# NEURALOS: TOWARDS SIMULATING OPERATING SYSTEMS VIA NEURAL GENERATIVE MODELS

# Anonymous authors

Paper under double-blind review

# **ABSTRACT**

We introduce NeuralOS, a neural framework that simulates graphical user interfaces (GUIs) of operating systems by directly predicting screen frames in response to user inputs such as mouse movements, clicks, and keyboard events. NeuralOS combines a recurrent neural network (RNN), which tracks the computer state, with a diffusion-based neural renderer that generates screen images. The model is trained on a dataset of Ubuntu XFCE recordings, which include both randomly generated interactions and realistic interactions produced by AI agents. Experiments show that NeuralOS successfully renders realistic GUI sequences, accurately captures mouse interactions, and reliably predicts state transitions like application launches. Beyond reproducing existing systems, NeuralOS shows that synthesized training data can teach the model to simulate applications that were never installed, as illustrated by a Doom application, and suggests a path toward learning user interfaces purely from synthetic demonstrations.

# 1 Introduction

"Chatting" with LLM feels like using an 80s computer terminal. The GUI hasn't been invented yet, but some properties of it can start to be predicted.

— Andrej Karpathy

Recent breakthroughs in generative models have transformed human-computer interaction, making it increasingly adaptive, personalized, and intuitive. Historically, computing interfaces were rigid and predefined, such as command-line terminals and static graphical menus (Engelbart, 1968). The emergence of large language models (LLMs) and multimodal AI systems expanded this paradigm by enabling interactions through natural language (Radford et al., 2019; Brown et al., 2020), images (Ho et al., 2020; Lipman et al., 2022; Radford et al., 2021; Song et al., 2020b), and videos (OpenAI, 2024). Recently, generative models have even begun simulating dynamic visual environments (Ha & Schmidhuber, 2018a; He et al., 2025), notably interactive video games (Alonso et al., 2024; Feng et al., 2024; Oh et al., 2015; Valevski et al., 2024). These advancements suggest a future where computing interfaces could become fully generative, dynamically adapting in real-time based on user inputs, contexts, and intentions (Deka et al., 2017).

In this paper, we introduce NeuralOS, a first step toward realizing this vision. NeuralOS simulates an operating system's graphical interface entirely using deep neural networks. By modeling the OS interface as a generative process, it directly predicts graphical frames from user input events, such as mouse movements, clicks, and keyboard interactions, without manually programmed kernels or applications. Figure 1 illustrates an example sequence generated by NeuralOS, demonstrating realistic cursor movements and window interactions predicted solely from user inputs.

NeuralOS integrates two complementary neural architectures, analogous to the traditional separation between OS kernels and desktop rendering programs: a recurrent neural network (RNN) (Hochreiter & Schmidhuber, 1997) maintains internal computer states (such as open applications, hidden windows, and recent actions), while a diffusion-based convolutional neural renderer generates screen images. We train NeuralOS end-to-end on interaction sequences recorded from Ubuntu XFCE environments, combining randomly generated and realistic AI-generated human-like interactions.

Developing NeuralOS posed several challenges. (1) Long-term state tracking was essential due to delayed interface responses (e.g., opening Firefox could take up to 30 frames); we addressed this by

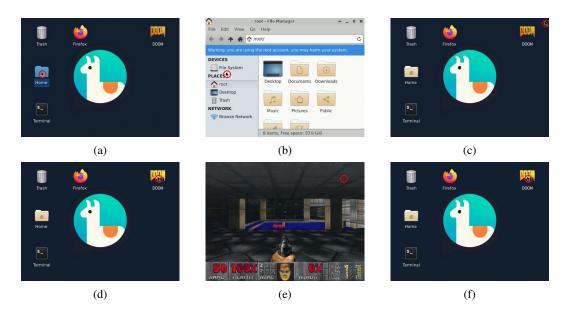


Figure 1: Real image sequence predicted by NeuralOS, illustrating the model's ability to simulate realistic GUI interactions. The sequence shows key frames as a user (a–c) opens and closes the "Home" folder, followed by (d–f) launches and closes a Doom application that was trained into the model using *synthetic* demonstrations. Cursor positions are highlighted with red circles. Frames are generated autoregressively, conditioned on previous frames and user inputs.

using an RNN-based state representation. (2) Precise cursor modeling required explicit positional encodings within our diffusion model. (3) Without pretrained encoders for GUI interactions, we developed a novel pretraining method in which the RNN outputs were pretrained via regression losses and subsequently integrated into the diffusion model via finetuning. (4) Exposure bias during inference was mitigated using scheduled sampling techniques (Bengio et al., 2015; Deng et al., 2023; Ranzato et al., 2015). (5) Extensive engineering was necessary for scalable data collection and real-time inference, leveraging parallel Docker environments and AI-generated user interactions.

Experiments show that NeuralOS can generate realistic screen sequences, accurately model mouse interactions, and reliably simulate transitions such as application launches. While computational constraints limit its ability to precisely model fine-grained keyboard inputs, NeuralOS represents a step toward neural operating systems that adapt interfaces in real time, potentially enabling users to interact through natural language or gestures rather than fixed menus.

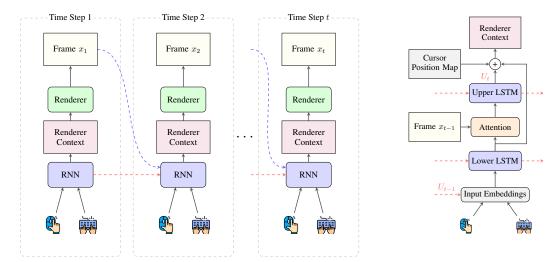
Beyond mimicking an existing operating system, NeuralOS can in principle learn user interfaces from demonstrations even if they are artificially constructed and do not exist in reality. As a proof of concept, we created a Doom application by combining fabricated desktop interactions with Viz-Doom gameplay recordings. NeuralOS learned to launch, play, and close Doom despite the application never being installed in the underlying system. This illustrates the broader principle that synthetic demonstrations, once distilled into a generative model, become usable user interfaces.

Our code, pretrained models, and an interactive demo are at https://neural-os.com.

# 2 GENERATIVE MODELING OF OPERATING SYSTEM INTERFACES

We formalize the task of simulating operating system (OS) graphical interfaces as an autoregressive generative modeling problem. At each discrete timestep t, the model predicts the next graphical frame  $x_t$  based on the sequence of previously observed frames  $x_{< t} = x_0, x_1, \ldots, x_{t-1}$  and user input events  $a_{\le t} = a_1, a_2, \ldots, a_t$  up to and including the current timestep.

Formally, each frame  $x_t$  is represented as an image tensor  $x_t \in \mathbb{R}^{H \times W \times C}$ , with H and W denoting image height and width, and C the number of color or feature channels. The input event  $a_t$  at



(a) High-level temporal architecture of NeuralOS.

(b) RNN structure at step t.

Figure 2: **NeuralOS Model Architecture.** (a) **High-level architecture of NeuralOS**. At each timestep, an RNN tracks the operating system's internal state based on user inputs (cursor positions, mouse clicks, keyboard events) and previously generated frames. This state is then passed as context to a diffusion-based renderer (UNet) that generates the next graphical frame. (b) **Detailed two-level RNN structure at timestep t**. The lower-level LSTM encodes user inputs, and then integrates visual information from the previous frame using attention. Its output is passed to the upper-level LSTM, which further processes these attention-informed representations. Feedback from the upper-level LSTM to the lower-level LSTM ( $U_{t-1}$ ) ensures that the lower-level LSTM maintains awareness of upper-level state context and previous attention results. The combined outputs of both LSTMs, and cursor position encoding, form the renderer context. This hierarchical structure maintains constant computational complexity per timestep and supports continuous state updates during inference, essential for real-time OS interface simulation.

timestep t includes cursor coordinates (x, y), binary indicators for mouse clicks (left or right), and a binary vector indicating which keyboard keys are pressed or released.

The probability distribution of an OS graphical sequence given user inputs can be expressed as:

$$P(x_{1:T} \mid a_{1:T}; \theta) = \prod_{t=1}^{T} P(x_t \mid x_{< t}, a_{\le t}; \theta), \tag{1}$$

where  $\theta$  represents the parameters of the neural generative model.

Unlike standard video generation, OS simulation must respond instantly to unpredictable user inputs, often causing abrupt changes in the interface, such as when a new application is launched. This contrasts with the smooth, predictable transitions typical in video generation. As a result, the model must maintain accurate and responsive state tracking. Next, we describe the NeuralOS architecture and training strategies designed for these requirements.

## 3 NEURALOS ARCHITECTURE

NeuralOS adopts a modular architecture inspired by the functional separation in traditional operating systems between kernel-level state management and graphical user interface (GUI) rendering. It comprises two primary components: a recurrent neural network (RNN) responsible for maintaining internal system states, and a diffusion-based renderer that generates graphical frames based on these states (see Figure 2a).

**Latent Diffusion Representation** NeuralOS uses a latent diffusion framework (Rombach et al., 2022). We train an autoencoder to compress high-resolution OS screen images into lower-dimensional latent representations, reducing spatial dimensions by a scaling factor s. All modeling

is performed within this latent space. At inference time, the generated latent frames are decoded back into pixel-level images only for users.

**Hierarchical RNN for State Tracking** NeuralOS employs a hierarchical two-level RNN architecture to track the system state (see Figure 2b). Unlike transformers (Vaswani et al., 2017), whose inference complexity increases with context length, the RNN maintains constant complexity per timestep, which is crucial for continuous, long-horizon OS simulation.

At each timestep t, user inputs  $a_t$  are encoded into embeddings. Specifically, cursor coordinates are discretized screen positions  $(a_t^x, a_t^y)$ , mouse clicks are binary indicators, and keyboard keys are binary press/release states. Each component is embedded separately and then concatenated:

$$\mathrm{embed}(a_t) = \mathrm{concat}(\mathrm{embed}(a_t^x), \ \mathrm{embed}(a_t^y), \ \mathrm{embed}(a_t^{\mathrm{L\,click}}) + \mathrm{embed}(a_t^{\mathrm{R\,click}}) + \sum_{\mathrm{kev}} \mathrm{embed}(a_t^{\mathrm{key}})).$$

These embeddings are processed by a lower-level LSTM, which also takes its previous hidden state  $l_{t-1}$  and feedback from the previous upper-level LSTM state  $U_{t-1}$  as inputs:

$$L_t, l_t = \text{LSTM}_{\text{lower}}(l_{t-1}, \text{concat}(\text{embed}(a_t), U_{t-1})),$$

where  $l_t$  denotes the hidden state and  $L_t$  denotes the corresponding output at timestep t.

To handle inherent uncertainties in OS behaviors, such as unpredictable application response times, the lower-level LSTM output  $L_t$  is used as a query vector to attend over the previous graphical frame using multi-headed attention (Vaswani et al., 2017):

$$c_t = \text{MultiHeadedAttention}(\text{query} = W_q L_t, \text{keys/values} = W_k x_{t-1} + E_{\text{pos}}),$$

where  $W_q, W_k$  are learnable projections and  $E_{pos}$  encodes positional information of the latent frame.

The attention output  $c_t$  is then merged with the original lower-level LSTM output  $L_t$ :

$$C_t = L_t + W_o c_t,$$

then processed by the upper-level LSTM:

$$U_t, u_t = \text{LSTM}_{\text{upper}}(u_{t-1}, C_t).$$

To ensure that the lower-level LSTM maintains awareness of higher-level contexts, the upper-level LSTM's output  $U_t$  is fed back as an input to the lower-level LSTM in the next timestep.

**Spatial Encoding of Cursor Positions** Precise cursor localization is critical for realistic OS interactions. NeuralOS explicitly encodes cursor positions using a Gaussian spatial map  $E_{\text{pos}} = M_t \in \mathbb{R}^{H \times W}$ . Instead of using a one-hot cursor position (which can lose precision due to latent resolution constraints), we construct a Gaussian map centered at the cursor's scaled coordinates:

$$M_t(i,j) = \exp(-\frac{(i-a_t^x/s)^2 + (j-a_t^y/s)^2}{2}).$$

As demonstrated in Section 7, this spatial encoding is vital for accurate cursor rendering. This map  $M_t$ , combined with LSTM outputs  $L_t$ ,  $U_t$ , forms the renderer context  $R_t \in \mathbb{R}^{H \times W \times C'}$ :

$$R_t = \operatorname{concat}(W_L L_t, W_U U_t, M_t).$$

**Diffusion-Based Renderer** To render the screen image, a UNet-based diffusion renderer generates the latent graphical frames conditioned on the renderer context  $R_t$  (Ronneberger et al., 2015).

$$x_t \sim P_{\theta}(\cdot \mid R_t).$$

We concatenate the noisy image with  $R_t$  as input to the UNet, and then predict the clean image.

#### 4 MULTI-STAGE TRAINING APPROACH

Training NeuralOS is challenging due to ineffective use of RNN outputs, error accumulation during inference, and difficulties in capturing long-term dependencies due to computational constraints. To address these challenges, we take a multi-stage training approach (Figure 3).

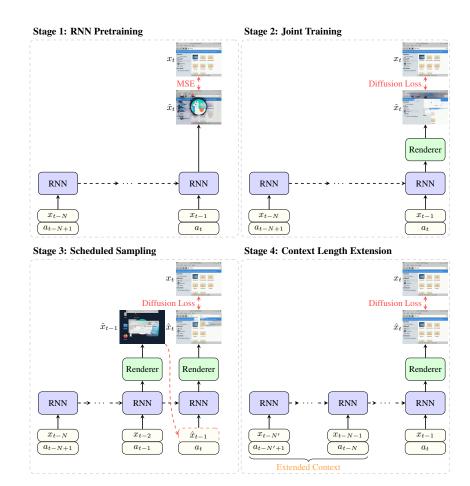


Figure 3: Multi-stage training pipeline for NeuralOS. (1) RNN Pretraining: The RNN is pretrained to predict latent frames using a mean squared error (MSE) loss. (2) Joint Training: The pretrained RNN and the diffusion-based renderer are jointly optimized using a standard diffusion loss. (3) Scheduled Sampling: To mitigate error accumulation caused by exposure bias, the most recent input frame is occasionally replaced by a previously generated frame. (4) Context Length Extension: The input context is extended to enable the model to capture long-term dependencies.

**Stage 1: RNN Pretraining** Unlike text-to-image diffusion models (Rombach et al., 2022), which use pretrained textual encoders, NeuralOS uses a customized RNN without pretrained checkpoints. Our initial experiments show that direct joint training leads to the renderer ignoring RNN outputs, as indicated by negligible gradient flow into the RNN. The diffusion-based renderer receives two streams of inputs: noisy latent frames and the RNN output; and without proper initialization of the RNN, it relies solely on the noisy image inputs, resulting in an under-trained RNN.

To address this, we first pretrain the RNN. We structure the RNN output  $R_t \in \mathbb{R}^{H \times W \times C'}$  to match the spatial dimensions of the latent frame  $x_t \in \mathbb{R}^{H \times W \times C}$ , but with more channels (C' > C). The RNN is pretrained to predict the latent frames  $x_t$  using a mean squared error (MSE) loss:

$$\mathcal{L}_{MSE} = ||R_t[:,:,:C] - x_t||_2^2.$$

After pretraining, the RNN-generated frames alone tend to be blurry due to averaging multiple plausible outcomes, but crucially provide a strong initialization for subsequent joint training.

**Stage 2: Joint Training** We jointly optimize the pretrained RNN and the diffusion-based renderer with a standard diffusion loss. The meaningful latent representations learned in RNN pretraining enable the renderer to use the RNN outputs, thus preventing the RNN outputs from being ignored.

**Stage 3: Scheduled Sampling** During inference, errors accumulate over time and progressively degrade the quality of the generated frames. This issue arises from exposure bias (Ranzato et al., 2015): a model trained exclusively on ground-truth frames becomes overly reliant on perfect inputs and struggles when forced to operate on its own imperfect predictions during inference.

To mitigate this, we introduce a scheduled sampling training stage (Ranzato et al., 2015; Deng et al., 2023): during training, we replace the most recent input frame  $x_{t-1}$  with the model-generated frame  $\hat{x}_{t-1}$  with a small probability p. This method makes the model robust against input noise, thus mitigating error accumulation and improving generation stability over extended interactions.

**Stage 4: Context Length Extension** Although NeuralOS can model sequences of arbitrary length, hardware memory limits require shorter sequences during training, which limits exposure to long-term dependencies. To address this, we introduce a final stage that extends training to longer contexts, following initial training on short context windows for efficiency.

Curriculum Training on Challenging Transitions In our collected OS interaction dataset, a substantial portion of frame transitions involve minor variations, such as slight cursor movements, which exhibit limited learning signals. To prioritize learning of significant OS state changes (e.g., opening menus or launching applications), we first train NeuralOS exclusively on challenging transitions. Specifically, we define challenging transitions as those frames whose pixel differences exceed a specified threshold:  $||x_t - x_{t-1}||_2^2 > \epsilon$ . Subsequently, we expand training to the full dataset. We apply this curriculum strategy for Stage 1 (RNN Pretraining) and Stage 2 (Joint Training).

# 5 DATA COLLECTION

**Agent-Based Demonstrations** To collect realistic user interactions, we use Anthropic's Claude-3.5-Sonnet computer-use agent (Anthropic, 2024), which processes screenshots and invokes provided interaction functions. To maximize interaction diversity, we structure the agent's exploration process around a state-space search tree representing various OS states (see Figure 11 in Section H).

Specifically, we prompt the agent to identify all interactable GUI elements by moving the cursor to each element's center and reporting its bounding box (center, width, height). Each identified GUI element becomes a node in a search tree rooted at the initial OS state. The agent is then guided through this tree: for each node, it moves the cursor to the corresponding GUI element and performs single or double clicks to transition to a new OS state, which then becomes a child node in the tree. We iteratively expand the tree until reaching a predefined maximum depth.

Next, we initiate further interactions from each leaf node, allowing the agent to explore freely from these distinct OS states for a fixed duration, thereby capturing diverse interaction sequences.

**Random Exploration** We find that relying exclusively on agent-generated demonstrations introduces spurious correlations. For example, the model incorrectly associates cursor movement toward a window's close button with the action of closing, even in the absence of a click. To mitigate such unintended associations, we supplement the dataset with random interaction data.

In generating these random interactions, we simulate mouse movements, clicks, and keyboard inputs (key presses and releases) stochastically. To improve realism, we introduce several constraints and heuristics iteratively developed through experimentation. Cursor movements are modeled using Bezier curves to emulate natural mouse trajectories. Double-click events, rare under purely random sampling, are explicitly generated.

## 6 EXPERIMENTAL SETUP

**Data** We collected data using 64 parallel Docker containers, each with Ubuntu 20.04 and XFCE desktops at a resolution of  $512 \times 384$ . To simplify the environment, the desktop was limited to four applications: *Home, Trash, Terminal*, and *Firefox*. Data consisted of 2K agent-based and 120K random exploration demonstrations, each 30 seconds long at 15 fps, resulting in 12TB of latent data after compression via an autoencoder. The autoencoder reduced the images by a factor of 8 to a latent resolution of  $64 \times 48$  with 16 channels. Details of the autoencoder are in Section G.

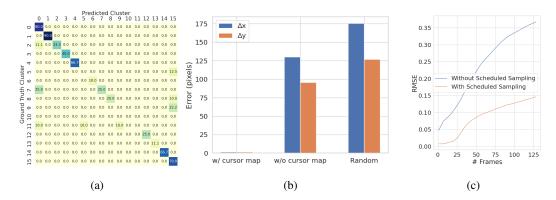


Figure 4: (a) Heatmap illustrating predicted vs. ground truth state transitions. Each cell represents the percentage of predictions assigned to a particular predicted cluster (x-axis), given a ground-truth cluster (y-axis). Only the top 16 clusters are displayed here; refer to Figure 9 for the complete heatmap. (b) Comparison of cursor position errors for NeuralOS (with cursor position map), NeuralOS without the cursor position map, and a random baseline. (c) Pixel RMSE of generated frames of using versus not using stage 3 scheduled sampling training.

**Model** The hierarchical RNN has two LSTM modules (each with hidden size 4,096) and a multiheaded attention module (8 heads, 1,024 total dimension). The RNN output is projected to 32 channels and concatenated with the noisy latent frame (16 channels). The UNet uses four resolution levels with channel multipliers of [1, 2, 3, 5], two residual blocks per level, and attention layers at resolutions 8, 4, and 2. It has a base model dimension of 192 and outputs 16 channels. The final model contains 2.2B parameters for the RNN and 263M parameters for the UNet renderer.

**Training and Inference** NeuralOS was trained using our proposed multi-stage training approach. See Section A for full details on hyperparameters. The total data processing and training took approximately 4 months, requiring 17,000 GPU hours on a server with 8 NVIDIA H200 GPUs (141GB memory per GPU) and an additional 6,000 GPU hours on a server with 8 NVIDIA H100 GPUs (80GB memory per GPU). At inference time, we used DDIM sampling (Song et al., 2020a) with 2 steps, achieving an inference speed of 18 fps on a single NVIDIA H100 GPU.

# 7 EXPERIMENTS

Given the substantial computational resources required to train NeuralOS, our evaluation focused on NeuralOS variants, ablations, and intermediate training phases. For all evaluations, we used a subset of 730 examples from a reserved evaluation dataset, unless mentioned otherwise.<sup>1</sup>

Cursor Position Modeling We evaluated cursor-position accuracy by training a regression model to detect cursor coordinates from the generated images. With the cursor position map, NeuralOS achieved highly accurate cursor localization, with an average position error of  $\Delta x=1.6$  and  $\Delta y=1.4$  pixels (Figure 4b). Given the  $512\times384$  resolution of the images, this corresponds to less than 0.5% of the frame width or height, indicating that cursor locations in generated images are very precise. This performance significantly outperformed a baseline without cursor position maps  $(\Delta x=130.0, \Delta y=95.8)^2$  and the random baseline  $(\Delta x=175.4, \Delta y=126.9)$ , confirming the importance of explicit spatial encoding for accurate cursor localization.

**State-Transition Modeling** Most OS interactions involve only cursor movements without significant visual changes. To focus evaluation on the crucial moments where the interface undergoes meaningful changes, such as opening or closing an application, we identified challenging frame

<sup>&</sup>lt;sup>1</sup>This number was chosen because our clustering procedure (detailed later) identified 73 clusters of challenging frame transitions, and we selected 10 examples per cluster, resulting in a total of 730 examples.

<sup>&</sup>lt;sup>2</sup>This baseline is an earlier model version trained for 700K steps under slightly different conditions.







Figure 5: Doom interactions generated by NeuralOS. Doom was never installed in the underlying operating system used for data collection; instead, the model learned to simulate the application from synthesized training data, producing realistic walking and shooting behavior from user inputs.

Table 1: Accuracy (%) of identifying the real OS, averaged over 30 video comparisons.

<b>Crop Setting</b>	10s	20s	30s	40s	50s	60s
No crop		55.0				
Bottom cropped (40px)	51.7	53.3	50.0	48.3	46.7	50.0

transitions from the evaluation set.<sup>3</sup> These transitions were clustered into 73 categories using DB-Scan. NeuralOS predictions were then matched against the nearest cluster labels and compared with the ground truth. As shown in Figure 4a, NeuralOS achieved an accuracy of 37.7% (diagonal alignment), substantially outperforming majority voting (1.4%). Note that off-diagonal predictions may still correspond to valid outcomes due to inherent timing variability in OS actions (see Section B).

**Human Evaluation** We conducted a human evaluation to test how indistinguishable NeuralOS outputs are from a real operating system. Human raters were shown side-by-side video clips of the same user interaction sequence, one generated by NeuralOS and one recorded from an Ubuntu XFCE desktop, and asked to identify the real system. Each clip lasted 10–60 seconds (sampled at 10s intervals) and was drawn from real human interactions collected through our online demo.<sup>4</sup>

To control for artifacts such as fluctuating disk-space counters (a known issue in diffusion), evaluation was performed under two conditions: (1) original videos and (2) videos with bottom 40 pixels (10% of height) cropped. In each condition, participants judged 30 randomly ordered comparisons.

As shown in Table 1, participants performed only slightly better than chance on up to 20 seconds of interactions, suggesting that NeuralOS is visually realistic for basic interactions.

**Learning from Synthetic Demonstrations** A benefit of an end-to-end learned user interface is its ability to learn from demonstrations that do not correspond to real implementations. To illustrate this, we constructed a Doom application entirely from synthetic data, even though Doom was never installed in the docker. Desktop screenshots were augmented with a Doom icon, synthetic double-click actions were inserted to simulate launching the game, VizDoom gameplay segments with corresponding input events were spliced in, and each sequence ended with an ESC action returning to the desktop. Trained on these demonstrations, NeuralOS learned to launch, play, and close Doom as if it was a native application (Figures 1 and 5).

**Long-Term Memory** We evaluated whether NeuralOS has long-term memory using a folder creation task. We navigated to  $Home \rightarrow Documents$ , created a "New Folder" (via right-click  $\rightarrow Create Folder$ ), closed the file manager, and later reopened Documents after a delay. During training, the model saw only short delays: we included 500 examples where the file manager was reopened after 8 frames in the training set. At test time, we extended the delay to 64, 128, and 256 frames. As a control, we also included 500 examples where no folder was created.

Folder presence was evaluated by checking the region where the folder icon should appear, comparing it to the ground truth with an MSE threshold. Table 2 shows the model can mostly recall whether a folder had been created even after 256 frames, which is far beyond both the 8-frame training delay

<sup>&</sup>lt;sup>3</sup>These were those with mean pixel distance greater than 0.1 from the previous frame to the target frame.

<sup>&</sup>lt;sup>4</sup>See Figure 19 in Section N for a screenshot of the human evaluation interface.

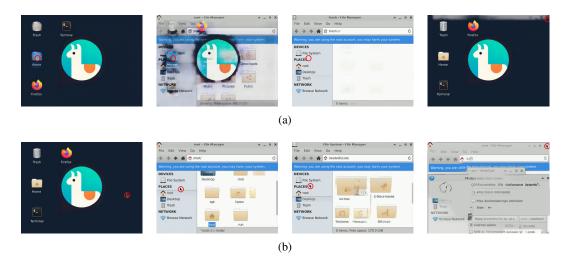


Figure 6: Ablation studies. (a) shows the limitations of directly using RNN outputs. (b) shows error accumulation without scheduled sampling, demonstrating its necessity for maintaining stability.

Table 2: Folder presence modeling accuracy after reopening file manager with increasing gaps.

Frame Gap	Ground Truth	Predicted		
Trume Sup	Ground frum	Has Folder	No Folder	
64	Folder Created No Folder Created	<b>0.52</b> 0.03	0.48 <b>0.97</b>	
128	Folder Created No Folder Created	<b>0.60</b> 0.00	0.40 <b>1.00</b>	
256	Folder Created No Folder Created		0.38 <b>0.98</b>	

and its maximum training context of 64 frames. This indicates that NeuralOS develops persistent memory representations of application state that generalize well beyond its training horizon.

**Ablations** Without joint training (relying solely on the pretrained RNN), the predictions are blurry (Figure 6a). This is caused by the MSE loss encouraging the model to predict averaged representations of multiple plausible outcomes rather than committing to a single clear target. Additionally, cursor positions were absent, despite the model correctly capturing state transitions (e.g., opening the home folder), indicating that the RNN still implicitly encoded cursor information.

Omitting scheduled sampling led to rapid deterioration in generated frame quality due to compounding prediction errors (Figure 6b). In contrast, incorporating scheduled sampling greatly improved the model's robustness (Figure 1), which is also evident from quantitative evaluation in Figure 4c.

#### 8 CONCLUSION AND FUTURE WORK

We introduced NeuralOS, a system that simulates OS graphical interfaces using generative models. Trained end-to-end on interaction sequences, NeuralOS produces realistic screen sequences, accurately predicts mouse interactions, and captures transitions such as opening applications.

A benefit of end-to-end trained interfaces is their ability to learn from *synthetic* demonstrations. Our Doom experiment illustrated this: the model simulated launching, playing, and closing Doom, even though the application was never installed in the underlying operating system. This points to a broader paradigm for future generative interfaces: if consistent demonstrations can be provided, even if they are manually edited, stylized, or fabricated, they can be internalized into usable interfaces.

# 9 REPRODUCIBILITY STATEMENT

We have made all code, model checkpoints, and data processing scripts available at https://anonymous.4open.science/r/neural-os. An interactive demo is also hosted on HuggingFace Spaces, where both code and trained models are transparent and can be easily duplicated and studied; we will add the link to our HuggingFace space in the final version if the paper is accepted. Details of training procedures and dataset construction are also provided in the Appendix.

#### 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.
- Anthropic. Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku, October 2024. URL https://www.anthropic.com/news/3-5-models-and-computer-use. Accessed: 2025-05-14.
- Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. *Advances in neural information processing systems*, 28, 2015.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.
- Haoxuan Che, Xuanhua He, Quande Liu, Cheng Jin, and Hao Chen. Gamegen-x: Interactive openworld game video generation. *arXiv preprint arXiv:2411.00769*, 2024.
- Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschman, Daniel Afergan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. Rico: A mobile app dataset for building data-driven design applications. In *Proceedings of the 30th annual ACM symposium on user interface software and technology*, pp. 845–854, 2017.
- Yuntian Deng, Noriyuki Kojima, and Alexander M Rush. Markup-to-image diffusion models with scheduled sampling. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=81VJDmOE2ol.
- Douglas C. Engelbart. The mother of all demos. https://www.youtube.com/watch?v=fhEh3tEL1V4, 1968. Demonstration at the Fall Joint Computer Conference, San Francisco, CA.
- Ruili Feng, Han Zhang, Zhantao Yang, Jie Xiao, Zhilei Shu, Zhiheng Liu, Andy Zheng, Yukun Huang, Yu Liu, and Hongyang Zhang. The matrix: Infinite-horizon world generation with real-time moving control. *arXiv preprint arXiv:2412.03568*, 2024.
- David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Advances in neural information processing systems*, 31, 2018a.
- David Ha and Jürgen Schmidhuber. World models. arXiv preprint arXiv:1803.10122, 2018b.
- Haoran He, Yang Zhang, Liang Lin, Zhongwen Xu, and Ling Pan. Pre-trained video generative models as world simulators. *arXiv preprint arXiv:2502.07825*, 2025.
  - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
  - Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.

- Seung Wook Kim, Yuhao Zhou, Jonah Philion, Antonio Torralba, and Sanja Fidler. Learning to simulate dynamic environments with gamegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1231–1240, 2020.
  - Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
  - Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard L Lewis, and Satinder Singh. Action-conditional video prediction using deep networks in atari games. *Advances in neural information processing systems*, 28, 2015.
  - OpenAI. Introducing Sora: OpenAI's Text-to-Video Model. https://openai.com/index/sora, February 2024. Accessed: 2025-04-22.
  - Virtuals Protocol. Video game generation: A practical study using mario, 2024. URL https://github.com/Virtual-Protocol/mario-videogamegen/blob/main/static/pdfs/VideoGameGen.pdf. Preprint.
  - Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
  - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
  - Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence level training with recurrent neural networks. *arXiv preprint arXiv:1511.06732*, 2015.
  - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
  - Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pp. 234–241. Springer, 2015.
  - Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv* preprint arXiv:2010.02502, 2020a.
  - Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv* preprint *arXiv*:2011.13456, 2020b.
  - Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. *arXiv preprint arXiv:2408.14837*, 2024.
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
  - Jiannan Xiang, Guangyi Liu, Yi Gu, Qiyue Gao, Yuting Ning, Yuheng Zha, Zeyu Feng, Tianhua Tao, Shibo Hao, Yemin Shi, et al. Pandora: Towards general world model with natural language actions and video states. *arXiv preprint arXiv:2406.09455*, 2024.
  - Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 1(2):6, 2023.

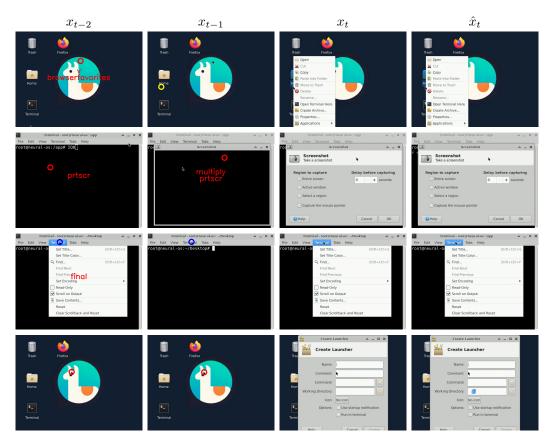


Figure 7: Correct prediction examples by NeuralOS. Each row shows two past frames (columns 1–2), ground-truth next frame (column 3), and NeuralOS's prediction (column 4). Cursor positions are marked one frame in advance with circles (red: move-only, blue: left-click, yellow: right-click). NeuralOS correctly captures various GUI transitions, including opening menus and launching applications.

Table 3: Detailed hyperparameters and dataset configurations used for each stage of NeuralOS's multi-stage training. "Challenging transitions" are where the target frame differs from the preceding input frame by a mean pixel difference greater than 0.1. These challenging transitions constitute approximately 2.8% of the full dataset.

Stage	Dataset	Batch	Steps	LR	Context	Sampling $p$
Stage 1: RNN Pretraining						
Challenging transitions	2.8% subset	256	50K	$8 \times 10^{-5}$	32	
Full dataset	100%	256	200K	$8 \times 10^{-5}$	32	
Stage 2: Joint Training (RN	N + Renderer)					
Challenging transitions	2.8% subset	64	100K	$8 \times 10^{-5}$	32	_
Full dataset	100%	64	1M	$8 \times 10^{-5}$	32	
Stage 3: Scheduled Samplin	g					
Full dataset	100%	256	500K	$8 \times 10^{-5}$	32	0.05
Full dataset (lr reduced)	100%	256	500K	$2 \times 10^{-5}$	32	0.05
Stage 4: Context Length Ex	tension					
Full dataset	100%	128	100K	$2 \times 10^{-5}$	64	0.05

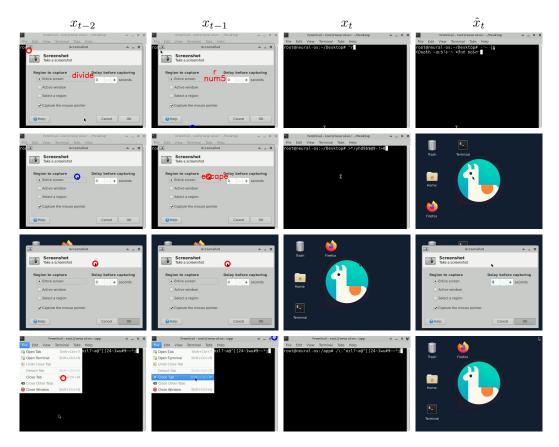


Figure 8: Prediction examples where the generated frame does not match the ground truth frame. The layout follows Figure 7. Note that not all mismatches represent errors. For example, the third row illustrates a case where the screenshot tool window is closed in the ground truth frame but remains open in the prediction. This discrepancy arises because the window-closing action (not shown due to the limited context window) can have variable timing. Thus, both the predicted and ground truth frames are valid outcomes in such scenarios.

#### A MULTI-STAGE TRAINING HYPERPARAMETERS

This section provides detailed hyperparameters and dataset configurations for each training stage of NeuralOS, which is summarized in Table 3.

In Stage 1 (RNN Pretraining), the RNN was trained first on the subset of challenging transitions, defined as examples whose target frame differs from the last input frame by an average pixel difference greater than 0.1. These challenging transitions constitute about 2.8% of the entire dataset. We used a batch size of 256 and an initial learning rate of  $8\times10^{-5}$ , training for 50K steps. Afterwards, the model was trained on the full dataset (100% of data) for an additional 200K steps, maintaining the same batch size and learning rate. The context window length was fixed at 32 frames during this stage.

In Stage 2 (Joint Training), the pretrained RNN and the renderer were jointly trained end-to-end, first focusing on the challenging transitions (2.8% subset) for 100K steps, then extended to the full dataset for an additional 1M steps. The learning rate remained at  $8\times 10^{-5}$ , with a reduced batch size of 64 to stabilize diffusion training. The context length remained 32 frames, and no scheduled sampling was applied in this stage.

In Stage 3 (Scheduled Sampling), we trained on the full dataset using scheduled sampling with probability p=0.05, where the most recent past frame was occasionally replaced by a model-generated frame during training. Initially, training was conducted for 500K steps at a batch size

of 256 and a learning rate of  $8 \times 10^{-5}$ . The learning rate was subsequently reduced to  $2 \times 10^{-5}$ , followed by an additional 500K training steps. The context window remained 32 frames.

Finally, in Stage 4 (Context Length Extension), we increased the context window length from 32 to 64 frames to enable NeuralOS to better capture long-term dependencies. Scheduled sampling probability was maintained at p=0.05. We used a lower learning rate of  $2\times 10^{-5}$ , reduced the batch size to 128 to fit GPU memory constraints, and finetuned the model for 100K additional steps.

# B QUALITATIVE ANALYSIS

We further analyze NeuralOS qualitatively by examining successful and unsuccessful generation examples, shown in Figure 7 and Figure 8, respectively. Each row illustrates two past frames, the ground-truth next frame, and NeuralOS's predicted frame. Cursor positions in past frames are annotated with colored circles to indicate the cursor's intended position at the next frame: red represents cursor movement only, blue denotes left-click actions, and yellow signifies right-click actions. Additionally, keys pressed at each frame are displayed in red text. Note that cursor annotations are shifted forward by one frame to clearly depict cursor positions expected in the immediate subsequent frame.

In Figure 7, NeuralOS accurately predicts various critical GUI transitions, such as launching applications and opening menus through both mouse clicks and keyboard inputs, demonstrating its ability to capture spatial and functional dynamics.

However, as shown in Figure 8, NeuralOS exhibits limitations, particularly for subtle actions like moving the cursor to a "Close Tab" button without clicking. Moreover, NeuralOS currently struggles to accurately represent fine-grained keyboard inputs, such as specific characters typed in a terminal.

It is worth noting that not all mismatches between predictions and ground truth constitute errors; some discrepancies arise from variable timing in GUI responses, exemplified in Figure 8.

## C FULL STATE TRANSITION HEATMAP

Due to space constraints, the main text presents only a truncated version of the state transition heatmap. In Figure 9, we provide the complete heatmap, showing NeuralOS's predictions across all identified clusters.

#### D INTERACTIVE WEB DEMO

To facilitate user interaction with NeuralOS, we developed a web-based frontend using FastAPI, accessible at https://neural-os.com/. Due to input rates (user actions) typically exceeding model inference speeds, we implemented a user-input queue. When the model finishes generating a frame, the system prioritizes processing recent meaningful inputs (clicks and keyboard events), discarding the redundant cursor movements if necessary. This approach maximizes the responsiveness of interactions.

## E TRANSFORMER VS RNN ENCODER

We initially implemented our system using a transformer-based encoder architecture. However, in interactive OS simulation, inference sequences can be arbitrarily long, as each user action generates a new frame that is appended to the model's context unless explicitly truncated. For long or continuous interactions, this results in a steady increase in transformer memory consumption during inference, ultimately leading to out-of-memory (OOM) failures for sufficiently long sequences.

To address this, we replaced the transformer with a hierarchical recurrent neural network (RNN) architecture. RNNs maintain a constant-size hidden state between steps, enabling inference-time GPU memory usage that is independent of sequence length. Table 4 shows the memory requirements of a RNN encoder vs a Transformer encoder. Moreover, computer state modeling is inherently Markovian, making the RNN a suitable choice.

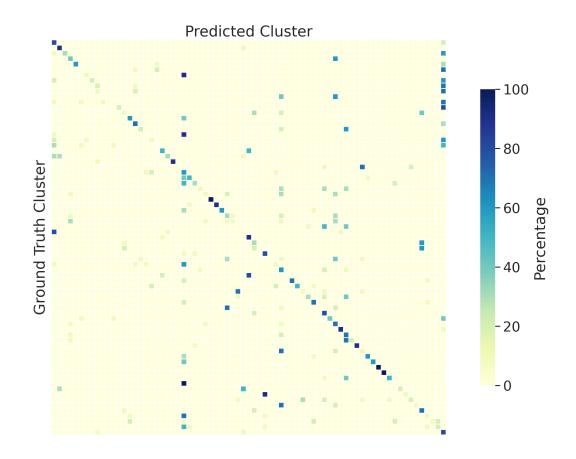


Figure 9: Complete heatmap of NeuralOS state transitions. Each cell represents the percentage of predictions corresponding to a predicted cluster (x-axis) given a ground-truth cluster (y-axis). Diagonal entries indicate exact cluster matches.

Table 4: Inference-time GPU memory usage (GB) for transformer and RNN architectures of similar parameter counts. The transformer encoder cannot accommodate 1024 frames on an L40 GPU.

Video Length (s)	Transformer (GB)	RNN (GB)	
4	21.61	11.96	
64	22.69	11.96	
128	23.81	11.96	
256	26.07	11.96	
512	30.60	11.96	
1024	>40 (OOM)	11.96	

Table 5: Application open–close accuracy. We report RMSE and mean absolute error (L1) between NeuralOS predictions and ground-truth frames, with pixel values normalized to [0, 1].

Application	Action	RMSE	L1
Terminal	Open	0.021	0.001
Firefox	Open	0.048	0.006
Trash	Open	0.017	0.001
Home	Open	0.019	0.001
Terminal	Close	0.015	0.001
Firefox	Close	0.016	0.001
Trash	Close	0.016	0.001
Home	Close	0.016	0.001

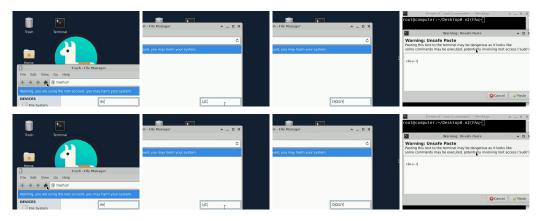


Figure 10: Examples of original images (top row) and their corresponding reconstructions (bottom row) from the trained autoencoder. Despite significant spatial compression (8× downsampling), the autoencoder preserves details.

## F APPLICATION OPEN-CLOSE TRANSITIONS

We evaluated the accuracy of application open–close transitions in NeuralOS. Specifically, we constructed 150 sequences for four main applications (Terminal, Firefox, Trash, Home) by sampling random cursor paths and click locations uniformly within the icon bounding boxes.

We compared the predicted and ground-truth frames at the end of each open and close event using root mean squared error (RMSE) and mean absolute error (L1). To account for rendering latency, we evaluated the frame after a fixed delay following each click: 40 frames after Firefox double-clicks, 10 frames after other application double-clicks, and 5 frames after single-click closing. These delays were determined empirically to align with the frame at which the application's visual state change completes.

As shown in Table 5, the low RMSE and L1 scores confirm that NeuralOS accurately performs application open and close transitions, consistent with observations from the interactive demo.

## G AUTOENCODER DETAILS

We trained a Variational Autoencoder to compress high-dimensional OS screenshots into low-dimensional latent representations suitable for efficient training of NeuralOS.

**Model Architecture** The architecture of the autoencoder is based on the model proposed by (Rombach et al., 2022) with some custom adjustments to improve reconstruction quality. The encoder consisted of four convolutional downsampling blocks with 128 base channels. Each downsampling stage contained two residual blocks and no attention layers. The latent channel dimension was set to 16.

**Training** The autoencoder was trained using a combined reconstruction and adversarial loss function. We trained the autoencoder using the Adam optimizer with a learning rate of 1e-6, a batch size of 10, and a total of 2 million training steps on our dataset. Training was conducted on a single NVIDIA H200 GPU.

After training, the encoder was able to compress each  $512 \times 384$  RGB frame into a latent representation of dimension  $16 \times 64 \times 48$  (downsampled by a factor of 8 in spatial dimensions), significantly reducing memory requirements for subsequent NeuralOS model training.

Examples of original and reconstructed images are provided in Figure 10.

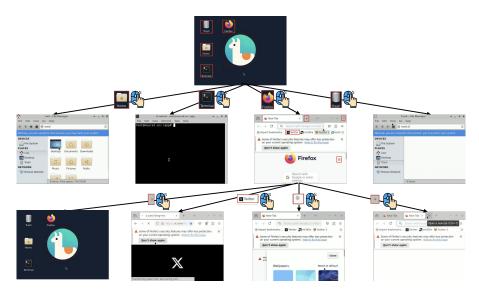


Figure 11: Illustration of Search-Tree-Based Data Collection. We construct a search tree representing OS states, starting from the initial desktop screen (root node). Each node corresponds to a distinct OS state, created by clicking or double-clicking interactable GUI elements identified by a computer-use agent. For clarity, only first-level transitions (opening applications) and one deeper exploration within Firefox are shown. This approach enables collecting diverse interaction data efficiently.

```
I need to generate some training data now. Please iteratively map out every application icon/button on the desktop screen by moving your mouse to it. For each icon, move your mouse EXACTLY to its center point. DO NOT click it though. Be thorough and identify every button or icon on the screen from the top to the bottom.

{{suffix}}
```

Figure 12: **Initial GUI Element Mapping (Root Node).** Prompt instructing the agent to identify interactable GUI elements on the initial desktop screen.

```
Now do a final check: Look carefully at every part of the screen for any unmapped buttons. If you find any, map their exact centers like before. If you are absolutely certain ALL buttons have been mapped, respond with ONLY the final coordinate list.

{{suffix}}
```

Figure 13: **Final Verification of GUI Elements (Root Node).** Prompt instructing the agent to perform a final check for any missed interactable GUI elements on the initial desktop screen.

# H COMPUTER-USE AGENT FOR DATA COLLECTION

To build the search tree illustrated in Figure 11, we used structured prompts to guide the computeruse agent. Starting from the initial desktop screen (root node), we sequentially instructed the agent to first map and then verify all interactable GUI elements (Figures 12 and 13). For subsequent GUI states (non-root nodes), we initially prompted the agent to transition to each new state (Figure 14). After transitioning, we issued follow-up prompts to map and verify any newly revealed GUI elements (Figures 15 to 17).

Each prompt included a standardized suffix, as shown in Figure 18, to ensure that the agent outputs the coordinates and dimensions of mapped GUI elements in a consistent, structured format. This allowed efficient parsing and automated processing of collected data.

I need to generate some training data now. First, please click/open the object at  $\{\{x\}\}$ ,  $\{\{y\}\}$ . Do not progress until it is open or clicked. Double click to open if necessary. Once you confirmed it is open or clicked, stop here.

Figure 14: **Transitioning to a New UI State (Non-Root Node).** Prompt instructing the agent to transition to the current state.

```
Now, iteratively map out all the buttons by moving your cursor to them. Focus on those related to the object you click/open at \{\{x\}\}, \{\{y\}\} - especially new buttons. For each button, move your mouse EXACTLY to its center point. DO NOT click it though. \{\{suffix\}\}
```

Figure 15: Initial Mapping of Newly Revealed GUI Elements (Non-Root Node). Prompt instructing the agent to identify GUI elements newly revealed after transitioning to the current UI state.

```
Please continue mapping out buttons on the screen but DO NOT click anymore buttons. Again, focus on mapping those related to the object you click/open at {{x}}, {{y}} - especially any new ones. As you exhaust the new buttons, move farther out. For each button, move your mouse EXACTLY to its center point. DO NOT click it though.

{{suffix}}
```

Figure 16: Continued Mapping of Remaining GUI Elements (Non-Root Node). Prompt instructing the agent to further map any remaining interactable GUI elements in the current UI state.

```
Now do a final check: Look carefully around the screen to see if there are any new unmapped buttons related to the object you click/open at \{\{x\}\}, \{\{y\}\}. If you find any, map their exact centers like before. DO NOT click more buttons though. If you are absolutely certain ALL related buttons have been mapped, respond with ONLY the final coordinate list. \{\{suffix\}\}\}
```

Figure 17: **Final Verification of GUI Elements (Non-Root Node).** Prompt instructing the agent to perform a final check ensuring all interactable GUI elements have been identified at the current UI state.

# I CURSOR POSITION PREDICTION MODEL TRAINING

To quantitatively evaluate cursor position accuracy in the generated frames (Figure 4b), we trained a regression model to predict cursor coordinates directly from screen images. The training procedure is detailed below.

**Model Architecture** We used a ResNet-50 convolutional backbone pretrained on ImageNet, with modifications for fine-grained spatial localization tasks. Specifically, we adjusted the stride and dilation parameters in the final convolutional layers to reduce downsampling from  $32 \times$  to  $16 \times$ , preserving more spatial resolution. The feature extractor output is passed through an additional intermediate convolutional layer followed by a fully-connected regression head, outputting continuous cursor coordinates (x, y).

**Training** We used the Adam optimizer with an initial learning rate of  $6 \times 10^{-5}$  and weight decay set to  $1 \times 10^{-5}$ . We clipped gradients at a maximum norm of 1.0. We optimized an L1 loss between predicted and ground-truth cursor positions. We trained with a batch size of 16 for a total of 2 epochs. Input images were used directly from collected data at the original resolution ( $512 \times 384$  pixels), normalized and rearranged to match the input format of ResNet-50. The training data consisted of randomly sampled frames from the full dataset, each labeled with the ground truth cursor positions captured during data collection. Training was performed on a single NVIDIA A6000 GPU. The test

```
If you find a keyboard input box, use the type tool specifying its center coordinates and
input box type via text from the following possibilities: terminal, browser, or other, with
coordinates appended in the format x:y.
At the end of your response, provide a structured list of ALL buttons you mapped using this
exact format:
For each button provide two pairs of numbers:
1. Location (L): x~y coordinates where you moved the mouse
{f 2.} Size (S): x:y dimensions of the button
Example format:
    L100^{2}00 S50:30, // Button 1: Located at (100,200) with size 50x30
    L300~400 S60:40
                      // Button 2: Located at (300,400) with size 60x40
Important:
- Use EXACTLY this format: Lx y Sx:y
- Separate each button with commas
 Include ALL buttons you mapped
 Numbers only, no text or descriptions
```

Figure 18: **Structured Output Suffix.** A standardized suffix appended to all prompts, guiding the agent to report mapped GUI elements using a consistent coordinate and dimension format.

error is 0.5 pixels for both x and y. Given that each image is  $512 \times 384$  pixels, this reflects extremely high localization precision. In other words, the regression model can predict the cursor location from a screen image with an average deviation of less than a single pixel, making it highly sensitive to small differences and suitable for evaluating fine-grained cursor accuracy in generated frames.

## J USE OF LLMS

We used large language models (LLMs) in two ways. First, we used an LLM-based computer-use agent (Claude-3.5-Sonnet) to generate interaction traces, which were incorporated into our training dataset alongside synthetic data. Second, we used an LLM-based writing assistant to polish grammar. All ideas, analyses, experiments, and scientific claims are our own, and we take full responsibility for the content of this work.

#### K SCHEDULED SAMPLING IMPLEMENTATION DETAILS

Scheduled sampling requires generating model-based frames during training, which incurs higher computational costs compared to using only ground-truth frames. In a multi-GPU training setting, naively performing scheduled sampling at random intervals could result in synchronization bottlenecks, as some GPUs would have to wait for others to complete computationally expensive sampling steps. To mitigate this issue, we implemented scheduled sampling at regular intervals across all GPUs simultaneously. Specifically, for a scheduled sampling probability of p=0.05, each GPU (and all examples within each GPU's batch) performs scheduled sampling exactly once every 20 steps. This synchronization approach ensures consistent training speed and prevents slowdown due to inter-GPU blocking.

## L RELATED WORK

NeuralOS is closely related to recent generative modeling approaches for simulating interactive environments conditioned on user inputs. "World Models" (Ha & Schmidhuber, 2018b) introduced latent-variable models for simulating reinforcement learning environments. GameGAN (Kim et al., 2020) used generative adversarial networks (GANs) for interactive game imitation, and Genie (Bruce et al., 2024) generated playable 2D platformer worlds. More recently, diffusion-based models have emerged as powerful real-time simulators: GameNGen (Valevski et al., 2024) simulated the game DOOM, MarioVGG (Protocol, 2024) simulated Super Mario Bros, DIAMOND (Alonso et al., 2024) simulated Atari and Counter-Strike, GameGen-X (Che et al., 2024) simulated open-

world games, and Matrix (Feng et al., 2024) simulated AAA games. Beyond video games, UniSim (Yang et al., 2023) developed simulators for real-world scenarios, and Pandora (Xiang et al., 2024) introduced controllable video generation using natural-language prompts.

Compared to these prior works, NeuralOS addresses distinct challenges unique to OS simulation: while most GUI frame transitions involve subtle changes, accurately modeling critical discrete events, such as opening applications or menus, is essential. Additionally, precise cursor position prediction is crucial for interactive usability. NeuralOS introduces targeted model and training innovations specifically addressing these challenges, paving the way toward fully generative OS simulations.

## M LIMITATIONS

 Our work represents an initial step toward a fully generative OS, but several limitations remain. Despite substantial training compute (17,000 H200 GPU hours and 6,000 H100 GPU hours), NeuralOS is still far from replicating the capabilities of a real operating system: screen resolution remains low, fine-grained keyboard interactions are not reliably supported (Section B). Additionally, many open challenges remain, including enabling the generative OS to interact with external resources (e.g., internet), and introduce controllability beyond traditional OS boundaries.

## N HUMAN EVALUATION

To evaluate how realistic NeuralOS appears to human observers, we conducted a perceptual study modeled after GameNGen (Che et al., 2024). Participants were shown pairs of short video clips—one generated by NeuralOS and one recorded from a real Ubuntu XFCE desktop—corresponding to the same underlying sequence of user interactions. Their task was to identify which clip came from the real operating system.

We evaluated six interaction durations: 10s, 20s, 30s, 40s, 50s, and 60s. To control for artifacts that could reveal the model's identity (such as flickering free-space counters, a known issue in diffusion models), the study was conducted under two conditions: (1) full-frame videos and (2) videos cropped to remove the bottom 40 pixels (10% of the height).

Each participant completed 30 randomized trials per condition. A screenshot of the evaluation interface is provided in Figure 19, and quantitative results are reported in Table 1 (main text). In the cropped condition, participants' accuracy fell close to chance, suggesting that NeuralOS simulations are visually convincing for short interaction sequences.

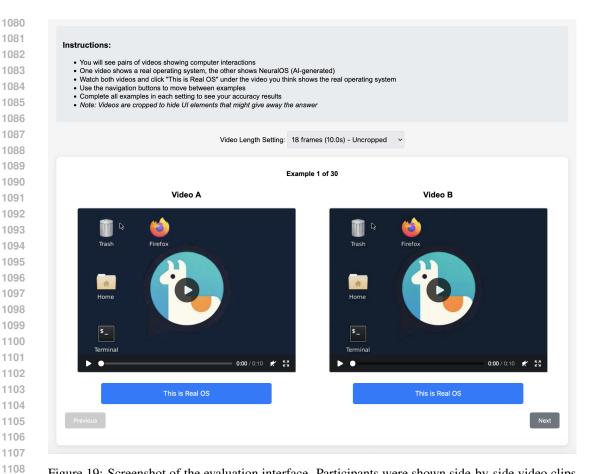


Figure 19: Screenshot of the evaluation interface. Participants were shown side-by-side video clips of the same user interaction sequence, one generated by NeuralOS and one recorded from a real operating system, and asked to identify the real operating system.