# ENVISION THE FUTURE IN OPEN-WORLD DYNAMIC TASKS BY A HIERARCHICAL WORLD MODEL WITH RESIDUAL ENHANCED FORESIGHT

**Anonymous authors**Paper under double-blind review

### **ABSTRACT**

Interacting with dynamic objects and even opponent agents in an open world remains a challenge for reinforcement learning. Task planning representations are crucial in such scenarios. Existing reasoning representations grounded in language or vision have demonstrated efficacy, yet most require pretraining and finetuning on domain-specific knowledge datasets. We argue that a reasoning representation purely learned from self-supervised environmental interactions, integrated with brain-like hierarchical structure, offers substantial value for openworld dynamic tasks. In this paper, we present ResDreamer, a hierarchical world model with residually connected visual planning representations. In ResDreamer, high-level world model observes lower level reconstruction residuals, aiming to capture more advanced world dynamics and form a more comprehensive internal world representation. Each layer of the world model employs enhanced environmental observations, which include visual foresight reconstructed from imagined trajectories. These foresight images are further calibrated by residuals predicted by the higher-level world model. Our approach demonstrates higher sampling efficiency, parameter efficiency, and scalability compared to state-of-the-art methods.

### 1 INTRODUCTION

In interaction or combat scenarios, task objectives are relevant with dynamic elements or even another active agent. This pose a significant challenge to decision-making agent. The vast state space of an open-ended environment exacerbates this challenge. The agent must construct an internal world representation based on partial information and make decisions accordingly.

Task planning representation is critical for achieving human-level intelligence in open-ended dynamic environments. World model has pushed the boundaries of reinforcement learning (RL) by reasoning in internal latent space (Schrittwieser et al., 2020; Robine et al., 2023; Zhang et al., 2023; Alonso et al., 2024). DreamerV3 (Hafner et al., 2025) achieved high performance generalization across over 150 diverse tasks with a unified hyperparameters. However, most existing model based RL (MBRL) methods only consider pixel reconstruction as one of the gradient signals for representation learning during training. The decision making process did not directly benefit from the world model's ability to predict future sensory signal.

Hierarchical methods can naturally integrate with task planning by transmitting plan representations between layers. Strategies at different levels allow for independent optimization on different time scales (Lee et al., 2022; Gumbsch et al., 2024). The latent space targets can provide guidance for the lower-level worker (Hafner et al., 2022; Vezhnevets et al., 2017). JARVIS-1 (Wang et al., 2023b), MC-Planner (Wang et al., 2023c), and RL-GPT (Liu et al., 2024) are open-ended embodied agents that integrates RL with Large Language Models (LLMs). Leveraging their generalized world knowledge, LLMs can provide advanced, interpretable language-based task plans through approaches such as task decomposition and policy-as-code. However, low-frequency target may become infeasible when the task relevant objects are in dynamic moving and active states.

We hold the view that multi-level reasoning is crucial for open-world general intelligence. Therefore, an ideal planning representation should be able to hierarchically capture the dynamic of the world at different levels of abstraction, and should be naturally scalable. Neuroscience evidence suggests

that the biological neural signals encode prediction error rather than the raw image (Rao & Ballard, 1999; Hosoya et al., 2005). Visual neurons employ a dynamic predictive coding strategy to filter out predictable components from the visual stream, transmitting only unexpected surprise or "report valuable" stimuli (Kok & de Lange, 2015).

Based on the above insights, we present ResDreamer, a hierarchical world model with residually connected visual planning representations. The residual of visual reconstruction serves as a signal for interlayer interaction in hierarchical world models. Information about prediction errors and feedback is transmitted between layers without the propagating gradients. The higher-level world model, by modeling visual residuals, not only constructs a comprehensive internal representation of the world but also refines the lower level's predictions through residual reasoning, thereby providing more accurate foresight.

In summary, the major contributions of this work are:

- We present a general hierarchical architecture for world models. Through the innovative design of enhanced visual observation, foresight prediction and sensory surprise are transmitted between neighboring layers in a brain-like manner.
- Experimental results validate the sample efficiency, parameter efficiency and scalability of our approach in MBRL context. This has enabled the world model to advance towards the scalable "ResNet era".

### 2 RELATED WORK

MBRL. Recurrent world dynamic models facilitate representation, simulation and policy improvement in MBRL (Ha & Schmidhuber, 2018). MuZero (Schrittwieser et al., 2020) conducts Monte Carlo tree search in the latent space by the learned state space model. DreamerV3 (Hafner et al., 2025) outperformed expert models tuned for specific domains and, for the first time, successfully collected diamonds from scratch in Minecraft. LS-Imagine (Li et al., 2024) breaks the limitations of single-step reasoning and uses the affordance map to trigger the cross-step jump prediction. It simulates jumping to the vicinity of high return targets in the future by magnifying specific areas in the observed image. In the field of visual MBRL, transformer (Micheli et al., 2022; Robine et al., 2023; Zhang et al., 2023), diffusion model (Alonso et al., 2024) are also known as effective world models. However, as far as we know, there is no MBRL method that naturally builds a hierarchical representation learning architecture based on the reconstruction residuals of sensory signals.

Hierarchical RL. Hierarchical RL is considered promising in alleviating the exploration stagnation caused by sparse rewards in complex and long-term tasks. Capturing task-relevant details across varying temporal scales is a key focus of current research efforts (McInroe et al., 2022; Schiewer et al., 2024; Lin et al., 2024; Li et al., 2024). Beyond static temporal scales, THICK (Gumbsch et al., 2024) adaptively discovers larger temporal scales by guiding lower-level world models to sparsely update their partial latent states. Automatic goal discovery is another critical aspect, enabling agents to autonomously identify and pursue meaningful objectives (Hafner et al., 2022; Hamed et al., 2024; Nicklas Hansen, 2025). Furthermore, hierarchical methods can benifit from internet-scale datasets to provide generalized prior knowledge for the lower-level policy (Baker et al., 2022; Yuan et al., 2024; 2023). However, none of the existing method exchange visual residual signal between layers. In our approach, higher-level models can be scalably stacked to achieve increasingly comprehensive representation learning.

World model. Recently, Vision-Language Model (VLM) (Cen et al., 2025) and Joint Embedding Predictive Architecture (JEPA) (LeCun, 2022) have emerged as competitive world model architectures. Web-scale pre-training and relatively sufficient expert demonstration are often prerequisite of VLA driven Minecraft agents (Wang et al., 2023a;c; Li et al., 2025). JEPA is a self-supervised representation learning framework. As an autoregressive generative architecture, it is pretrained on Internet scale multimedia data in the absence of pixel-level reconstruction. Various instances of JEPA have demonstrated its potential across a wide range of domains, including images (Assran et al., 2023), videos (Bardes et al., 2024; Assran et al., 2025), optical flow(Bardes et al., 2023), point clouds (Saito et al., 2025), and graph data (Skenderi et al., 2023). Our approach not only enables plug-and-play training from scratch but also maintains pathways for acquiring internet-scale knowledge guidance through methods such as MineClip (Fan et al., 2022).

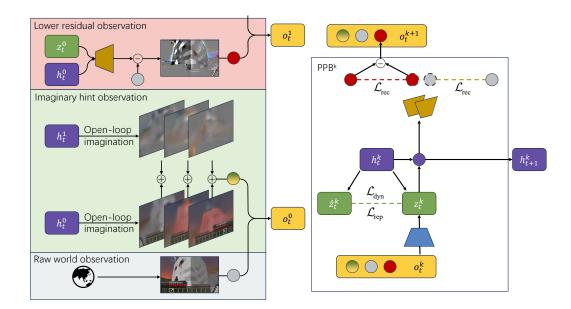


Figure 1: Overview of ResDreamer a model base RL algorithm based on hierarchical world model. The left side shows the structure of enhanced visual observations, through which the world model layers communicate. The right side shows the modules and training process of the k-th layer world model. The Encoder reads enhanced visual observations and gives the posterior  $z_t^k$ . The dynamic predictor learns to estimate  $z_t^k$  with  $\hat{z}_t^k$  without accessing the observation. The sequence model updates internal state  $h_t^k$  by  $z_t^k$ . The Decoder reconstructs the observation signal which generates reconstruction loss and residual visual signal for upper layer.

### 3 Method

In this section, we present the details of ResDreamer. We introduce ResDreamer from the perspectives of representation learning and behavior learning. First, we describe the basic module of each layer in our Hierarchical Recurrent State-Space Model (HRSSM). Next, we present our primary innovation in representation learning architecture, namely the enhanced observation through residual connection. Finally, we formalize the loss functions and the training algorithm.

### 3.1 HIERARCHICAL WORLD MODEL

We implement the HRSSM based on Predictive Processing Blocks (PPBs). Predictive Processing or Predictive coding is a paradigm to explain hierarchical reciprocally connected organization of the cortex (Huang & Rao, 2011).

In the k-th layer block PPB<sup>k</sup>, recurrent state contains the deterministic state  $h_t^k$  and the stochastic state  $z_t^k$ . The sequence model is used to represent the state transitions conditioned by action taken. The Encoder extracts useful information from the new input observations to guide the recurrent state update, while the Predictor attempts to predict the stochastic state without accessing the observations.

$$\text{PPB}^k \begin{cases} \text{Sequence model:} & h_t^k = S_\phi \left( z_{t-1}^k, h_{t-1}^k, a_{t-1} \right) \\ \text{Encoder:} & z_t^k \sim q_\phi \left( z_t^k \mid h_t^k, o_t^k \right) \\ \text{Predictor:} & \hat{z}_t^k \sim p_\phi \left( \hat{z}_t^k \mid h_t^k \right) \\ \text{Decoder:} & \hat{o}_t^k \sim D_\phi \left( \hat{o}_t^k \mid h_t^k, z_t^k \right). \end{cases}$$

where  $\hat{z}_t^k$  is the predicted stochastic state,  $o_t^k$  and  $\hat{o}_t^k$  are true and reconstructed observations. Layer index  $k=0,1,\cdots L-1$  and L is the number of HRSSM layers. Each layer's PPB module contains all the components of the dreamerV3 (Hafner et al., 2025).

### 3.2 VISUAL HINT STRUCTURE AND RESIDUAL MODELING

Figure 1 gives an overview of ResDreamer, a hierarchical world model in which layers communicate through error feedback and predictive visual hints. Enhanced observation  $o_t^k$  become the channel for feedforward and feedback information between adjacent hierarchical world models. They have become the key to our hierarchical world model, enabling it to scale up with linearly increasing communication bandwidth.

$$o_t^k = \left\{ o_{\text{imag}}^k, o_{\text{raw}}, o_{\text{res}}^k \right\}_t,$$

$$\hat{o}_t^k = \left\{ \hat{o}_{\text{raw}}^k, \hat{o}_{\text{res}}^k \right\}_t.$$
(2)

where subscripts  $(\cdot)_{imag}$ ,  $(\cdot)_{raw}$  and  $(\cdot)_{res}$  stands for **imaginary hint**, raw world and **lower residual observation** respectively.  $o_{raw}$  is always original sensory input during environmental interaction. None of these observations propagate gradients.

The **lower residual observation**  $o_{\text{res}}^k$  is reconstruction error from lower level, thus is empty for bottom layer.

$$o_{\text{res}}^{k} = \begin{cases} \text{empty set,} & k = 0, \\ \text{Norm}^{k} \left( o_{\text{raw}} - \hat{o}_{\text{raw}}^{0} \right), & k = 1, \\ \text{Norm}^{k} \left( o_{\text{res}}^{k-1} - \hat{o}_{\text{res}}^{k-1} \right), & k = 2, 3, \dots, L-1. \end{cases}$$
(3)

where the omitted time indices are all t, and the same applies hereafter. Norm<sup>k</sup>(·) computes the mean and variance across the pixel dimension and updates them using an exponential moving average. The **lower residual observation** and the raw environmental observation are sensory signals of equivalent status, both requiring the Decoder to reconstruct.

It is worth noting that any layer of the well-trained PPB can generate future imaginary trajectories by replacing the posterior with the prior without observation. The **imaginary hint observation**  $o_{\text{imag}}^k$  is the channel concatenation of imagined video frames on the planned action trajectory. Due to incomplete modeling, the reconstructed video from the imagined model state can be fuzzy and distorted. The upper-level world model is precisely trained to reproduce the reconstruction error of the current layer. Therefore, we add the residual video from the upper layer onto current video prediction as correction.

$$o_{\text{imag}}^{k} = \begin{cases} \left\{ \hat{o}_{\text{raw}}^{0} + \hat{o}_{\text{res}}^{1} \right\}_{t+1:t+H}, & k = 0, \\ \left\{ \hat{o}_{\text{res}}^{k} + \hat{o}_{\text{res}}^{k+1} \right\}_{t+1:t+H}, & k = 1, 2, \cdots, L-2, \\ \left\{ \hat{o}_{\text{res}}^{k} \right\}_{t+1:t+H}, & k = L-1, . \end{cases}$$

$$(4)$$

The **imaginary hint observation** utilities the HRSSM's reasoning capabilities to directly envision the future, thereby providing additional hints for the current moment. From the perspective of convolutional neural networks (CNN), the imaginary hint effectively generates dynamic CNN kernels based on predictive visual foresight. This process is similar to "gaze control" in cognitive science, which refers to the fact that attention is determined by knowledge-driven prediction (Jovancevic-Misic & Hayhoe, 2009; Henderson, 2017).

Specifically, if the raw image shape is (h, w, 3), then  $o_{\text{res}}^k$  has shape  $(h, w, 3 \times H)$  and  $o_{\text{imag}}^k$  has shape  $(h, w, 3 \times H)$ . Figure 3 shows the raw observation and imaginary hint observation on world model bottom layer while the agent is combating a ghast. The complete process of constructing enhanced observations from the bottom layer to the top layer and updating the recursive state in sequence is shown in Algorithm 1.

At this point, we have established the feedforward and feedback information channels of the hierarchical world model based on the enhanced observation (see Figure 2). This architecture combines the bandwidth advantage of inter-layer communication and the computational efficiency advantage within layers.

The visual hint does incur necessary computational cost, but from the perspective of parameter scale, the above architecture introduces almost no overhead. Within each layer, although the number of image channels has significantly increased due to the addition of video hint, the distribution of the visual hint is highly matched with the original image distribution benefits from the residual modeling, thus allowing for the sharing of major convolutional features. Therefore, in practice,

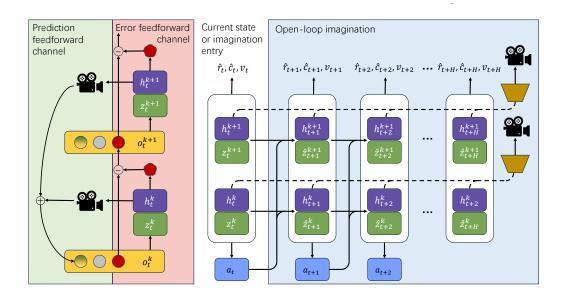


Figure 2: The information channel between world model layers is bidirectional. Only errors and predicted pixel values are transmitted between layers, with no gradients being passed. On one hand, each layer of the PPB generates predictions about the external world and transmits visual planning representations to lower layers accordingly. On the other hand, the PPB learns to reconstruct signals from lower layers, and the error signals are fed back to higher layers.

we have not expanded the depth of the encoder and dimensions of stochastic state compared to dreamerV3 (Hafner et al., 2025).



Figure 3: Visualization of Residual Enhanced Visual Observation on layer 0. It can be seen that at timestep 0, the agent had inferred in imagination that the opponent would turn red and enter the attack state at timestep 2, and drop the bomb at timestep 4. **Green**: raw observation. **Blue**: imaginary hint observation. **Yellow**: imaginary at time-step of next row. Visual hints is reconstructed from imaginary internal states. Upper layer residual hint is added to imagined future image. Each row shows the complete observation of the input encoder, with an interval of 2 timesteps. The video segment shows a ghast faces the agent and shoots a fireball. The agent reasons in imagination and makes a retreating evasive move.

### 3.3 WORLD MODEL AND BEHAVIOR LEARNING

In the context of RL, the final goal is to improve the policy. The actor-critic method is employed for policy optimization and state value learning.

Actor: 
$$a_t \sim \pi_\theta (a_t \mid s_t)$$
  
Critic:  $v_t \sim v_\psi (v_t \mid s_t)$  (5)

where  $s_t = \{s_t^0, s_t^1, \cdots, s_t^{k-1}\}$ . The actor-critic is conditioned on the concatenation of  $s_t^k$  from all layers.

Additional information such as rewards and episode continuation flags are predicted from the recurrent state.

Reward head: 
$$\hat{r}_t \sim p_\phi \left( \hat{r}_t \mid s_t \right)$$
  
Continue head:  $\hat{c}_t \sim p_\phi \left( \hat{c}_t^k \mid s_t \right)$  (6)

The world environment generously provides continuous stream of sensory signals. Reconstructing sensory inputs serves as a critical training signal for world models. This drives the model to encode as much environmental information as possible in deterministic state.

$$\mathcal{L}_{\text{rec}}^{k}\left(\phi\right) = -\ln p_{\phi}\left(o_{\text{raw}} \mid s_{t}^{k}\right) - \ln p_{\phi}\left(o_{\text{res}}^{k} \mid s_{t}^{k}\right) \tag{7}$$

where  $k = 0, 1, \dots, L - 1$ .

The stochastic state serves as the information channel through which new observation guide model updates. For instance, in a typical configuration of 32 categorical variable with 32 classes, the encoder extracts only 256 bits of information from observations at each step. Therefore, the encoder has to retain only the most critical information for updating the internal model. The enforced sparsity makes the stochastic state more feasible to predict, while the representation loss ensures that it tends to converge to a more predictable representation.

$$\mathcal{L}_{\text{dyn}}^{k}(\phi) = \max\left(1, \text{KL}\left[\text{sg}\left(q_{\phi}\left(z_{t}^{k} \mid h_{t}^{k}, o_{t}^{k}\right)\right) \mid\mid p_{\phi}\left(z_{t}^{k} \mid h_{t}^{k},\right)\right]\right)$$

$$\mathcal{L}_{\text{rep}}^{k}(\phi) = \max\left(1, \text{KL}\left[q_{\phi}\left(z_{t}^{k} \mid h_{t}^{k}, o_{t}^{k}\right) \mid\mid \text{sg}\left(p_{\phi}\left(z_{t}^{k} \mid h_{t}^{k}\right)\right)\right]\right)$$
(8)

The prediction heads are similarly trained in a self-supervised manner, with the only difference being that they are conditioned on a stack of recurrent states from all layers.

$$\mathcal{L}_{\text{heads}}(\phi) = -\ln p_{\phi} \left( r_t \mid s_t \right) - \ln p_{\phi} \left( c_t \mid s_t \right) \tag{9}$$

Assuming that the world dynamic and the task-related experiences can be stably represented by the world model, the actor-critic can learn from the imaginary state trajectories, thereby significantly improving the sample efficiency.

The value distribution may span multiple orders of magnitude. Therefore, we parameterize the critic as a categorical distribution with exponentially spaced bins. We compute the bootstrapped  $\lambda$ -return  $R_t^\lambda$  to train the critic.  $R_t^\lambda$  accounts for  $r_t$  within the trajectory horizon T and incorporates the critic's expected value for returns beyond the horizon. The reward signal  $r_t$  may originate from the environment or be estimated by the reward prediction head from imagined trajectories.

$$\mathcal{L}(\psi) = -\sum_{t=1}^{T} \ln p_{\psi} \left( R_{t}^{\lambda} \mid s_{t} \right)$$

$$R_{t}^{\lambda} = \begin{cases} r_{t} + \gamma c_{t} \left( (1 - \lambda) v_{t} + \lambda R_{t+1}^{\lambda} \right), & t < T \\ \mathbb{E} \left[ v_{\psi} \left( \cdot \mid s_{t} \right) \right], & t = T \end{cases}$$

$$(10)$$

The actor learns to maximize returns with entropy regularizer. To remain robust to outliers, we track the range between the 5th and 95th percentiles of returns using an exponential moving average. For further details on the loss function, please refer to Hafner et al. (2025).

$$\mathcal{L}(\theta) = -\sum_{t=1}^{T} \frac{R_t^{\lambda} - \operatorname{sg}(v_{\psi}(s_t))}{\operatorname{max}(1, S)} \operatorname{log} \pi_{\theta}(a_t \mid s_t) + \eta \operatorname{H}[\pi_{\theta}(a_t \mid s_t)]$$

$$S = \operatorname{EMA}\left(\operatorname{Per}(R_t^{\lambda}, 95) - \operatorname{Per}(R_t^{\lambda}, 5)\right)$$
(11)

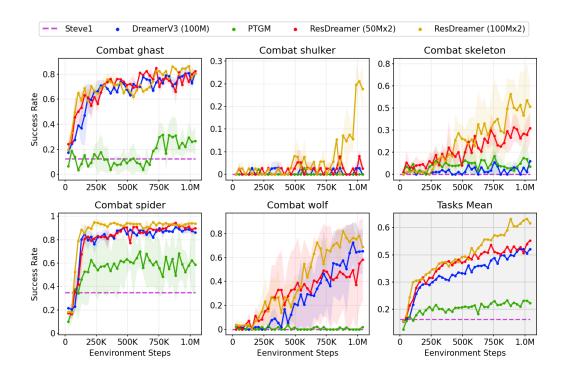


Figure 4: Comparison of ResDreamer against Steve-1 (Lifshitz et al., 2023), DreamerV3 (Hafner et al., 2025), PTGM (Yuan et al., 2024). We introduce the compared models in Appendix C.

### 4 EXPERIMENTS

Engaging in combat within open-ended worlds presents significant challenges including terrain comprehension, the utilization of weapons and defensive tools, and dynamic anticipation of enemy movements. We evaluate ResDreamer on 5 combat tasks in MineDojo (Fan et al., 2022) as is introduced in Table 2.

The agent is equipped with iron armors and iron sword shield at initialization in all tasks. We adopt sparse reward from MineDojo at episode termination and dense reward from MineCLIP (Fan et al., 2022). Each MineCLIP reward is computed of video segment of 16 time-steps, with calculations taking place every 8 frames. In addition, the agent is rewarded at any valid attack and punished for losing health points. The agent is trained for  $1\times10^6$  environment steps. Image input for both agent and MineCLIP model is  $160\times256$  pixels. All experiments can be reproduced with VRAM less than 29 GB.

### 4.1 MAIN COMPARISON

We measure the performance for all the methods with success rates during training. Our implementation is based on DreamerV3 (Hafner et al., 2025) and provides a brain-inspired hierarchical scaling method for it. To make a fair comparison with it in terms of parameter efficiency, ResDreamer is tested with two parameter configurations. Further details are provided in Appendix A.

ResDreamer (100Mx2) adopts residual enhanced observations to extend the Dreamer world model to 2 layers. It demonstrates the best sample efficiency and convergence speed across all tasks. ResDreamer (100Mx2) is also the only method that have significant probability of defeating a shulker in  $1\times10^6$  environment steps. As is shown in Table 1, despite using sparse hierarchical connections and only 84% of the parameter size, ResDreamer (50Mx2) has still surpassed the average performance of the DreamerV3.

The training curves in Figure 4 and comparisons in Figure 6 suggest that ResDreamer is an effective method for hierarchical scale-up of world models. The residual-enhanced visual reasoning represen-

Table 1: ResDreamer and DreamerV3 baseline model sizes

Configurations	DreamerV3	ResDreamer (50Mx2)	ResDreamer (100Mx2)
D	C1.4.4	4006	C144
Recurrent $h_t$ size	6144	4096	6144
Recurrent $z_t$ size	$32 \times 48$	$32 \times 32$	$32 \times 48$
Hidden size	768	512	768
Encoder CNN channels	48	32	48
Decoder CNN channels	32	32	32
hierarchies	1	2	2
Total parameters	109.5M	92.0M	192.7M

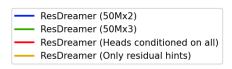
tation can leverage the inference and reconstruction capabilities of world models to achieve dynamic visual reference encoding for current observations, while also enabling error propagation upward to form a more comprehensive world representation across multiple layers of the world model.

### 4.2 MODEL ANALYSIS

The enhanced observation through residual connection is a key feature of ResDreamer, enabling the flow of predictive and error information across the layers of the world model. Figure 5 show the results of the following alternative setups.

ResDreamer (50Mx3): the ResDreamer model from the main comparison is extended to three layers. ResDreamer aims to learn more comprehensive world representations in a scalable manner. The mean task success rate of the three-layer ResDreamer surpasses that of the two-layer version. This indicates that ResDreamer provides an effective method for scaling up world models, achieving model parallelism with linearly increasing communication bandwidth consumption.

ResDreamer (Heads conditioned on all): the actor, critic, and prediction heads in ResDreamer are conditioned on the recursive states of all layers. Experimental results show a performance decline under equivalent environmental interaction steps. Theoretically, complete recursive states contain more information, but the distribution of lower residual observations shifts during training process of lower-layer models, leading to relatively unstable representations in the upper-layer world model before convergence. Further studies could compare performance across additional environmental interaction step settings.



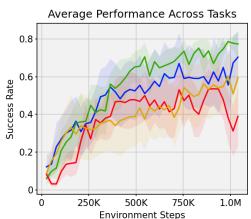


Figure 5: Alternative ResDreamer setups study results.

**ResDreamer (Only residual hints)**: the imaginary hint observations in ResDreamer consist solely of residual signals from the upper layer, without incorporating the current layer's predictive reconstruction. Although the current layer's recursive state already encompasses complete information from open-loop predictions, we find that imaginary hint observations with residual connections yield superior performance. This suggests that visual foresight corrected by residuals facilitates the encoder's learning of a more favorable posterior distribution.

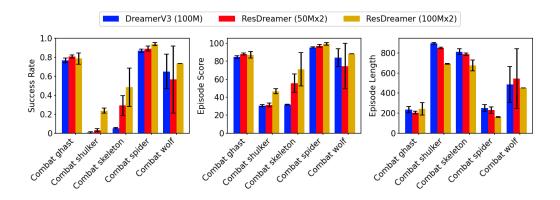


Figure 6: Comparisons of success rate, episode score and episode length across tasks. It can be seen that ResDreamer achieves higher scores and success rates with fewer steps. Although the ResDreamer (50Mx2) has slightly fewer total parameters than DreamerV3 (100M), it performs better in almost all tasks.

### 5 CONCLUSION

In this paper, we present ResDreamer, a hierarchical world model featuring residual-connected visual planning representations. Residual enhanced observations establish an information channel between layers. Those explainable sensory signals are stripped, and the remaining novel stimuli are passed on to higher-level models for learning. The high-level predictions will modify the visual-based planning representation with residual signal, helping the encoder perform gaze control based on accurate temporal predictions. Through comparisons with baselines and model analysis, we demonstrate that ResDreamer achieves superior sample efficiency with fewer parameters compared to baselines. ResDreamer facilitate excellent world model scalability. Data exchange occurs only between adjacent layers, and the communication bandwidth consumption increases linearly with the number of parameters.

The primary limitation of ResDreamer lies in its static image foresight horizon length. Long-horizon goal image increases computational cost, whereas overly short visual hints may fail to provide sufficient environmental dynamics information. We leave the development of adaptive-length image foresight to future work.

### REPRODUCIBILITY STATEMENT

We submit the source code as part of supplementary materials. Following the experiment setup in Table 1, all the results can be reproduced on publicly available RL environments and open source code repositories.

### ETHICS STATEMENT

In this work, we are adhere to the code of ethics. This work does not involve human subjects, personal data, or sensitive information. All training data are synthesized by publicly available environment simulator. Our MBRL approach is task-agnostic, introducing no prior biases. We advocate for thorough testing and safety evaluations before deploying this reinforcement learning system in broader applications, especially physical systems.

### REFERENCES

Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos J Storkey, Tim Pearce, and François Fleuret. Diffusion for world modeling: Visual details matter in atari. *Advances in Neural Information Processing Systems*, 37:58757–58791, 2024.

- Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat,
   Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding
   predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15619–15629, 2023.
  - Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, et al. V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *arXiv* preprint arXiv:2506.09985, 2025.
  - Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022.
  - Adrien Bardes, Jean Ponce, and Yann LeCun. Mc-jepa: A joint-embedding predictive architecture for self-supervised learning of motion and content features. *arXiv preprint arXiv:2307.12698*, 2023.
  - Adrien Bardes, Quentin Garrido, Jean Ponce, Xinlei Chen, Michael Rabbat, Yann LeCun, Mahmoud Assran, and Nicolas Ballas. Revisiting feature prediction for learning visual representations from video. *arXiv preprint arXiv:2404.08471*, 2024.
  - Shaofei Cai, Zhancun Mu, Anji Liu, and Yitao Liang. Rocket-2: Steering visuomotor policy via cross-view goal alignment. *arXiv preprint arXiv:2503.02505*, 2025a.
  - Shaofei Cai, Zhancun Mu, Haiwen Xia, Bowei Zhang, Anji Liu, and Yitao Liang. Scalable multitask reinforcement learning for generalizable spatial intelligence in visuomotor agents, 2025b. URL https://arxiv.org/abs/2507.23698.
  - Jun Cen, Chaohui Yu, Hangjie Yuan, Yuming Jiang, Siteng Huang, Jiayan Guo, Xin Li, Yibing Song, Hao Luo, Fan Wang, et al. Worldvla: Towards autoregressive action world model. *arXiv* preprint arXiv:2506.21539, 2025.
  - Jingwen Deng, Zihao Wang, Shaofei Cai, Anji Liu, and Yitao Liang. Open-world skill discovery from unsegmented demonstrations. *arXiv preprint arXiv:2503.10684*, 2025.
  - Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. URL https://openreview.net/forum?id=rc8o\_j818PX.
  - Christian Gumbsch, Noor Sajid, Georg Martius, and Martin V. Butz. Learning hierarchical world models with adaptive temporal abstractions from discrete latent dynamics. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=TjCDNssXKU.
  - David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Advances in neural information processing systems*, 31, 2018.
  - Danijar Hafner, Kuang-Huei Lee, Ian Fischer, and Pieter Abbeel. Deep hierarchical planning from pixels. *Advances in Neural Information Processing Systems*, 35:26091–26104, 2022.
  - Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse control tasks through world models. *Nature*, pp. 1–7, 2025.
- Hany Hamed, Subin Kim, Dongyeong Kim, Jaesik Yoon, and Sungjin Ahn. Dr. strategy: Model-based generalist agents with strategic dreaming. In *International Conference on Machine Learning*, 2024.
  - John M Henderson. Gaze control as prediction. *Trends in cognitive sciences*, 21(1):15–23, 2017.
    - Toshihiko Hosoya, Stephen A Baccus, and Markus Meister. Dynamic predictive coding by the retina. *Nature*, 436(7047):71–77, 2005.

- Yanping Huang and Rajesh PN Rao. Predictive coding. Wiley Interdisciplinary Reviews: Cognitive Science, 2(5):580–593, 2011.
- Jelena Jovancevic-Misic and Mary Hayhoe. Adaptive gaze control in natural environments. *Journal* of Neuroscience, 29(19):6234–6238, 2009.
  - Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.
  - Peter Kok and Floris P de Lange. Predictive coding in sensory cortex. In *An introduction to model-based cognitive neuroscience*, pp. 221–244. Springer, 2015.
  - Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62(1):1–62, 2022.
  - Seungjae Lee, Jigang Kim, Inkyu Jang, and H Jin Kim. Dhrl: A graph-based approach for long-horizon and sparse hierarchical reinforcement learning. *Advances in Neural Information Processing Systems*, 35:13668–13678, 2022.
  - Jiajian Li, Qi Wang, Yunbo Wang, Xin Jin, Yang Li, Wenjun Zeng, and Xiaokang Yang. Open-world reinforcement learning over long short-term imagination. arXiv preprint arXiv:2410.03618, 2024.
  - Muyao Li, Zihao Wang, Kaichen He, Xiaojian Ma, and Yitao Liang. Jarvis-vla: Post-training large-scale vision language models to play visual games with keyboards and mouse. *arXiv:2503.16365*, 2025.
  - Shalev Lifshitz, Keiran Paster, Harris Chan, Jimmy Ba, and Sheila McIlraith. Steve-1: A generative model for text-to-behavior in minecraft. *Advances in Neural Information Processing Systems*, 36: 69900–69929, 2023.
  - Haoxin Lin, Yu-Yan Xu, Yihao Sun, Zhilong Zhang, Yi-Chen Li, Chengxing Jia, Junyin Ye, Jiaji Zhang, and Yang Yu. Any-step dynamics model improves future predictions for online and offline reinforcement learning. *arXiv* preprint arXiv:2405.17031, 2024.
  - Shaoteng Liu, Haoqi Yuan, Minda Hu, Yanwei Li, Yukang Chen, Shu Liu, Zongqing Lu, and Jiaya Jia. Rl-gpt: Integrating reinforcement learning and code-as-policy. *Advances in Neural Information Processing Systems*, 37:28430–28459, 2024.
  - Trevor McInroe, Lukas Schäfer, and Stefano V Albrecht. Multi-horizon representations with hierarchical forward models for reinforcement learning. *arXiv preprint arXiv:2206.11396*, 2022.
  - Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world models. *arXiv preprint arXiv:2209.00588*, 2022.
  - Vlad Sobal Yann LeCun Xiaolong Wang Hao Su Nicklas Hansen, Jyothir S V. Hierarchical world models as visual whole-body humanoid controllers, 2025.
  - Rajesh PN Rao and Dana H Ballard. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature neuroscience*, 2(1):79–87, 1999.
  - Jan Robine, Marc Höftmann, Tobias Uelwer, and Stefan Harmeling. Transformer-based world models are happy with 100k interactions. *arXiv preprint arXiv:2303.07109*, 2023.
  - Ayumu Saito, Prachi Kudeshia, and Jiju Poovvancheri. Point-jepa: A joint embedding predictive architecture for self-supervised learning on point cloud. In 2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pp. 7348–7357. IEEE, 2025.
  - Robin Schiewer, Anand Subramoney, and Laurenz Wiskott. Exploring the limits of hierarchical world models in reinforcement learning. *Scientific Reports*, 14(1):26856, 2024.
  - Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.

- Geri Skenderi, Hang Li, Jiliang Tang, and Marco Cristani. Graph-level representation learning with joint-embedding predictive architectures. *arXiv preprint arXiv:2309.16014*, 2023.
- Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International conference on machine learning*, pp. 3540–3549. PMLR, 2017.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv: Arxiv-2305.16291*, 2023a.
- Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, Xiaojian Ma, and Yitao Liang. Jarvis-1: Openworld multi-task agents with memory-augmented multimodal language models. *arXiv preprint arXiv:* 2311.05997, 2023b.
- Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*, 2023c.
- Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing Lu. Skill reinforcement learning and planning for open-world long-horizon tasks. *arXiv preprint arXiv:2303.16563*, 2023.
- Haoqi Yuan, Zhancun Mu, Feiyang Xie, and Zongqing Lu. Pre-training goal-based models for sample-efficient reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- Weipu Zhang, Gang Wang, Jian Sun, Yetian Yuan, and Gao Huang. Storm: Efficient stochastic transformer based world models for reinforcement learning. *Advances in Neural Information Processing Systems*, 36:27147–27166, 2023.

### THE USE OF LARGE LANGUAGE MODELS

We only use LLMs as a tool for improving the quality of writing. We manually write the complete version of the paper. The output of the LLMs is merely used as a synonym replacement for some parts of our manually-written version. LLMs are not used for any creative work such as the research ideation or the design of experiments.

### A MODEL DETAILS

Figure 2 intuitively illustrates the data flow diagram of open-loop imagination and how it constructs enhanced visual observations during training. Our hierarchical model extends the process of updating the internal recurrent state based on observations. See Algorithm 1 for details.

The sequence of environmental interactions stored in the replay buffer is utilized only for training the representational learning of the world model, while policy improvement relies exclusively on imagined trajectories. Consequently, the training pipeline and the environment interaction are entirely asynchronous. For a detailed description of the training pipeline, refer to Algorithm 2.

### B ENVIRONMENT DETAILS

- MineDojo agent's initial inventory includes a iron sword, shield, and a full suite of iron armors across all tasks. The maximum number of time-steps for one episode is 1000. For other specifications, see Table 2.
- As shown in Table 2, the five Mobs each possess distinct characteristics. Each episode terminates upon timeout or when the agent's health reaches zero, which implies that the agent must not only explore and approach enemies but also learn to evade attacks or defend with a shield. The rich interaction mechanisms thoroughly test the generalization capabilities of RL algorithms.

### 648 Algorithm 1 Update the recurrent state of ResDreamer upon observation 649 **Input:** recurrent state $s_t$ , raw observation $o_{\text{raw}}$ . 650 **Output:** recurrent state $s_{t+1}$ , world model losses $\mathcal{L}_{dyn}(\phi)$ , $\mathcal{L}_{rep}(\phi)$ , $\mathcal{L}_{rec}(\phi)$ . 651 1: Open-loop roll out imaginary state-action trajectory $\{\hat{s}^{0:L-1}, a\}_{t+1:t+H}$ 652 2: initiate $o_{res}^k$ with empty set. 653 3: **for** each $k = 0, 1, \dots, L - 1$ **do** 654 4: Compute $o_{\text{imag}}^k$ with Eq. (4). 655 Compute $o_t^k$ with Eq. (2). 5: 656 $z_t^k \leftarrow \text{sample} \left[ q_\phi \left( z_t^k \mid h_t^k, o_t^k \right) \right].$ ▷ Encoder 6: 657 $\hat{z}_t^k \leftarrow \text{sample} \left[ p_\phi \left( z_t^k \mid h_t^k \right) \right].$ ▶ Predictor 658 Compute prediction loss $\mathcal{L}_{dyn}^k(\phi)$ and representation loss $\mathcal{L}_{rep}^k(\phi)$ with Eq. (8). 659 660 9: $h_{t+1}^k \leftarrow S_\phi\left(z_t^k, h_t^k, a_t^k\right).$ 661 10: Compute sensory signal reconstruction $\hat{o}_t^k = \{\hat{o}_{\text{raw}}^k, \hat{o}_{\text{res}}^k\}_t$ . ▷ Decoder 662 Compute reconstruction loss $\mathcal{L}^k_{\text{rec}}(\phi)$ with Eq. (7). 11: 663 12: Compute $o_{res}^k$ with Eq. (3). 664 13: **end for** 665 14: **return** $s_{t+1}, \mathcal{L}_{dyn}(\phi), \mathcal{L}_{rep}(\phi), \mathcal{L}_{rec}(\phi)$ . 666 667 **Algorithm 2** The training pipeline of ResDreamer 668 669 1: initiate parameters $\phi, \theta, \psi$ . 670 2: initiate carried state $s_{carry}$ . 671 while not converged do ▶ World model representation learning 4: 672 Sample a environmental interaction sequence $\{o_{\text{raw}},a\}_{0:T-1}$ from replay buffer. 5: 673 **for** each $t = 0, 1, \dots, T - 1$ **do** 6: 674 Update the $s_{\text{carry}}$ upon $\{o_{\text{raw}}\}_t$ with Algorithm 1. Store trajectory feature $\{h_t^{0:L-1}, z_t^{0:L-1}\}$ and losses $\mathcal{L}_{\text{dyn}}(\phi), \mathcal{L}_{\text{rep}}(\phi), \mathcal{L}_{\text{rec}}(\phi)$ . 7: 675 8: 676 9: end for 677 10: 678 Stack feature sequence $F \leftarrow \{h_{0:T-1}^{0:L-1}, z_{0:T-1}^{0:L-1}\}.$ 11: 679 Compute the bootstrapped $\lambda$ -return $R_t^{\lambda}$ and critic loss $\mathcal{L}(\theta)$ with Eq. 10. 12: 680 13: View F as a batch of entry points sized T. 681 Open-loop roll out imaginary state-action trajectory of B time-steps starting at entry points 14: 682 batch F. 683 15: for each imaginary trajectory $\{\hat{s}_{0:B-1}, a_{0:B-1}\}$ do 684 Compute the normalized return and actor loss $\mathcal{L}(\psi)$ with Eq. 11. 16: 685 17: 686 18: Back propagate losses $\mathcal{L}_{dyn}(\phi)$ , $\mathcal{L}_{rep}(\phi)$ , $\mathcal{L}_{rec}(\phi)$ , $\mathcal{L}(\theta)$ , $\mathcal{L}(\psi)$ . 19: Optimize parameters $\phi$ , $\theta$ , $\psi$ . 687 20: end while 688

### C BASELINE INTRODUCTION

# C.1 SELECTED METHODS

689 690 691

692 693

694

696

697

698

699

700

701

We compare ResDreamer with strong Minecraft RL algorithms, including:

DreamerV3 (Hafner et al., 2025): A model-based RL foundation model. DreamerV3 is trained from scratch without demonstrations and domain knowledge. It generates future latent states recurrently with a non-hierarchical world model.

STEVE-1 (Lifshitz et al., 2023): An finetuned Video Pretraining (VPT) model for open-ended text and visual instructions following. It is post trained through self-supervised behavioral cloning. We test its zero-shoot text instructions following performance in MineDojo tasks.

Table 2: MineDojo tasks specifications.

Mobs	Biome	Mob Features	MineClip prompt
Spider	extreme hills	Fast movement	combat a spider in night extreme hills with a iron sword, shield, and
Shulker	end	Shoots guided bullets which causes floating	a full suite of iron armors combat a shulker in the end with a iron sword, shield, and a full suite of iron armors
Wolf	taiga	More agile, group attacks	combat a wolf in taiga with a iron sword, shield, and a full suite of iron armors
Skeleton	extreme hills	Accurate ranged attacks with arrows	combat a skeleton in night extreme hills with a iron sword, shield, and a full suite of iron armors
Ghast	nether	Flying, ranged attacks with explosive fireball, terrain destruction	combat a ghast in nether with a iron sword, shield, and a full suite of iron armors

PTGM (Yuan et al., 2024): A hierarchical approach integrating a high-level task goal generation strategy and a low-level goal-conditioned RL strategy. The high-level goal strategy is pretrained on large-scale, task-agnostic datasets, while the low-level strategy is learned online through RL. We utilize the open-source upper-layer strategy parameters of PTGM and evaluate its online training performance on MineDojo tasks using the default configuration of PTGM code-base.

### C.2 Unselected Methods

We provide introductions of other strong Minecraft agents and the reasons we do not compare Res-Dreamer with them.

LS-Imagine (Li et al., 2024): An MBRL method that achieves arbitrary time-span reasoning through dual-branch prediction. It is based on DreamerV3, but it supports long-term prediction by simulating jumping to the vicinity of navigation targets through cropping observation. However, combat missions are different from navigation and exploration. Factors such as terrain, enemy reactions, etc. have a significant impact on the expected return, and cutting the images disrupts the data distribution. For instance, it is not reasonable to jump to flying enemies like ghasts by cropping the image.

Voyager (Wang et al., 2023a), JARVIS-1 (Wang et al., 2023b), MC-Planner (Wang et al., 2023c), RL-GPT (Liu et al., 2024): Open-Ended embodied agents that integrates RL with LLM. They adopt heterogeneous hierarchical models, leveraging the prior knowledge of LLMs to achieve task decomposition, long-term planning, code as strategy, and lifelong skill accumulation. Their focus lies in the integration and interaction methods between LLMs and RL, emphasizing the evaluation of an agent's efficiency in accumulating atomic skills and activating technological milestones. Our proposed ResDreamer is a model-based RL foundation model, focusing on evaluating the data efficiency, scalability, and interpretability. ResDreamer can work together with all kinds of upper layer LLMs as a more powerful RL algorithm.

ROCKET-2 (Cai et al., 2025a), ROCKET-3 (Cai et al., 2025b) SkillDiscovery (Deng et al., 2025), JarvisVLA (Li et al., 2025): Open-world VLA agents powered by imitation learning (IL) and prior knowledge of visual foundation model such as SAM (Kirillov et al., 2023). VLA agents focus on following open instructions within a broader range of atomic skills and their combinations. However, ResDreamer is a MBRL foundation model trained without any prior knowledge. ResDreamer focuses on developing a task-agnostic and domain general hierarchical world model method.

## D ADDITIONAL VISUALIZATION

Residual Enhanced Visual Observation is the main innovation of this work. To visually demonstrate the structure of this visual foresight and the planning information it provides, we visualize the observation sequence of the agent during its combat with the wolf in Figure 7.

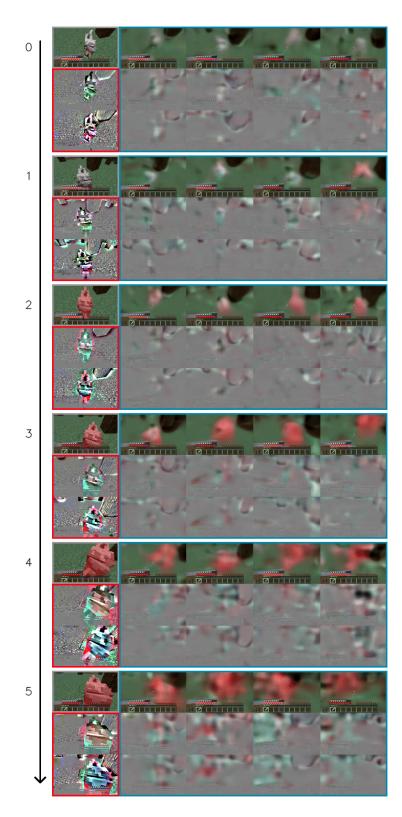


Figure 7: Visualization of all the observations of a ResDreamer with three hierarchies. **Gray**: raw observation. **Red**: residual observation. **Blue**: original open-loop imagination.