$G_L^\lambda \doteq V_\psi(s_L) \quad (2b)$$

# Agent Learning: Training Details

- ▶ Normalization ratio  $S$  for actor loss:

$$S = \text{percentile}(G_t^\lambda, 95) - \text{percentile}(G_t^\lambda, 5) \quad (3)$$

- ▶ Value function regularization using EMA:

$$\psi_{t+1}^{\text{EMA}} = \sigma \psi_t^{\text{EMA}} + (1 - \sigma) \psi_t \quad (4)$$

where:

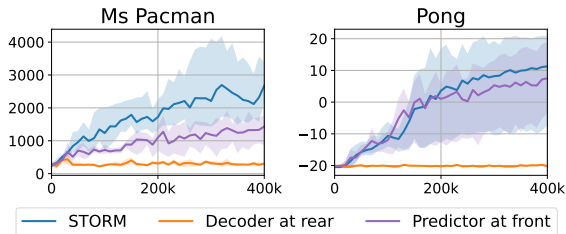
- ▶  $\sigma$  is the decay rate
- ▶  $\psi_t$  represents current critic parameters
- ▶  $\psi_{t+1}^{\text{EMA}}$  denotes updated critic parameters

# Game Scores

Game	Random	Human	SimPLe	TWM	IRIS	DreamerV3	STORM
Alien	228	7128	617	675	420	<b>959</b>	<b>984</b>
Amidar	6	1720	74	122	143	139	<b>205</b>
Assault	222	742	527	683	<b>1524</b>	706	801
Asterix	210	8503	<b>1128</b>	<b>1116</b>	854	932	1028
Bank Heist	14	753	34	467	53	<b>649</b>	<b>641</b>
Battle Zone	2360	37188	4031	5068	<b>13074</b>	12250	<b>13540</b>
Boxing	0	12	8	<b>78</b>	70	<b>78</b>	<b>80</b>
Breakout	2	30	16	20	<b>84</b>	31	16
Chopper Command	811	7388	979	1697	1565	420	<b>1888</b>
Crazy Climber	10780	35829	62584	71820	59234	<b>97190</b>	66776
Demon Attack	152	1971	208	350	<b>2034</b>	303	165
Freeway	0	30	17	24	<b>31</b>	0	<b>34</b>
Freeway w/o traj	0	30	17	24	<b>31</b>	0	0
Frostbite	65	4335	237	<b>1476</b>	259	909	1316
Gopher	258	2413	597	1675	2236	3730	<b>8240</b>
Hero	1027	30826	2657	7254	7037	<b>11161</b>	<b>11044</b>
James Bond	29	303	101	362	463	445	<b>509</b>
Kangaroo	52	3035	51	1240	838	<b>4098</b>	<b>4208</b>
Krull	1598	2666	2204	6349	6616	7782	<b>8413</b>
Kung Fu Master	256	22736	14862	24555	21760	21420	<b>26182</b>
Ms Pacman	307	6952	1480	1588	999	1327	<b>2673</b>
Pong	-21	15	13	<b>19</b>	15	<b>18</b>	11
Private Eye	25	69571	35	87	100	882	<b>7781</b>
Qbert	164	13455	1289	3331	746	3405	<b>4522</b>
Road Runner	12	7845	5641	9109	9615	15565	<b>17564</b>
Seaquest	68	42055	683	<b>774</b>	661	618	525
Up N Down	533	11693	3350	<b>15982</b>	3546	7667	7985
Human Mean	0%	100%	33%	96%	105%	112%	<b>126.7%</b>
Human Median	0%	100%	13%	51%	29%	49%	<b>58.4%</b>

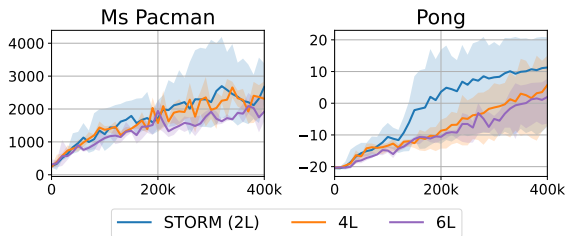
# World Model Architecture Variants

- ▶ Two key architectural variants explored:
  - ▶ “Decoder at rear”: Using  $z_t \sim \hat{\mathcal{Z}}_t$  instead of  $z_t \sim \mathcal{Z}_t$  for observation reconstruction. It shows that reconstruction loss should be applied directly to encoder output
  - ▶ “Predictor at front”: Using  $z_t$  instead of  $h_t$  as input for prediction functions - minimal impact on single-frame reward tasks (e.g. *Pong*), but performance drops on multi-frame reward tasks (e.g. *Ms. Pacman*)



## Model Size and Layer Ablation Study

- ▶ Default configuration uses 2 Transformer layers
  - ▶ Much smaller than 10-layer models in IRIS and TWM
- ▶ Layer ablation study findings:
  - ▶ More layers do not improve performance
  - ▶ In *Pong*, 4 and 6-layer models reach max reward even with 4x samples
- ▶ Three key factors explaining this scaling behavior:
  - ▶ Frame prediction is simplified by:
    - ▶ Small frame-to-frame differences, residual connections aiding prediction
  - ▶ Limited data in Atari 100k:
    - ▶ Lacks image diversity, therefore insufficient samples for larger models
  - ▶ End-to-end training effects:
    - ▶  $\mathcal{L}^{\text{rep}}$  loss directly impacts encoder: Large sequence models may overly influence encoder





# Agent State Representation Study

- Options for agent's state  $s_t$  include:
  - Predicted observation  $\hat{o}_t$ , hidden state  $h_t$ , latent state  $z_t$ , or combined state  $[h_t, z_t]$
- Key findings from ablation studies:
  - Context-dependent performance:
    - Including  $h_t$  improves performance in context-heavy games (e.g. *Ms. Pacman*)
    - Minimal impact in simpler games like *Pong*, *Kung Fu Master*
  - Limitations of using only  $h_t$ :
    - Can lead to catastrophic forgetting in policy-dependent environments
    - World model's non-stationarity creates instability
  - Benefits of randomness:
    - Including distributions like  $\mathcal{Z}_t$  helps mitigate  $h_t$ -only issues

