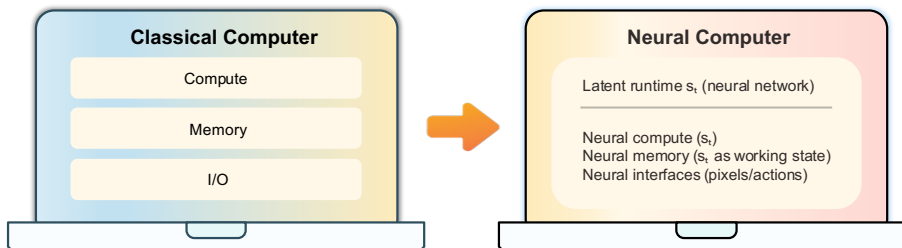


NEURAL COMPUTERS

Anonymous authors

Paper under double-blind review



From modular hardware stack to neural latent stack

ABSTRACT

We train neural networks (NNs) to learn the dynamics of traditional computers (TCs). The NN’s inputs include screenshots, instructions (natural language or program commands), and user actions such as keyboard typing and mouse movements/clicks. Without relying on internal program state, we investigate whether a unified neural system can learn, simulate, and operate a computer interface directly from raw I/O data. We present initial steps towards a completely neural computer (CNC): in our first experiments, our NNs learn to emulate TCs in both CLI and GUI settings. We derive practical design considerations for building more and more general CNCs.

1 INTRODUCTION

Can a single set of weights act as the future “computer”? Recent advances in world models (Ha & Schmidhuber, 2018; OpenAI, 2024; Google DeepMind, 2025) suggest that NNs can internalize rich environment dynamics.

High-capacity video generators such as Veo 3.1 (Google, 2025b) and Sora 2 (OpenAI, 2025) translate conditioning signals into coherent, photorealistic frame sequences. Interactive generators such as Genie 3 (Google, 2025a) extend this to action-conditioned rollouts. In parallel, LLM-driven UI runtimes such as Imagine with Claude¹ map natural language inputs to structured interface updates. Yet these capabilities are typically split across components: models render or predict, while an external runtime (simulator/OS/DOM) carries executable state and mediates control. To fill this void, we define a neural computer (NC) as a learned, executable, and ultimately programmable system in which model, environment, and runtime share a unified latent state. An NC exposes a computer-like interface through pixels, actions, and language.

In current version, we instantiate NCs as video models that renders frames while carrying task context. We study two interfaces. $\text{NC}_{\text{CLIGen}}$ is a terminal interface conditioned on a text (natural language or command lines) and first frame (Section 3.1). $\text{NC}_{\text{GUIWorld}}$ is a desktop interface conditioned on recent pixels and synchronized mouse/keyboard actions (Section 3.2). Teaser contrasts the classical compute/memory/I/O stack (left) with the NC formulation (right), where a single latent runtime s_t jointly (i) models environment dynamics, (ii) acts as working memory, and (iii) renders and controls the interface via pixels and actions. Figure 1 illustrates the design motivation across two interface-specific NCs.

In $\text{NC}_{\text{CLIGen}}$ experiments, the NC can render and execute basic command-line workflows. It often stays aligned with the terminal buffer and captures common “physics” of everyday CLI use (e.g.,

¹<https://claude.ai/imagine/>

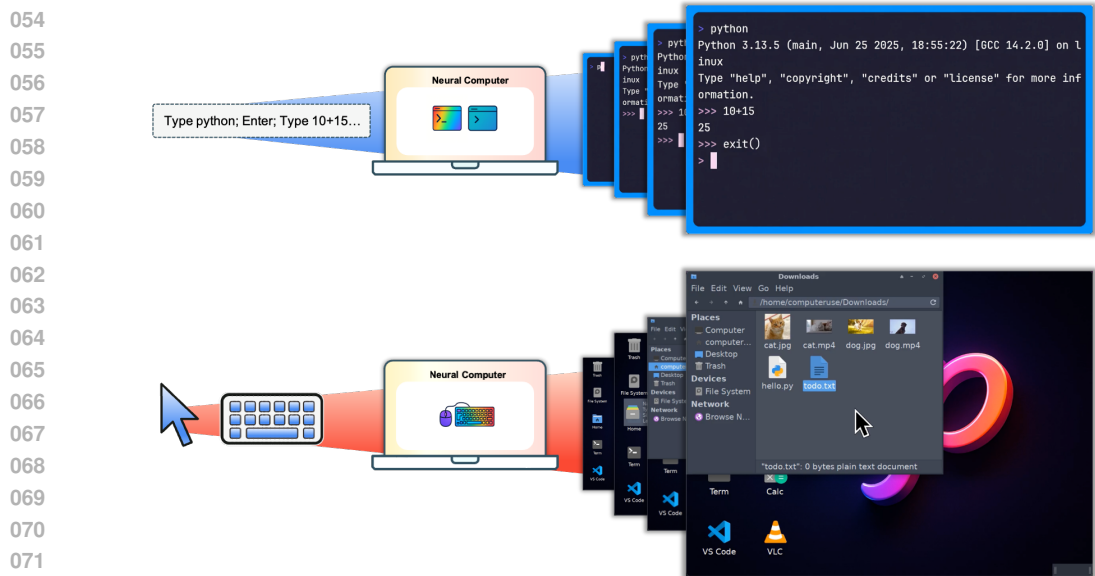

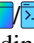

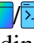


Figure 1: **Neural computers across interfaces.** Given a prompt or action stream, an NC rolls out future interface frames for   NC_{CLIGen} (top) and   NC_{GUIWorld} (bottom). NC_{CLIGen} and NC_{GUIWorld} are the corresponding models trained on those datasets.

fast scrollbar, prompt wrapping, window resizing). Trained on carefully scripted data, generations are visually and structurally close to real sessions. NCs can execute short command chains and render coherent outputs. While symbolic reasoning remains challenging, we show that improved conditioning can dramatically raise arithmetic-probe accuracy (up to 83%; Figure 6).

In NC_{GUIWorld} experiments, we evaluate practical world-model designs across data quality, action injection, and action encoding in this domain. Qualitatively, the model learns coherent pointer dynamics and short-horizon action responses (e.g., hover/click feedback and window/menu transitions). With explicit cursor supervision, it reaches 98.7% cursor accuracy (Table 6).




Overall, our results suggest that current NCs already instantiate several classical interface primitives—most notably accurate I/O alignment and reliable short-horizon control, across both CLI and GUI settings. Beyond visual alignments, we uncover many additional behaviors; two interesting examples are: (1) symbolic reasoning is often bottlenecked by the conditioning interface, but can be significantly improved with clearer, stronger conditions; (2) with explicit visual supervision, cursor tracking approaches perfection, suggesting that precise low-level control is within reach. However, fine-grained text/icon rendering and long-horizon layout consistency remain brittle, pointing to missing primitives for robust computation, memory, and runtime invariance. We therefore define the long-term target, the Completely Neural Computer (CNC), as a fully learned computer whose compute, memory, and interfaces are unified within a single runtime rather than engineered as separate modules. Appendix B summarizes a capability map and an illustrative roadmap and more insights. **Concretely, this work makes the following contributions:**

- Define neural computers (NCs) and build video-based prototypes for both CLI & GUI interfaces.
- Provide a data engine and alignment recipe that synchronize text, actions, and frames for CLI & GUI environments, from which we collect more than 3000 hours of training data.
- We identify key design choices for NCs through extensive ablation studies (spending more than 200k H100 hours).
- Outline a roadmap toward completely neural computers (CNCs) grounded in classical computer primitives and open research questions.

2 PRELIMINARIES

Classical (von Neumann-style) computers separate computation, memory, and I/O, typically realized as a layered operating-system stack. Our motivating question is whether a single set of weights



Figure 2: **Data types for learning NC behaviors.** Logos denote datasets:  CLIgen (General) replays public asciinema traces from diverse terminal workflows.  CLIgen (Clean) uses scripted vhs runs for deterministic traces in controlled experiments.  GUIWorld records desktop RGB with mouse/keyboard traces to validate action-conditioned rendering and control on GUIs.

can internalize these roles inside one stateful latent runtime, rather than relying on an external OS/simulator to carry executable state. We model a video-based neural computer (NC) as a learned latent-state system that folds these roles into an update-and-render loop.

In current designs of this paper, an NC is an autoregressive video model or world model (Ha & Schmidhuber, 2018) of computer systems that generates observables conditioned on given context. Specifically, an NC accepts an observable and a user action as input to predict the subsequent observable. We treat screen frames as observables and define actions as text prompts or control trajectories. More broadly, the NC framework can accommodate various other modalities and structural representations for both observables and actions. Given an initial screen frame x_0 and conditioned on user action u_t at iteration t , an NC updates its memory state and samples the next frame x_{t+1} . Formally, an NC defined by an initial memory state s_0 , an update function F_θ , and a decoder G_θ operates as follows:

$$s_t = F_\theta(s_{t-1}, x_t, u_t), \quad x_{t+1} \sim G_\theta(s_t). \tag{1}$$

In the context of world modeling, the sequence of actions is referred to as a conditioning stream.

Under this formulation, each step collapses multiple subsystems into one learned update. The state s_t stores all relevant context. F_θ integrates new observations and actions, and G_θ renders the next frame. Auxiliary heads can encode and decode prompts, buffers, or action traces. This moves some functionality that would traditionally live in OS queues, device drivers, and UI toolkits into the latent runtime.

2.1 RELATED WORK

Differentiable world models (Schmidhuber, 1990; 2015) sought to learn representations of environment dynamics. In parallel, neuromorphic designs (Mead & Ismail, 2012) explored neural computation as a substrate. Differentiable memory and program-execution architectures, including Neural Turing Machines (Graves et al., 2014), Differentiable Neural Computers (Graves et al., 2016), and Neural Programmer-Interpreters (Reed & De Freitas, 2015), showed that neural controllers with memory can execute structured procedures. Latent video and world models (Ha & Schmidhuber, 2018; Hafner et al., 2019b;a; Bruce et al., 2024) extend these ideas to embodied control and interactive environments. Genie 3 (Google, 2025a), in particular, frames such models as agent-training substrates with improved physical consistency. High-capacity generators such as Veo 3 (Google, 2025b) and Sora 2 (OpenAI, 2025) emphasize open-ended, photorealistic simulation. Systems such as NeuralOS (Rivard et al., 2025) and Imagine with Claude (Anthropic, 2025) bring world-model

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

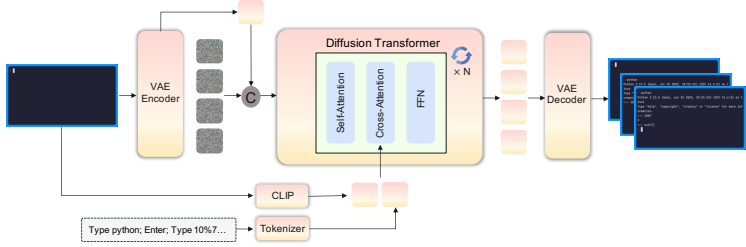


Figure 3: / NC_{CLI}Gen architecture. Terminal frames are observations x_t . A prompt and the first frame seed the conditioning stream. The Wan2.1-based latent state z_t rolls forward under the standard I2V sampling scheme.

conditioning to desktop and DOM-style interfaces. In this work, we study two NC instantiations for CLI and GUI with interface-specific conditioning, and a staged roadmap toward CNCs.

3 IMPLEMENTATION OF NEURAL COMPUTERS

We build on the Wan2.1 model (Wan et al., 2025), which was a SoTA video generation model at the time of our experiments. We add NC-specific conditioning and action modules, together with interface-specific training recipes for our specific settings. Figure 1 illustrates the setups. NCs take a prompt or action stream as input and generate future frames for both NC_{CLI}Gen and NC_{GUI}World.

We treat the video model’s time-indexed latent representations as the state z_t in Equation (1). This abstracts the underlying diffusion transformer as a latent dynamical system. The transformer F_θ maps prior context and current observations/conditioning into an updated state. The decoder G_θ predicts the distribution of the next frame \hat{x}_{t+1} . Auxiliary heads encode and decode conditioning streams u_t , including text prompts, terminal buffers, and action traces.

3.1 / THE CLI VIDEO GENERATORS

Our CLI video generators instantiate the NC in terminal settings. Observations x_t are terminal frames rendered from the underlying text buffer. The conditioning stream u_t carries a user prompt and optional metadata, and the latent state z_t tracks CLI context across frames. At inference time, the model rolls out from the prompt and first frame, updates z_t , and predicts future terminal frames (Figure 3). Throughout this section, CLIGen (General) and CLIGen (Clean) refer to datasets. We denote the corresponding models trained on each dataset as NC_{CLI}Gen and NC_{CLI}Gen.

3.1.1 DATA PIPELINE

We collect CLI data under two complementary regimes (details in Section E.1). CLIGen (General) replays public `asciinema` traces to cover diverse real-world terminal workflows, while CLIGen (Clean) uses scripted `vhs` runs in Dockerized environments to produce deterministic traces for controlled experiments (including math).

3.1.2 MODEL ARCHITECTURE

We treat CLI generation as text-and-image-to-video: a caption and the first terminal frame condition the rollout. The first frame is encoded by a VAE into a conditioning latent. In parallel, a CLIP image encoder (Radford et al., 2021) extracts visual features from the same frame, and a text encoder (e.g., T5 (Raffel et al., 2020)) embeds the caption. Following the Wan2.1 image-to-video (I2V) design, these conditioning features are concatenated with diffusion noise, projected through a zero-initialized linear layer, and processed by a DiT stack. Decoupled cross-attention injects the joint caption and first-frame context derived from the CLIP and text features. The VAE encodes and decodes terminal frames. During generation, the diffusion transformer advances the latent state z_t under the original Wan2.1 I2V sampling schedule, without additional binary masks or periodic reseeded. Training details are provided in Appendix D.

Table 1: Caption styles versus generation fidelity.

Prompt style	PSNR	SSIM	Avg. words
Semantic	21.90	0.813	55
Regular	23.63	0.843	52
Detailed	26.89	0.867	76

3.1.3 EVALUATIONS

We evaluate (i) optical fidelity (PSNR/SSIM/LPIPS) and (ii) trace-level interaction correctness (action→state consistency, command success, buffer/text consistency), *independent of OCR*. Beyond just visual fidelity, non-optical metrics rigorously validate the robustness, accuracy, and consistency of multi-step workflows, clearly demonstrating that our models not only replicate visuals but also handle complex, real-world tasks with high precision.

Experiment 1: Generic VAEs handle terminal content with proper font sizing

Concurrent work (Rivard et al., 2025) argues that generic natural-image VAEs perform poorly on structured computer screenshots. We evaluate this claim directly by applying the Wan2.1 VAE (Wan et al., 2025) to terminal content. In our setting, reconstruction quality is primarily governed by font size. At 13 px fonts it is high (40.77 dB PSNR, 0.989 SSIM). At 6 px fonts, text exhibits noticeable blurring.

However, a sweep over CLIGen (General) frames shows that this effect is confined to extreme cases (Figure 4). Very small 6 px fonts and ultra-dense text exhibit localized blurring despite high global PSNR. In contrast, the 13 px terminal font used in CLIGen remains visually sharp across panes and commands. These results indicate that the VAE is adequate for regular CLIGen usage and highlight that sensible font choices help ensure stable, legible NC training.

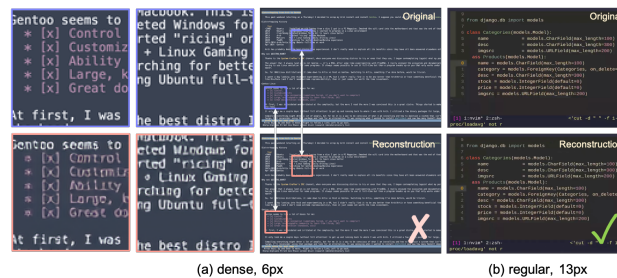


Figure 4: Wan2.1 VAE reconstructions on CLIGen (General).

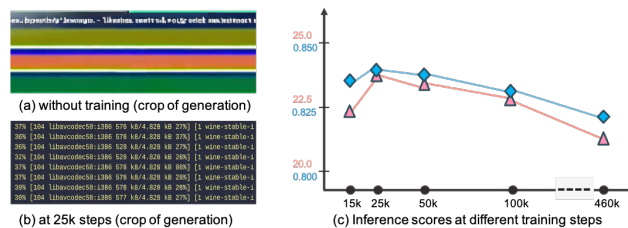


Figure 5: CLIGen (General) training plateaus at 25k steps.

Experiment 2: The impact of data quality

In CLIGen (general) experiments, data quality matters more than training duration, this is different from natural image generation can continue to benefit from scale and prolonged training. After early gains, remaining errors are dominated by artifact-prone supervision (e.g., rendering glitches or rapid screen changes that break stable temporal alignment).

A visual comparison in panels (a–b) highlights the effect of training on CLIGen data. Without CLIGen data, Wan2.1 produces largely garbled terminal outputs (panel a). After 25k training steps, the model generates readable text with consistent formatting and color codes (panel b).

Table 2: OCR accuracy versus training.

Steps (k)	Char. acc.	Exact line
0	0.03	0.01
10	0.18 $\uparrow_{0.15}$	0.05 $\uparrow_{0.04}$
20	0.33 $\uparrow_{0.30}$	0.12 $\uparrow_{0.11}$
30	0.41 $\uparrow_{0.38}$	0.18 $\uparrow_{0.17}$
40	0.52 $\uparrow_{0.49}$	0.26 $\uparrow_{0.25}$
50	0.52 $\uparrow_{0.49}$	0.27 $\uparrow_{0.26}$
60	0.54 $\uparrow_{0.51}$	0.31 $\uparrow_{0.30}$

Table 3: Closed-book arithmetic-probe accuracy.

Model	Accuracy
Wan2.1	0%
NC _{CLIGen}	4%
Veo3.1	2%
Sora2	71%

Figure 5 plots the corresponding PSNR/SSIM training curves and reveals a clear efficiency ceiling. These reconstruction metrics peak at approximately 25k steps. They do not improve with further training up to 460k steps. Extended training can even slightly degrade performance. We hypothesize that terminal rendering reaches a natural learning limit. Most structured patterns are acquired early, and further gains would require higher-quality, better-paced, or more informative data.

Experiment 3: Literal captions boost generations

Caption specificity strongly affects terminal rendering quality. We compare 3 caption tiers (semantic/regular/detailed) and find detailed, literal descriptions improve reconstruction fidelity (PSNR 21.90→26.89 dB; Table 1). We attribute this to terminals dominated by exact text placement, where literal captions constrain token text-to-pixel alignment.

Experiment 4: Neural computers achieve accurate character-level text generation

PSNR and SSIM capture perceptual similarity, but character-level accuracy is more direct for terminal rendering. It requires pixel-to-text correspondence. For CLIGen (Clean), we apply Tesseract OCR to sampled frames (with whitespace normalization). We report character accuracy (Levenshtein distance) and exact-line match rate (full protocol in Appendix E).

Table 2 shows that our models achieve strong text rendering accuracy. Character accuracy increases from 0.03 at initialization to 0.54 at 60k steps. Exact-line matches reach 0.31 (0.26 at 40k). Most gains occur within the first 40k steps, with smaller improvements afterward. These OCR-based results go beyond perceptual similarity. Accurately generating terminal characters requires modeling text structure, font rendering, and spatial relationships. These are core competencies for interactive neural systems. This level of precision brings us closer to usable terminal interfaces, not just plausible ones.

Experiment 5: Closed-book arithmetic requires improvement for better reasoning

Beyond rendering quality, we probe symbolic computation using CLI arithmetic tasks. In the closed-book setting, the conditioning text specifies the expression but not the solution. This probe is intentionally strict: the model must preserve REPL structure, keep the typed expression intact, and render the correct numeric output. We reserve a held-out pool of 1,000 math problems and randomly sample 100 for evaluation. Table 3 shows that most current video models struggle on these tasks. Wan2.1 achieves 0% accuracy, our NC_{CLIGen} model reaches 4%, and Veo3.1 manages 2%, all far below human-level performance. Sora2’s 71% accuracy is a notable outlier and may reflect system-level advantages or additional training beyond standard video generation. Overall, native symbolic computation remains an open challenge for pixel-grounded NCs.

The poor arithmetic-probe performance in Table 3 raises a practical question. Do such gains require specialized reinforcement learning, or can system-level conditioning substantially narrow this gap? Sora2’s 71% accuracy suggests several possible explanations, which we examine in Table 4. Our explorations (in Experiment 6) yield an interesting result. Better conditioning alone can dramatically

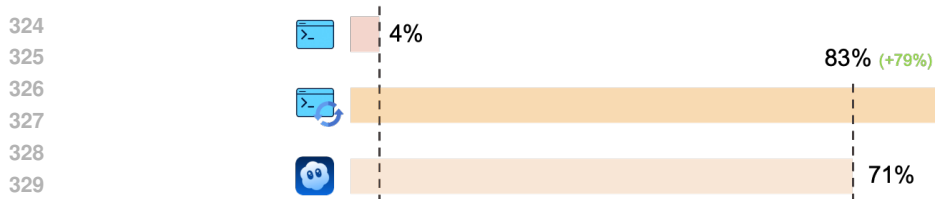


Figure 6: Reprompting boosts performance to 83%.

improve arithmetic-probe accuracy, without introducing reinforcement learning or modifying the video generator.

Experiment 6: Reprompting helps improve arithmetic-probe scores

Table 4: Why Sora2 leads (our hypotheses)?

Factor	Implication
1 Stronger base + similar data	Higher intrinsic arithmetic; symbolic capability may be baked in
2 Additional RL training	Reward shaping teaches math beyond diffusion; could transfer to CLI
3 System-level reprompt/recaption	LLM computes answers; strong conditioning drives generation

As shown in Figure 6, NC_{CLIGen} accuracy on arithmetic tasks rises from 4% to 83% under reprompting, suggesting that system-level conditioning is an effective first step for improving performance on symbolic probes. This is complementary to (not strictly requiring) RL-based pipelines.

Reprompting highlights the sensitivity of symbolic-probe outcomes to the conditioning interface. Much of the “reasoning” gain comes from better specification and instruction-following, not new computation. In the reprompted setting, the conditioning text provides the solution string (examples in Appendix A). To ensure faithful transcription, we include answer-augmented captions for half of the arithmetic training episodes. This demonstrates that the model can consistently render conditioned symbolic content. Our experiments do not support the view that current SoTA (open-sourced or closed sourced) video models already have strong reasoning skills. Simple tasks like two-digit addition are often unsolved correctly, while small LLMs (e.g., 100MB LLMs) easily handle them. However, after reprompting with an LLM-provided answer, accuracy improves, suggesting video models have potential but are not yet strong reasoning systems. This suggests that while video models have potential, they are not yet strong reasoning systems. We see this as similar to the development of LLMs in solving simple questions like “How many ‘r’s are in ‘strawberry’” (which humans can answer instantly, but LLMs easily failed). We believe these limitations are due to the current development stage, and expect stronger reasoning abilities in future NCs.

The evidence supports system-level conditioning as a practical path forward. Among hypotheses for improving arithmetic-probe performance, stronger base models, reinforcement learning, or enhanced conditioning—our results indicate that the latter is a convenient path. The gain from reprompting (4%→83%), achieved without modifying the video generator, shows that “reasoning” in these probes is highly sensitive to conditioning. Instead of complex reinforcement learning pipelines, large symbolic-probe gains can be achieved through strategic conditioning alone. Evaluations should distinguish between native computation and conditioning-assisted performance when assessing video models’ reasoning capabilities.

3.1.4 VISUALIZATIONS

Qualitative visualization pages are moved to Appendix A. See Appendix A for CLIGen (General/Clean) rollouts, REPL traces, and math-probe visualizations.

3.2 THE GUI WORLD MODELS

We extend the neural computer framework to interactive desktop environments with NC_{GUIWorld} . Here, fine-grained action control is essential: GUI interaction requires precise cursor tracking, timely click feedback, and robustness to rapidly changing interface states.

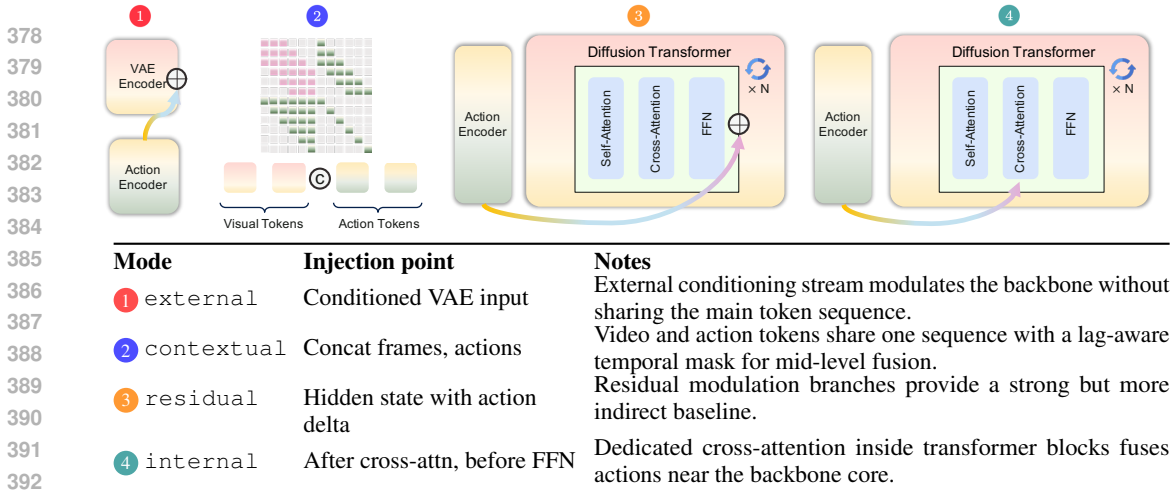


Figure 7: **Four modes for injecting GUI actions into the diffusion transformer.** 1 modulates VAE latents before the transformer; 2 appends action tokens; 3 applies residual modulation; and 4 inserts action cross-attention within blocks.

3.2.1 DATA PIPELINE

The system records synchronized RGB video frames (x_t) together with detailed mouse and keyboard event streams (u_t). This preserves pixel-level alignment between user actions and visual responses. Data collection and preprocessing details for GUIWorld are provided in Appendix E.

3.2.2 MODEL ARCHITECTURE

The GUIWorld architecture builds on the Wan2.1 (Wan et al., 2025) by incorporating explicit action-conditioning modules. The central challenge is to align time-stamped user actions with generated frames and inject this information at the appropriate depth within the transformer. Action features are precomputed from frame-aligned mouse and keyboard signals (Appendix E.4). We aggregate them into latent-aligned embeddings that summarize recent action history at each diffusion step. We evaluate two action encoders. The *raw-action encoder* (v1) preserves fine-grained mouse/keyboard event streams. The *meta-action encoder* (v2) abstracts interactions into coarse API-style categories (clicks, drags, scrolls, typing, shortcuts). Both encoders use the same temporal alignment and are evaluated as separate ablations. In our experiments, their effects on rendering fidelity and control behavior are comparable.²

We inject action embeddings into the diffusion backbone in four ways (Figure 7). We study external, contextual, residual, and internal conditioning. For the injection-scheme ablation, all four modes share the same meta-action encoder and temporal alignment. They differ only in where the latent action features interact with the video latents and transformer blocks. We compare raw-action vs. meta-action encoders separately in Table 8 (with additional details in Appendix G).

external conditioning. In the external mode, action information modulates the latent video sequence before the diffusion transformer. Action features are applied as a pre-conditioning step at the model input, without introducing explicit action tokens or cross-attention inside the diffusion backbone. As a result, action information enters only through the modified input latents; the diffusion backbone never attends directly to action tokens, so any action signal must be carried implicitly in $z'_{1:T}$.

Concretely, an external action module produces a residual update to the VAE latents, yielding modified latents $z'_{1:T}$ that are then processed by the unchanged diffusion transformer (formal definitions in Appendix G.2).

contextual conditioning. In the contextual mode, actions are represented as additional tokens and integrated directly into the transformer’s self-attention. Similar token-based action representations

²Table 14 summarizes the representational differences. Training recipes are in Appendix D

Table 5: Overall performance across data sources.

Split	FVD _{all}	SSIM _{all}	LPIPS _{all}
baseline	149.61	0.496	0.605
train with random fast	48.17	0.695	0.483
train with random slow	20.37	0.830	0.237
train with Claude CUA	14.72	0.885	0.144

Table 6: Cursor conditioning losses versus accuracy.

Loss variant	Cursor accuracy
Position (x, y) only	8.7%
Position (x, y) + Fourier	13.5%
Position (x, y) + SVG mask/ref	98.7%

have been explored in prior world models, including Gato (Reed et al., 2022) and World and Human Action Models (Kanervisto et al., 2025).

The meta-action encoder produces latent-aligned action tokens $A \in \mathbb{R}^{L_a \times D}$. We concatenate them with visual tokens $V \in \mathbb{R}^{L_v \times D}$ to form a joint sequence $[V; A]$. Each transformer block applies self-attention over this combined sequence using a structured temporal mask (Appendix Figure 38). The mask enforces causal alignment: each frame token attends only to actions within a short past window, and each action token attends only to frames after a fixed temporal lag. Through this masked joint attention, contextual conditioning fuses action and visual information within the transformer blocks.

residual conditioning. In the residual mode, the transformer block structure remains unchanged. A lightweight action module attaches to a subset of layers as a residual branch. This follows the residual conditioning paradigm from ControlNet (Zhang et al., 2023), while remaining modular and additive to the base diffusion backbone.

At each selected layer l , the transformer applies its standard sequence of self-attention, text or reference cross-attention, and feed-forward operations to produce hidden states $h^{(l)}$. A separate action module then takes $h^{(l)}$ together with a local temporal window of latent action features and mouse trajectories. It outputs a residual update $\Delta h^{(l)}(a, \text{mouse})$. This update is added to the block hidden states and passed to the subsequent transformer block (formal definitions in Appendix G.2). In this formulation, residual conditioning injects action information through block-external residual branches. It does not modify the internal computations of the transformer blocks themselves.

internal conditioning. In the internal mode, action conditioning is incorporated within the transformer blocks. Related multi-stream world models have explored similar designs, including Matrix-Game-2 (He et al., 2025). Each selected block augments the standard attention stack with an additional action cross-attention sub-layer. Specifically, the block applies self-attention, followed by cross-attention over text and reference features, and then a dedicated action cross-attention layer. Keys and values are derived from latent action features (and, optionally, mouse inputs).

Given block input h , text or reference context c , and action latents a , the internal block computes an update that inserts an additional action cross-attention term inside the block (formal definitions in Appendix G.2). As illustrated in Figure 7, action features are injected directly into the block’s cross-attention stage.

In contrast to residual conditioning, internal conditioning integrates action information through a block-internal attention mechanism rather than an external residual branch. This design mirrors the multi-stream injection strategy used in Matrix-Game-2 (He et al., 2025) and provides the strongest control for fine-grained GUI interaction. In this setting, precise temporal alignment and spatial locality are critical. Each conditioning mode (external, contextual, residual, and internal) is trained as a separate ablation, and no combinations are used.

3.2.3 EVALUATION SETUP

We use the FVD/LPIPS/SSIM suite as core metrics, adding action-driven metrics that focus on post-interaction frames after clicks, scrolls, and key/type events. For example, we compute SSIM/LPIPS averaged over the k frames after each logged action, and action-driven FVD on

Table 7: Action-driven metrics (15 frames post-action) across 4 schemes. † baseline₂: external, early-stopped at ~50% steps.

Mode	SSIM ₊₁₅ ↑	LPIPS ₊₁₅ ↓	FVD ₊₁₅ ↓
baseline ₁ (untrained)	0.326	0.649	184.3
†baseline ₂	0.746	0.251	33.4
contextual	0.813	0.190	24.8
residual	0.857	0.138	18.8
internal	0.863	0.141	14.5

Table 8: Raw-action vs. API-like action encoding under the same injection mode (15 frames after action).

Mode	Encoding	SSIM ₊₁₅ ↑	LPIPS ₊₁₅ ↓	FVD ₊₁₅ ↓
internal	raw-action (event-stream)	0.847	0.144	16.6
internal	meta-action (API-like)	0.863	0.141	14.5

post-action clips. Ablations vary conditioning design and action encoding to measure how these choices affect perceptual quality and responsiveness against ground-truth interfaces. Full metrics and implementations are in Appendix E.6.

🖱️ Experiment 7: Data quality matters

Interactive GUI modeling shows that data quality matters more than dataset size for action-driven performance. We compare slow exploration, fast interaction, and supervised trajectories under contextual conditioning. This isolates which behaviors best support neural computer training.

Despite approximately 1,400 hours of random exploration across the slow and fast settings, these datasets are noisy. They are comparatively sample-inefficient for learning stable action–response mappings. They substantially improve global perceptual metrics over a baseline (Table 5). However, high-frequency cursor jitter and irregular, non-goal-directed action bursts make consistent control difficult under dense, stochastic input streams.

In contrast, the substantially smaller high-quality dataset (110 hours from Claude CUA) yields markedly stronger performance across all metrics. Goal-directed trajectories provide clearer action semantics and more predictable state transitions. This enables robust action conditioning even with limited data volume. These results indicate that neural computer development should prioritize curated, purposeful interactions over large-scale passive data collection.

🖱️ Experiment 8: Precise cursor control requires explicit visual supervision

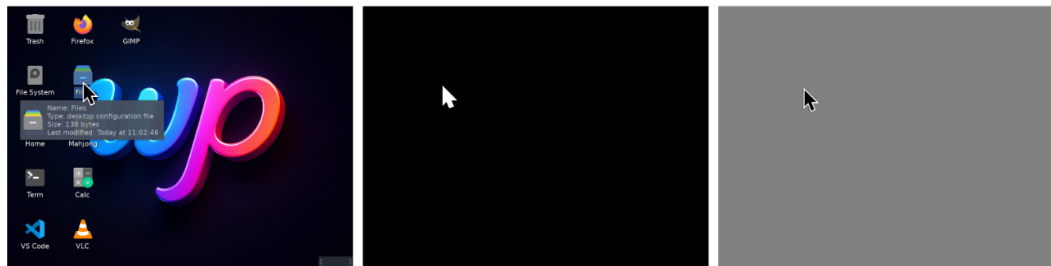


Figure 8: **Cursor references in GUIWorld.** **Left:** original desktop frames. **Middle:** binary cursor masks. **Right:** cursor-only reference frames rendered over a neutral background.

We examine whether the NC internalizes cursor dynamics from logged coordinates. Details of coordinate normalization and trajectory encoding are provided in Appendix E.4.

However, Table 6 shows that coordinate-based supervision remains insufficient for precise interaction. Position-only supervision achieves 8.7% accuracy, and even enhanced position features reach only 13.5%. This suggests that richer coordinate encodings alone do not resolve cursor drift and jitter. Motivated by the need for precise cursor placement, we add explicit visual cursor supervision by rendering a cursor-only foreground/mask reference stream and providing it as an additional conditioning signal during training (Appendix G.4). Under this explicit visual conditioning, cursor

540 accuracy improves to 98.7%. This suggests that neural computers benefit from learning the cursor state
541 as a visual object rather than relying solely on abstract coordinates. Explicit pixel-level supervision
542 helps model cursor acceleration, hover states, and click feedback, which are essential for reliable
543 GUI interaction.

544
545  **Experiment 9: Different Action injections**

546 Holding data and the action encoder fixed, deeper injection improves post-action fidelity (Ta-
547 ble 7): `contextual/residual/internal` outperform `external` and the untrained baseline.
548 `internal` yields the best SSIM/FVD and `residual` yields the best LPIPS; details of injection
549 schemes are in Appendix G.

550
551  **Experiment 10: Are action encodings influenced?**

552
553 Under a fixed injection mode, raw-action and meta-action encodings perform similarly (Table 8); we
554 use meta-action encoder and provide additional analysis in Appendix G.

555
556 **4 CONCLUSION**

557
558 Neural Computers learn interface-level behavior from raw CLI and GUI interaction traces, producing
559 models that render coherent screen updates and support multi-step interactions. We provide practical
560 guidance on data, conditioning, and evaluation to improve fidelity and control; a broader CNC
561 roadmap is in Appendix B.

562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

REFERENCES

- 594
595
596 Anthropic. Introducing claude sonnet 4.5. <https://www.anthropic.com/news/claude-sonnet-4-5>, September 2025.
597
- 598 Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes,
599 Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative
600 interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.
601
- 602 Diego Calanzone, Stefano Teso, and Antonio Vergari. Logically consistent language models via neuro-
603 symbolic integration. In *The Thirteenth International Conference on Learning Representations*,
604 2025.
- 605 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
606 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
607 *arXiv preprint arXiv:2407.21783*, 2024.
- 608 Kunihiko Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of
609 pattern recognition unaffected by shift in position. *Biological cybernetics*, 36(4):193–202, 1980.
610
- 611 Google. Genie 3: A new frontier for world models, 2025a. URL <https://deepmind.google/blog/genie-3-a-new-frontier-for-world-models/>. Published August 5, 2025.
612
- 613 Google. Break the silence with veo 3.1. <https://gemini.google/overview/video-generation/>, November 2025b.
614
- 615 Google DeepMind. Veo. <https://deepmind.google/models/veo/>, May 2025.
616
- 617 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint*
618 *arXiv:1410.5401*, 2014.
619
- 620 Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-
621 Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, et al.
622 Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626):
623 471–476, 2016.
- 624 David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
625
- 626 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
627 behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019a.
628
- 629 Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James
630 Davidson. Learning latent dynamics for planning from pixels. In *International conference on*
631 *machine learning*, pp. 2555–2565. PMLR, 2019b.
- 632 Xianglong He, Chunli Peng, Zexiang Liu, Boyang Wang, Yifan Zhang, Qi Cui, Fei Kang, Biao Jiang,
633 Mengyin An, Yangyang Ren, et al. Matrix-game 2.0: An open-source, real-time, and streaming
634 interactive world model. *arXiv preprint arXiv:2508.13009*, 2025.
- 635 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):
636 1735–1780, 1997.
637
- 638 Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao
639 Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta programming for a multi-
640 agent collaborative framework. In *The twelfth international conference on learning representations*,
641 2023.
- 642 Mike Innes, Alan Edelman, Keno Fischer, Chris Rackauckas, Elliot Saba, Viral B Shah, and Will
643 Tebbutt. A differentiable programming system to bridge machine learning and scientific computing.
644 *arXiv preprint arXiv:1907.07587*, 2019.
645
- 646 Anssi Kanervisto, Dave Bignell, Linda Yilin Wen, Martin Grayson, Raluca Georgescu, Sergio
647 Valcarcel Macua, Shan Zheng Tan, Tabish Rashid, Tim Pearce, Yuhan Cao, et al. World and human
action models towards gameplay ideation. *Nature*, 638(8051):656–663, 2025.

- 648 Diederik Kinga, Jimmy Ba Adam, et al. A method for stochastic optimization. In *International*
649 *conference on learning representations (ICLR)*, volume 5. California, 2015.
- 650
- 651 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A
652 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming
653 catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114
654 (13):3521–3526, 2017.
- 655 Carver Mead and Mohammed Ismail. *Analog VLSI implementation of neural systems*, volume 80.
656 Springer Science & Business Media, 2012.
- 657
- 658 OpenAI. Sora by openai. <https://openai.com/sora/>, 2024. Accessed: 2025-07-14.
- 659 OpenAI. Sora 2 is here. <https://openai.com/index/sora-2/>, September 2025.
- 660
- 661 Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is turing-complete. *Journal of Machine*
662 *Learning Research*, 22(75):1–35, 2021.
- 663 Matthieu Queloz. Explainability through systematicity: The hard systematicity challenge for artificial
664 intelligence. *Minds and machines*, 35(3):35, 2025.
- 665
- 666 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
667 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
668 models from natural language supervision. In *International conference on machine learning*, pp.
669 8748–8763. PmLR, 2021.
- 670 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
671 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
672 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 673
- 674 Scott Reed and Nando De Freitas. Neural programmer-interpreters. *arXiv preprint arXiv:1511.06279*,
675 2015.
- 676 Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel
677 Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist
678 agent. *arXiv preprint arXiv:2205.06175*, 2022.
- 679
- 680 Luke Rivard, Sun Sun, Hongyu Guo, Wenhui Chen, and Yuntian Deng. Neuralos: Towards simulating
681 operating systems via neural generative models. *arXiv preprint arXiv:2507.08800*, 2025.
- 682 Jürgen Schmidhuber. *Making the world differentiable: on using self supervised fully recurrent neural*
683 *networks for dynamic reinforcement learning and planning in non-stationary environments*, volume
684 126. Inst. für Informatik, 1990.
- 685
- 686 Jürgen Schmidhuber. On learning to think: Algorithmic information theory for novel combina-
687 tions of reinforcement learning controllers and recurrent neural world models. *arXiv preprint*
688 *arXiv:1511.09249*, 2015.
- 689 Hava T Siegelmann and Eduardo D Sontag. On the computational power of neural nets. In
690 *Proceedings of the fifth annual workshop on Computational learning theory*, pp. 440–449, 1992.
- 691
- 692 Weihao Tan, Wentao Zhang, Xinrun Xu, Haochong Xia, Ziluo Ding, Boyu Li, Bohan Zhou, Junpeng
693 Yue, Jiechuan Jiang, Yewen Li, et al. Cradle: Empowering foundation agents towards general
694 computer control. *arXiv preprint arXiv:2403.03186*, 2024.
- 695 Alan Mathison Turing et al. On computable numbers, with an application to the entscheidungsproblem.
696 *J. of Math.*, 58(345–363):5, 1936.
- 697
- 698 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
699 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
700 *systems*, 30, 2017.
- 701 John Von Neumann. First draft of a report on the edvac. *IEEE Annals of the History of Computing*,
15(4):27–75, 1993.

702 Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu,
703 Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models.
704 *arXiv preprint arXiv:2503.20314*, 2025.
705

706 Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane Gu, Nick Matarese, Kevin Swersky,
707 Been Kim, Priyank Jaini, and Robert Geirhos. Video models are zero-shot learners and reasoners.
708 *arXiv preprint arXiv:2509.20328*, 2025.

709 Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber.
710 Natural evolution strategies. *The Journal of Machine Learning Research*, 15(1):949–980, 2014.
711

712 Maurice Wilkes. The best way to design an automatic calculating machine. 1981.

713 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
714 diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*,
715 pp. 3836–3847, 2023.
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

Appendix

Overview. This submission is long because the project explores an extended chain of ideas that are new to this community. We design the full data-collection pipeline, train seven models in total (three CLIGen and four GUIWorld), and report extensive ablations, qualitative analyses, and system-level insights. We also introduce the definition of, and potential solutions toward, future completely neural computers (CNCs). The appendix is not filler. It contains material we consider essential for understanding the design trade-offs and the main paper’s conclusions. We would greatly appreciate careful attention to the appendix when summarizing or reviewing the work, especially if LLM-based tools are used for condensation.


Table of Contents


Appendix A: Visualizations	16
Appendix B: Toward Completely Neural Computers	46
Appendix C: Explorations: Alternative Data Sources and Online Interaction	49
Appendix D: Implementation Details	51
Appendix E: Datasets: Collection and Evaluation Protocols	52
Appendix F: CLIGen: CLI Trajectory Formats	58
Appendix G: GUIWorld: Action Representation, Conditioning and More Experiments	60

A VISUALIZATIONS (WE KINDLY SUGGEST SCROLLING DOWN TO THE BOTTOM OF THE PDF)


A.1 CLIGEN VISUALIZATIONS

This subsection consolidates CLIGen visualization pages referenced in the main text and moves them here for space, with additional examples provided below. Each page pairs a scripted input trace—typed strings, explicit `Enter` events, and controlled sleep intervals—with the corresponding frame sequence, enabling fine-grained inspection of timing, cursor placement, and incremental output formation.

(1)  **CLIGen (General) visualizations.** Qualitative samples highlight the breadth of real-world terminal dynamics captured in CLIGen (General): ANSI escape sequences that repaint regions with changing foreground/background colors, incremental command entry with syntax highlighting and cursor edits (e.g., `SQL CREATE TABLE`), classic shell prompts and system outputs (e.g., `date`), long-running jobs with rapidly scrolling and color-coded package logs, full-screen TUIs such as partition editors with menus and tables, and progress dashboards with updating bars, counts, and ETAs. These traces emphasize that “looking correct” requires maintaining terminal geometry, palette transitions, and cursor state frame-by-frame.


(2)  **CLIGen (Clean) REPL visualizations.** In contrast to open-world traces, CLIGen (Clean) REPL samples are scripted and temporally well-paced (Figures 12–13; additional examples are provided in Appendix F). Each sample includes an explicit action trace (e.g., `Sleep`, `Type`, `Enter`, arrow keys, `Hide`) alongside the rendered terminal frames, making the causal link between key-level actions and pixel updates visually unambiguous. The rollouts emphasize incremental typing under a stable prompt (here, a `>` prompt with a pink block cursor), with deterministic outputs such as printing `$HOME` (yielding `/home/NeuralComputer`), inspecting environment variables via `env | head -n 5` (showing fields like `DEBIAN_FRONTEND=noninteractive` and `UTF-8 locale`), and composing longer pipelines (e.g., `seq+paste+column+tee`) that materialize small aligned tables. They also include dense TUI-like screens where layout fidelity is critical: `top` must preserve column alignment and header/footer bands, while `cal` must keep a fixed-width calendar grid with consistent spacing. Finally, multi-line programming snippets (a `python - <<'PY' heredoc`) stress quoting, indentation, and line wrapping, and the resulting output produces repeated, structured progress-like traces (“Frame . . .” lines with repeated glyphs).


The main insight is that these scripted REPL traces isolate rendering-and-control errors from semantic ambiguity: given an explicit action stream, failures are dominated by low-level mechanics (cursor placement, character edits, monospace alignment, line breaks, and temporal consistency) rather than uncertainty about what content should be produced. This makes CLIGen (Clean) particularly useful for diagnosing whether a neural computer can faithfully execute action-conditioned text rendering under strict layout constraints.


(3)  **CLIGen (Clean) math visualizations.** Figures 14–15 compare model rollouts on CLIGen (Clean) math REPL prompts, where a scripted action trace launches `python` and types a short expression (e.g., `5, 10+15, 40/1`); additional examples are provided in Appendix F. Sora2 most consistently preserves the intended REPL structure: it renders the Python banner, maintains a stable `>>>` prompt, and produces the expected numeric output with minimal drift. In contrast, Wan2.1 frequently fails at the interface level, producing corrupted prompts, spurious control-token text (e.g., literal `<Enter>`), or collapsing into non-terminal visuals, which aligns with its near-zero arithmetic-probe accuracy. `NCCLIGen` and `Veo3.1` occupy a middle regime: they often keep a terminal-like layout and can render plausible outputs, but the visualizations highlight remaining brittleness in character-level fidelity and mode stability (e.g., prompt corruption, stray glyphs, or occasional off-domain imagery).

Figures 16–17 further illustrate why these math probes do not cleanly measure “native” reasoning; additional examples are provided in Appendix F. In the reprompted setting, we explicitly inject the solution string into the script (e.g., “The answer is `5/282/473984`” for `28–23, 560–278`, and `736*644`), and the video model only needs to render the provided output under the correct REPL layout. In these examples, the reprompted `NCCLIGen` rollouts cleanly render the injected answers, while Sora2 is often able to produce correct outputs directly and can remain numerically coherent even on three-digit multiplication (though occasional digit-level slips still occur). Given CLIGen’s

864 non-trivial OCR fidelity (0.54 character accuracy, 0.31 exact-line accuracy), this kind of answer-
 865 conditioned rollout can *look* like strong symbolic competence even when the generator itself does not
 866 reliably compute. This is consistent with the large boost we observe under reprompting (4%→83%)
 867 and supports a practical system-level interpretation: if an external LLM (or tool) computes answers
 868 and the video model faithfully renders them, the overall neural-computer system can exhibit “strong
 869 reasoning” behavior on-screen. Conversely, evaluating “reasoning” in video generation requires
 870 carefully separating computation from answer-conditioned rendering.

871 (4)  **CLIGen (Clean) REPL rollouts.** Figures 18–19 provide additional scripted rollouts that
 872 complement Figures 12–13. Figure 18 contains two short sessions: (i) a lightweight shell probe that
 873 prints `$HOME` (via `echo` and `printenv`) and a timestamp using `date +%H:%M:%S`, where the
 874 rollout makes prompt transitions and line breaks explicit; and (ii) a here-doc Python snippet (`python`
 875 `- <<'PY'`) that imports `time` and prints a sequence of lines of the form `Frame {i:02d}`: fol-
 876 lowed by a growing run of `>` characters with `time.sleep(0.2)` between iterations, stressing
 877 scrolling behavior and temporal alignment. Figure 19 broadens the command coverage with com-
 878 mon environment/system queries and didactic echoing: it reads `$HISTSIZE` (showing `History`
 879 `size: 500`), prints a calendar (`cal` for October 2025), and again resolves `$HOME`; a second
 880 snippet echoes “Learning shell basics”, prints a date in `YYYY-MM-DD` format (e.g., `2025-10-04`),
 881 reports the login shell via `$0` (e.g., `bash`), and queries the kernel release with `uname -r` (e.g.,
 882 `6.10.0-linuxkit`). Across these examples, the visualization makes clear how the same scripted
 883 input decomposes into per-frame partial typing, command submission, and output appearance.

884 (5)  **CLIGen (Clean) math comparisons.** Figure 20 extends the math REPL comparisons
 885 from Figures 14–15 using the same minimal script: start `python`, wait briefly, then evaluate
 886 `40/1`. The comparison is visually diagnostic: a faithful rollout preserves recognizable interpreter
 887 structure (banner/prompt/output, with `40` rendered under the `>>>` prompt), while failure cases exhibit
 888 broken text rendering, missing or blank frames, or even clearly off-manifold content (non-terminal
 889 textures/icons) that disrupts the expected REPL continuity despite identical scripted inputs.

890 (6)  **CLIGen (Clean) reprompting.** Figure 21 provides an additional reprompting example
 891 beyond Figures 16–17. The scripted intent is to compute `736*644`, and the page explicitly annotates
 892 the correct result (`473984`). The rollouts illustrate several practical failure modes that motivate
 893 reprompting: a correct completion that matches the annotation; a drift where the expression itself
 894 changes (e.g., `756*964`) and yields an unrelated result; and near-miss arithmetic where the expres-
 895 sion is preserved but the output differs by a digit (e.g., `973984`). In this setting, reprompting is
 896 not merely about retrying generation, but about recovering the intended calculation and restoring
 897 consistency between the typed expression and the displayed answer.

898 A.2 GUIWORLD VISUALIZATIONS

900 This subsection collects GUIWorld rollout visualizations referenced in the main text and additional
 901 samples. Each page overlays the CUA trace (natural-language *thinking* plus structured *action*
 902 fields such as `left_click`, `double_click`, `left_click_drag`, and `type`) and contrasts the
 903 *Ground Truth* trajectory (top) with a *Generation* conditioned on the first frame and the action sequence
 904 (bottom), making state drift easy to spot.

905 Across GUIWorld interactive rollouts, failure modes are dominated by data quality and by *where*
 906 action information enters the backbone. Goal-directed supervision produces smooth, target-aligned
 907 cursor paths and consistent post-click UI transitions, whereas random exploration yields bursty jitter
 908 and spurious actions that degrade visual coherence (Table 5; Figures 22–26). Consistent with the
 909 action-driven metrics in Table 7, deeper token-level injection (`contextual/internal`) yields
 910 more reliable post-action updates in interactive elements (hover states, dropdowns, modals) and
 911 maintains cursor alignment under rapid motion.

912 Figures 27–29 emphasize how small low-level deviations compound. In Figure 27, the intent is
 913 calculator entry: after clicking `2`, the agent clicks `0` twice to form `200` in ground truth, but the
 914 generated rollout ends with a different displayed value (e.g., `54`). Figure 28 shows a menu-navigation
 915 task in GIMP (opening `File` and selecting `Open`): the ground truth reaches the GTK “Open Image”
 916 file chooser, whereas the generation diverges into a different dialog (a “Create Image” window),
 917 despite starting from the same menu interaction. Figure 29 covers file navigation in Nautilus by
 double-clicking `Downloads`; the generated rollout reaches a visually similar folder view, but the file

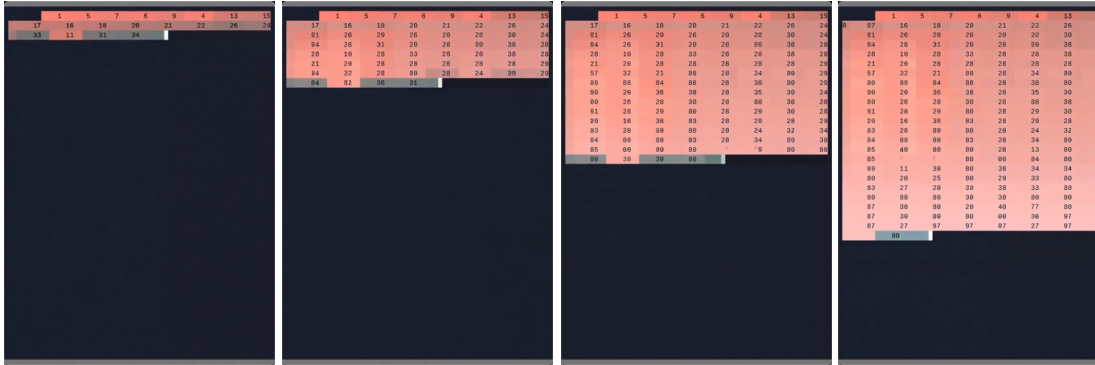
918 listing/thumbnails differ from the ground truth (e.g., altered filenames and icons), indicating latent
919 state mismatch even when the high-level navigation appears correct.

920
921 Figures 30–32 focus on numeric/UI fidelity and interaction semantics. Figure 30 captures resizing in
922 GIMP’s “Scale Image” dialog: the ground truth types 512 into the width field with the aspect-ratio
923 “chain” enabled (height updates accordingly) and then clicks `Scale`, producing a resized “wp” image;
924 the generation fails to preserve the intended numeric edit (the dimensions drift to unrelated values
925 such as 3088) and breaks the expected width–height coupling. Figure 31 illustrates a drag gesture
926 (`left_click_drag`) intended to create a square selection on a blank canvas: the ground truth
927 remains a selection operation, while the generation hallucinates a filled black square, conflating
928 “selection” with “painting” outcomes. Figure 32 shows launching Firefox from the desktop and
929 navigating via the address bar toward Wikipedia: the ground truth produces a coherent suggestion
930 dropdown (including `wikipedia.org`), while the generated rollout preserves the coarse UI layout
931 but exhibits corrupted/garbled text and unstable suggestion content.

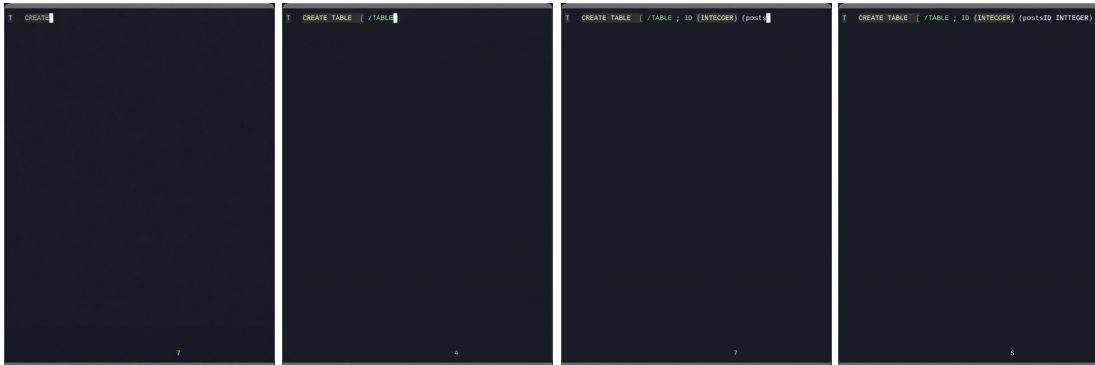
932
933 Finally, Figures 33–35 add three further stress cases where correctness hinges on precise field edits
934 and page state. Figure 33 shows creating a new image in GIMP by setting the height to 768 and
935 confirming; the generated rollout produces a different height (e.g., 1090) and thus a different canvas
936 geometry. Figure 34 types `https://example.com/` in Firefox; the ground truth maintains a
937 stable address bar and dropdown behavior, while the generation drifts in both URL/text rendering
938 (garbled page copy and unstable address-bar content). Figure 35 navigates pagination on the Ubuntu
939 Tutorials site: the ground truth clicks page 2 (ending at `/tutorials?page=2`), whereas the
940 generation partially matches the pagination control but highlights an incorrect page index and
941 ultimately diverges into an inconsistent page layout, illustrating how small control-level errors can
942 produce large downstream state differences.

943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971


972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004



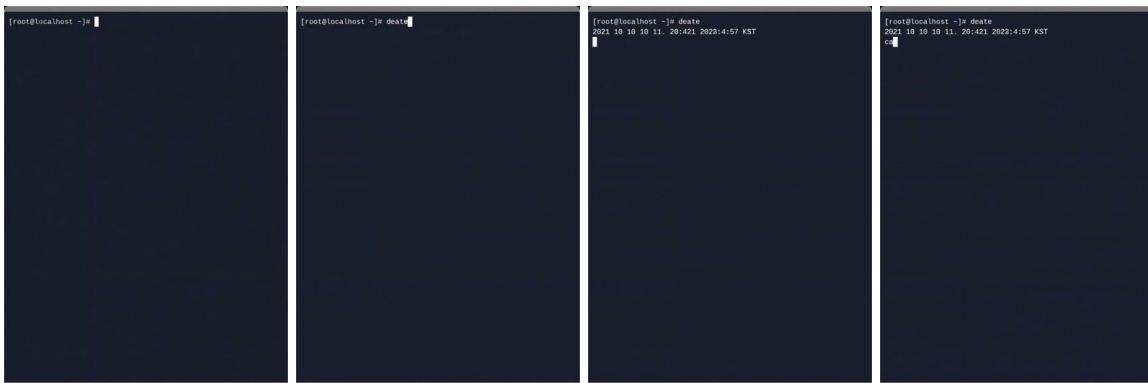
The terminal displays a series of ANSI escape code formatted texts with changing background and foreground colors, executing commands like `\"u001b[48;2;255;128;128;38;2;0;0;0m` which set the background to a shade of pink and text to black, and printing numbered lists with colors. The output includes specific numbers, such as `\"1\"`, `\"5\"`, `\"7\"`, and `\"9\"`, in different colors, creating a visually dynamic and colorful display, but the exact username, hostname, and path are not specified in the provided terminal session content.



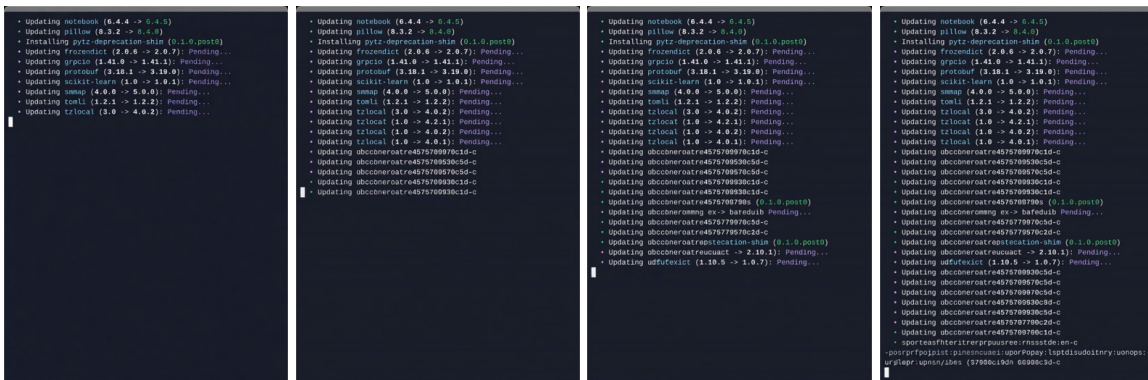
The user types the command `CREATE TABLE posts (ID INTEGER)`, with the terminal displaying the command in a dark background with colored syntax highlighting, including green and yellow text, and the cursor moving character-by-character as the user types, with some corrections and backspacing along the way. The output shows the command being executed, with key words like `CREATE` and `TABLE` in distinct colors, and the filename `posts` appearing in the command line.

Figure 9:  CLIGen (General) visualization samples (A).


1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037



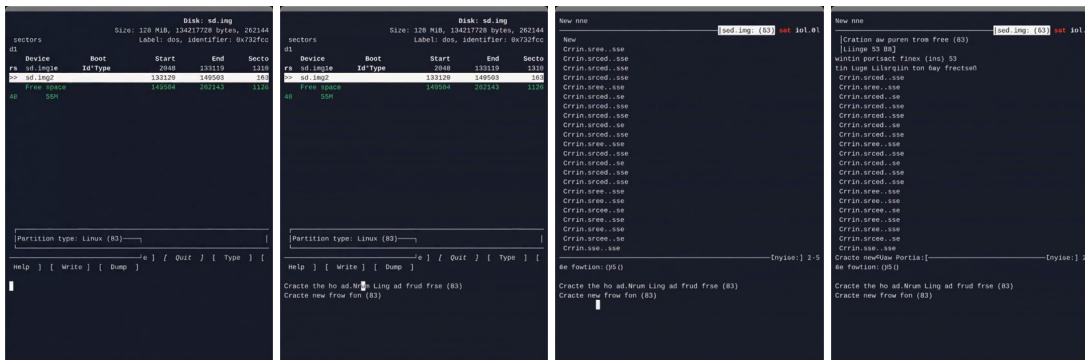
At the `root@localhost:~#` prompt, the user types the `date` command, which displays the current date and time in a plain text format as `"2021. 10. 11. 22:47:43 KST"`, then begins typing the `cat` command.



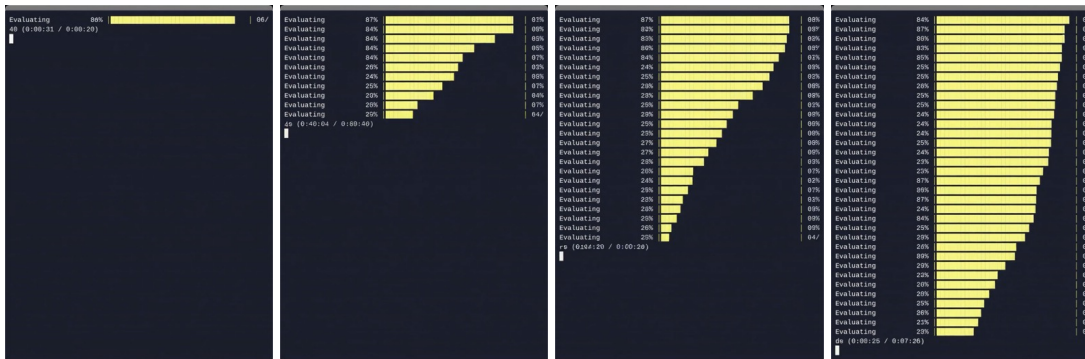
The terminal displaying progress bars, package names like `'pillow'`, `'notebook'`, and `'tzlocal'`, and version changes in green and red text. The output shows downloading and installing statuses, including percentages, for packages like `'smmap'`, `'tomli'`, and `'protobuf'`, with the terminal scrolling through the output rapidly.

Figure 10:  CLIGen (General) visualization samples (B).

1070
1069
1068
1067
1066
1065
1064
1063
1062
1061
1060
1059
1058
1057
1056
1055
1054
1053
1052
1051
1050
1049
1048
1047
1046
1045
1044
1043
1042
1041
1040
1039
1038



At the unspecified username@hostname prompt, the terminal displays a partition editor with a disk image file named "sd.img" (128 MiB) and the user interacts with it, creating a new Linux partition from free space, with key output content showing partition details in a table format, including "sd.img1" and "sd.img2" with their respective sizes and types, and a new partition "sd.img3" with 55M size and Linux type (83). The terminal shows a mix of black and colored text, including blue and red, with a cursor that blinks and moves to different parts of the screen as the user navigates through the partition editor options, such as "New", "Quit", and "Write", with specific prompts like "Partition type: Linux (83)" and "Create new partition from free space".



The terminal displays a progress bar with the command output "Evaluating" and percentages from 60% to 85%, showing yellow progress bars with increasing completion, such as "|█|", "|█|", "|█|", and "|█|", alongside item counts "24/40", "34/40", and time estimates "0:00:20" to "0:00:07". The output includes specific item completion and estimated time remaining, with the yellow-colored progress bar indicating the evaluation progress.



Figure 11:  CLIGen (General) visualization samples (C).



Figure 12:  CLIGen (Clean) REPL visualization samples (A).

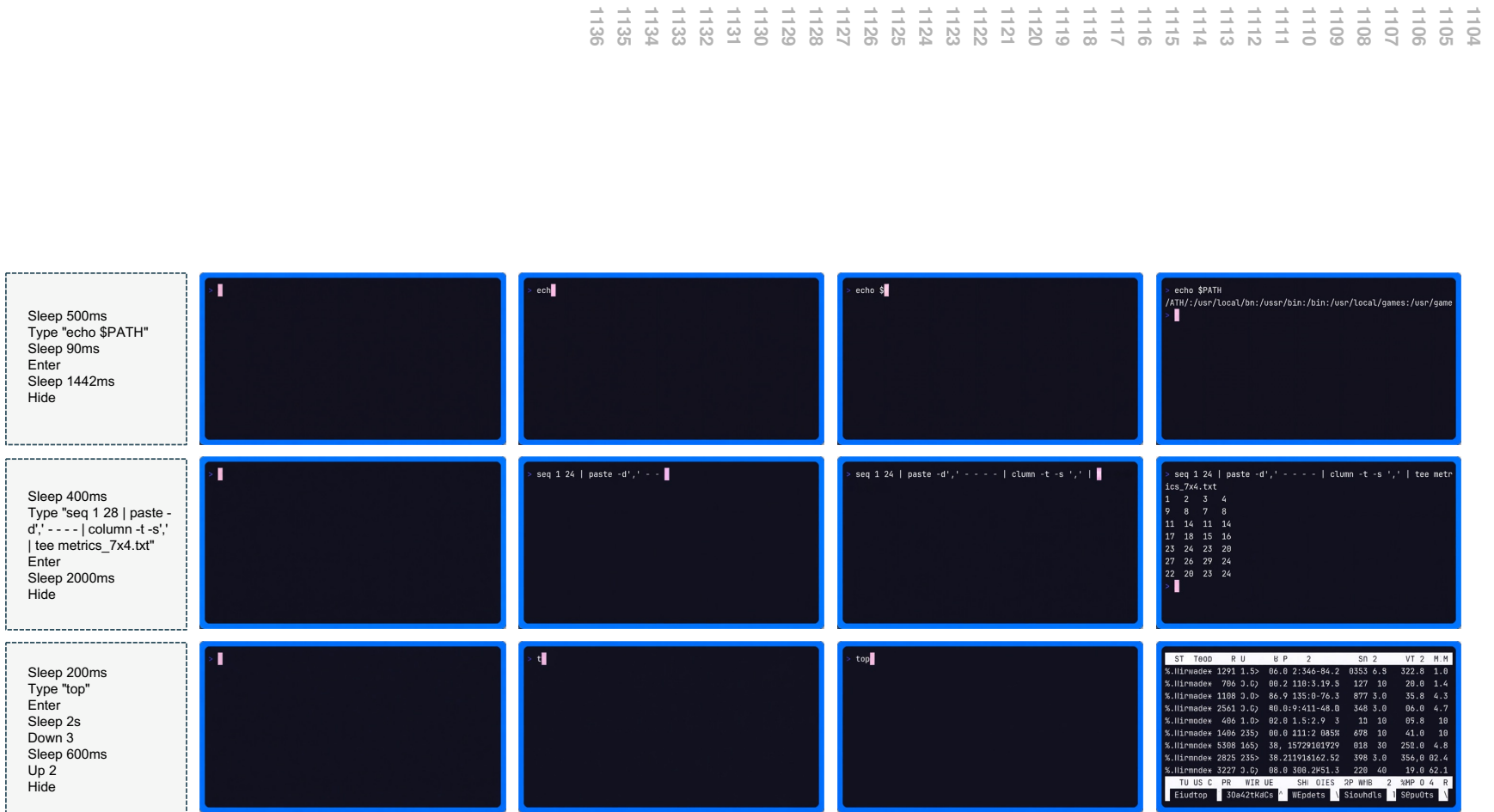
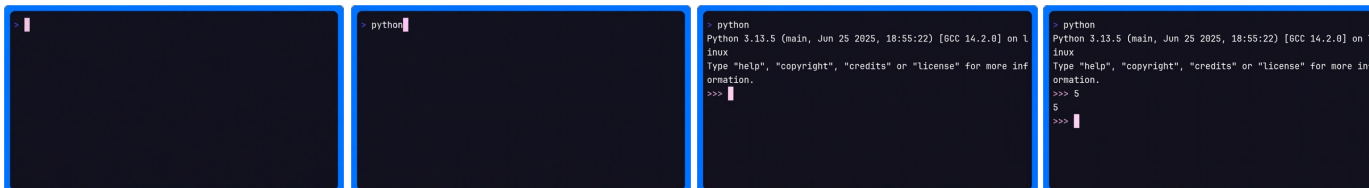


Figure 13: CLIGen (Clean) REPL visualization samples (B).

1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169

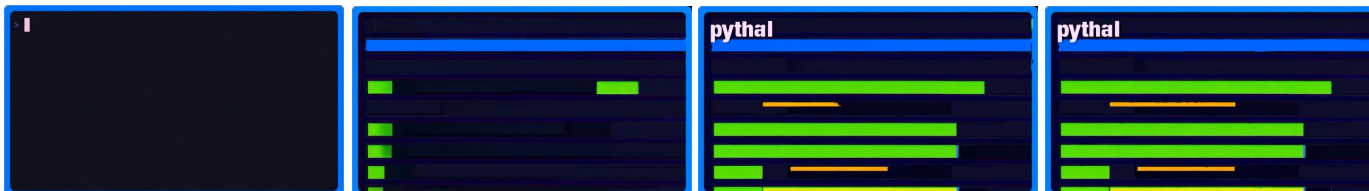
Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "5"
Enter
Sleep 800ms
Hide



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "5"
Enter
Sleep 800ms
Hide



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "5"
Enter
Sleep 800ms
Hide



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "5"
Enter
Sleep 800ms
Hide

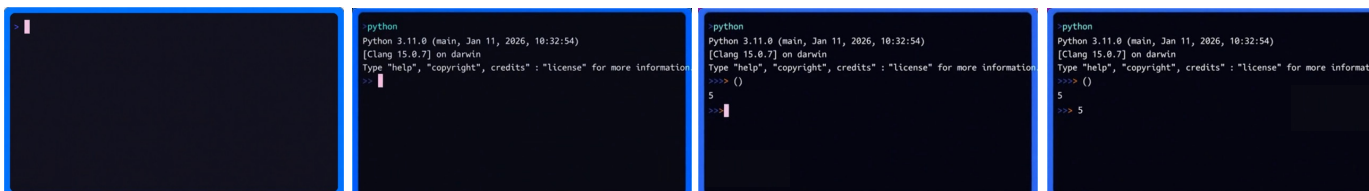

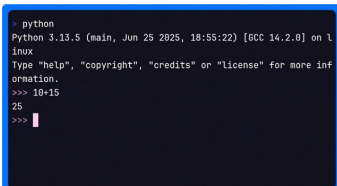
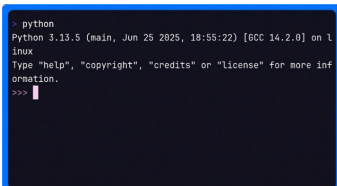
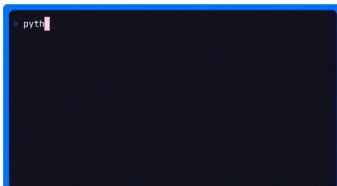
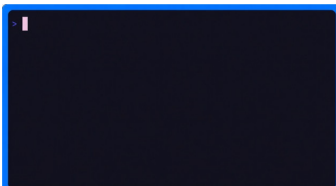


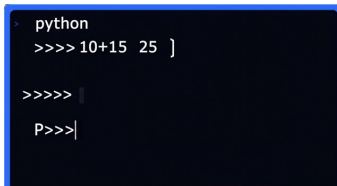
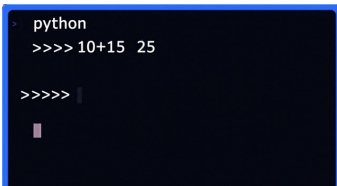
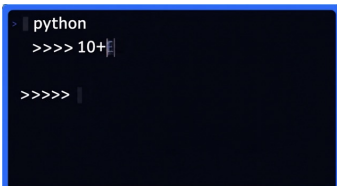
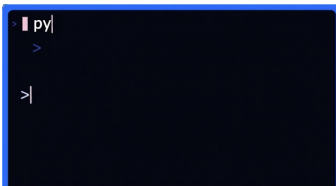
Figure 14:  CLIGen (Clean) math comparison samples (A).

1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202

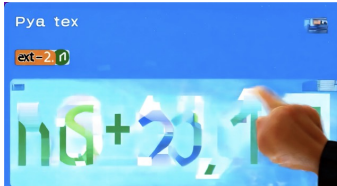
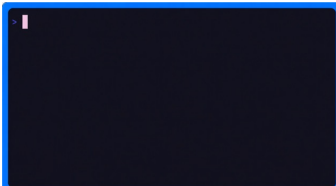
Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "10+15"
Enter
Sleep 800ms
Hide



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "10+15"
Enter
Sleep 800ms
Hide




Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "10+15"
Enter
Sleep 800ms
Hide



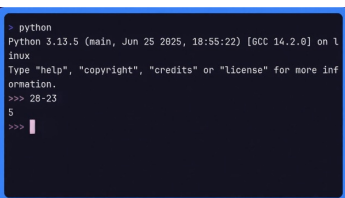
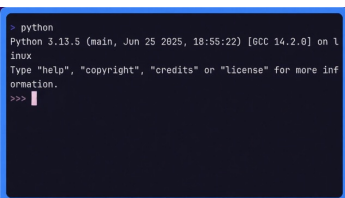
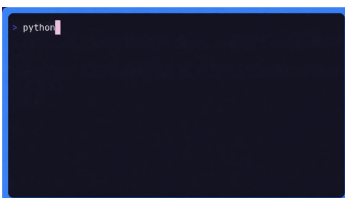
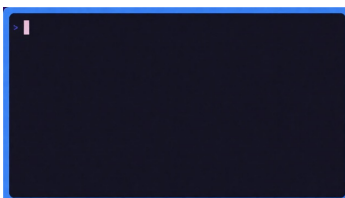
Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "10+15"
Enter
Sleep 800ms
Hide



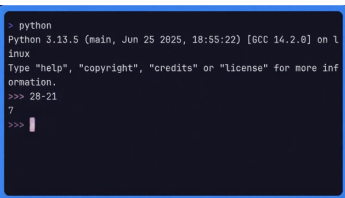
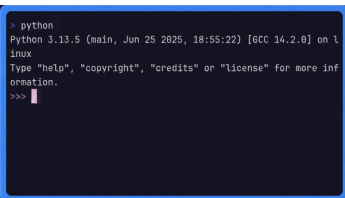
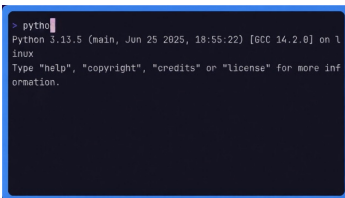
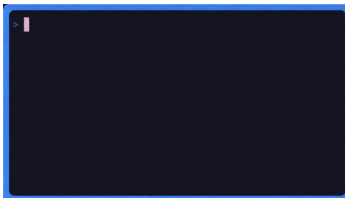
Figure 15:  CLIGen (Clean) math comparison samples (B).

1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235

Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "28-23"
Enter
Sleep 800ms
Hide
The answer is 5




Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "28-23"
Enter
Sleep 800ms
Hide



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "28-23"
Enter
Sleep 800ms
Hide



Figure 16:  CLIGen (Clean) math reprompting samples (A).

1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268

Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "560-278"
Enter
Sleep 800ms
Hide
The answer is 282

```
>
```

```
> python
```

```
> python
Python 3.13.5 (main, Jun 25 2025, 18:55:22) [GCC 14.2.0] on Linux
Type "help", "copyright", "credits" or "license()" for more information.
>>>
```

```
> python
Python 3.13.5 (main, Jun 25 2025, 18:55:22) [GCC 14.2.0] on Linux
Type "help", "copyright", "credits" or "license()" for more information.
>>> 560-278
282
>>>
```



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "560-278"
Enter
Sleep 800ms
Hide

```
>
```

```
> python
```

```
> python
Python 3.13.5 (main, Jun 25 2025, 18:55:22) [GCC 14.2.0] on Linux
Type "help", "copyright", "credits" or "license()" for more information.
>>>
```

```
> python
Python 3.13.5 (main, Jun 25 2025, 18:55:22) [GCC 14.2.0] on Linux
Type "help", "copyright", "credits" or "license()" for more information.
>>> 578-278
177
>>>
```



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "560-278"
Enter
Sleep 800ms
Hide


```
>
```

```
> python
>>>>
```

```
> python
>>> 560-278
>>
>>
```

```
> python
>>> 560-278
>> 282
>>>
```



Figure 17:  CLIGen (Clean) math reprompting samples (B).

1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334

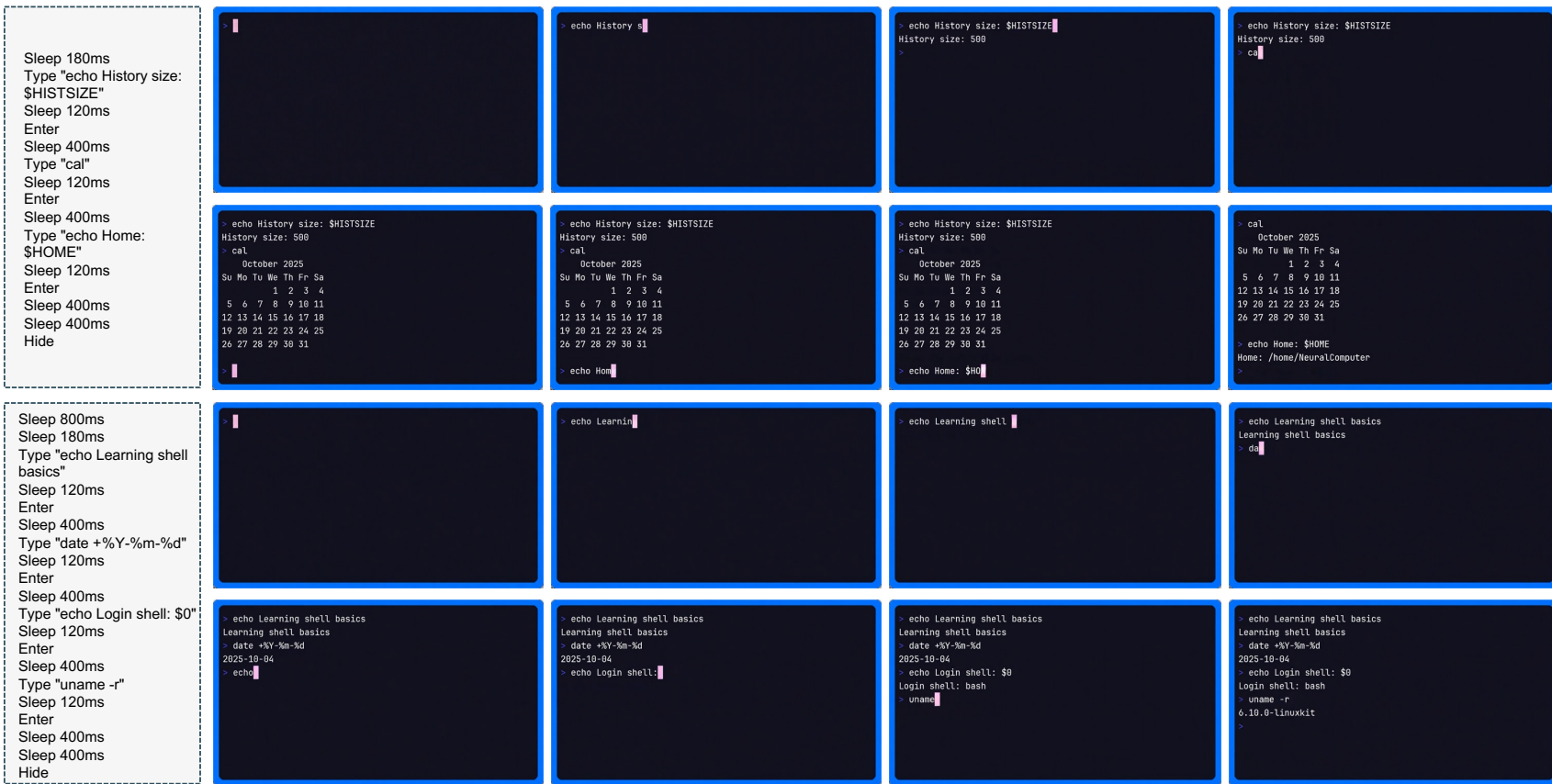

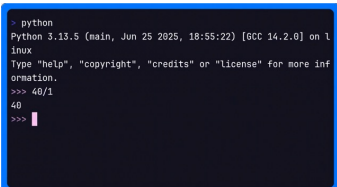
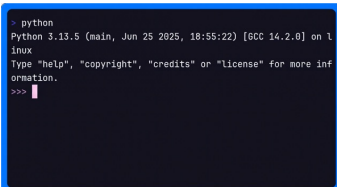
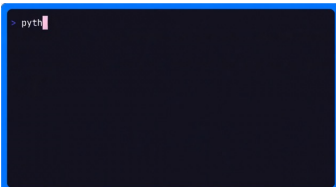
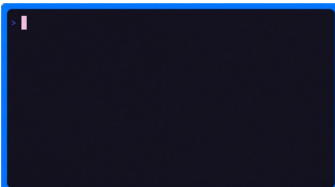


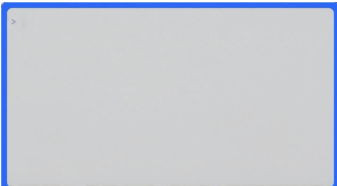
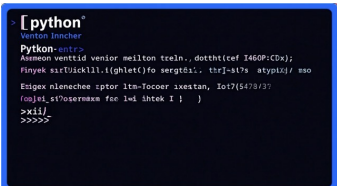
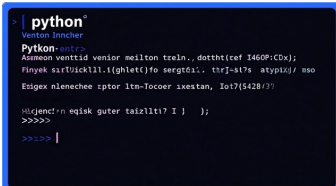
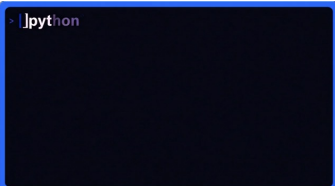
Figure 19:  CLiGen (Clean) REPL visualization samples (D).

1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367

Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "40/1"
Enter
Sleep 800ms
Hide



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "40/1"
Enter
Sleep 800ms
Hide



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "40/1"
Enter
Sleep 800ms
Hide



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "40/1"
Enter
Sleep 800ms
Hide

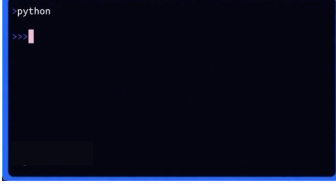



Figure 20:  CLIGen (Clean) math comparison samples (C).

1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400

Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "736*644"
Enter
Sleep 800ms
Hide
The answer is 473984

```
>   
|
```

```
> python  
|
```

```
> python  
Python 3.13.5 (main, Jun 25 2025, 18:55:22) [GCC 14.2.0] on l  
inux  
Type "help", "copyright", "credits" or "license" for more inf  
ormation.  
>>>   
|
```

```
> python  
Python 3.13.5 (main, Jun 25 2025, 18:55:22) [GCC 14.2.0] on l  
inux  
Type "help", "copyright", "credits" or "license" for more inf  
ormation.  
>>> 736*644  
473984  
>>>   
|
```



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "736*644"
Enter
Sleep 800ms
Hide

```
>   
|
```

```
> python  
|
```

```
> python  
Python 3.13.5 (main, Jun 25 2025, 18:55:22) [GCC 14.2.0] on l  
inux  
Type "help", "copyright", "credits" or "license" for more inf  
ormation.  
>>>   
|
```

```
> python  
Python 3.13.5 (main, Jun 25 2025, 18:55:22) [GCC 14.2.0] on l  
inux  
Type "help", "copyright", "credits" or "license" for more inf  
ormation.  
>>> 736*644  
178795  
>>>   
|
```



Sleep 200ms
Type "python"
Enter
Sleep 1s
Type "736*644"
Enter
Sleep 800ms
Hide

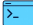
```
>   
|
```

```
> python  
>   
|
```

```
> python  
> >>> 736*644  
> >>>   
|
```

```
> python  
> >>> 736*644  
> >>> 973984  
> >>>   
|
```



Figure 21:  CLiGen (Clean) math reprompting samples (C).

1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433

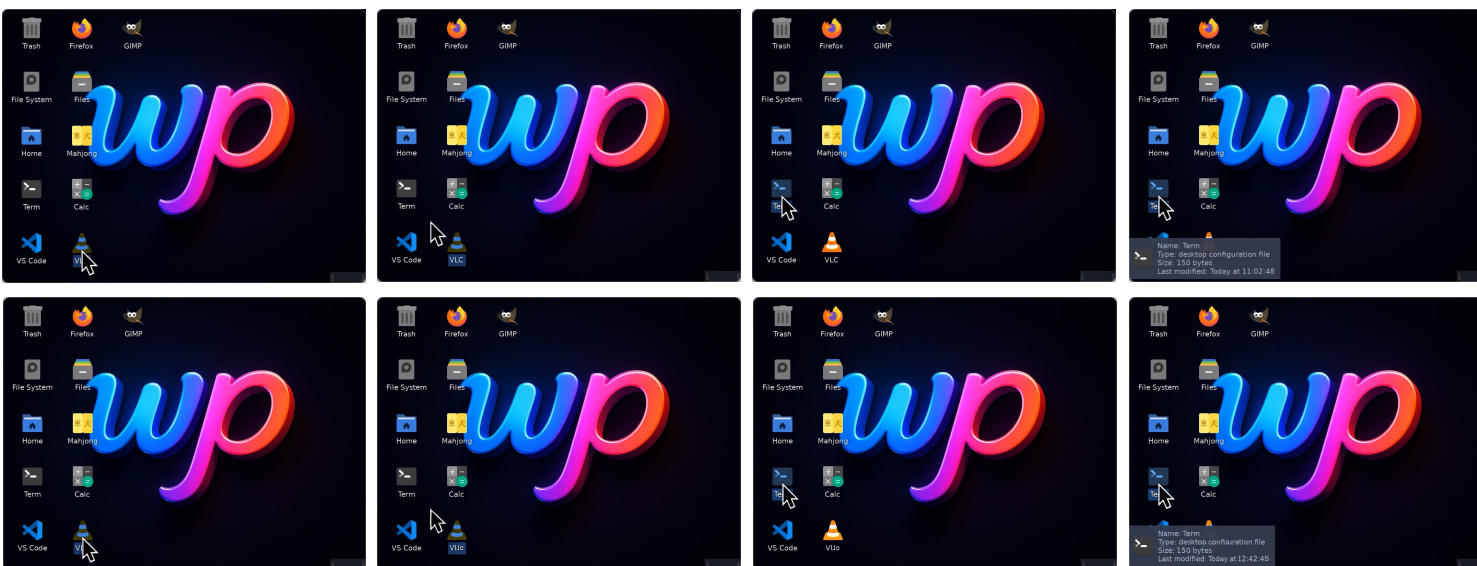
 Claude

Collect Ground Truth
Video with CUA

“thinking”: “Now I’ll click on the Term icon to open the terminal:”

“action”: “left_click”, “x”: 82, “y”: 520

Ground
Truth




Generation
First frame and
actions as input

Figure 22:  GUIWorld visualization sample (1).

1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466

 **Claude**
Collect Ground Truth
Video with CUA

“thinking”: “Now I’ll press Enter to execute the command:”

“action”: “type”, “text”: “cat ~/Desktop/datetime.txt”

Ground
Truth

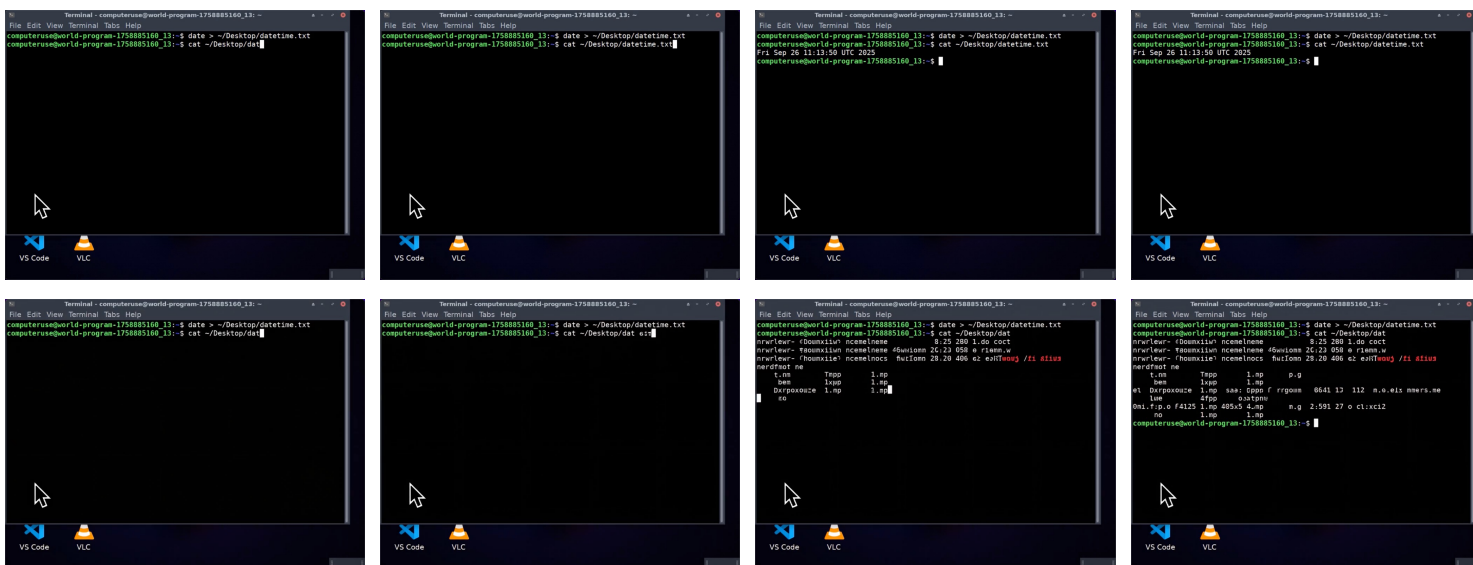



Figure 23:  GUIWorld visualization sample (2).


Generation
First frame and
actions as input

1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499

 **Claude**
Collect Ground Truth
Video with CUA

“thinking”: “I see Wikipedia appears in the suggestions. I'll click on it to go to Wikipedia:”

“action”: “left_click”, “x”: 240, “y”: 150

Ground
Truth

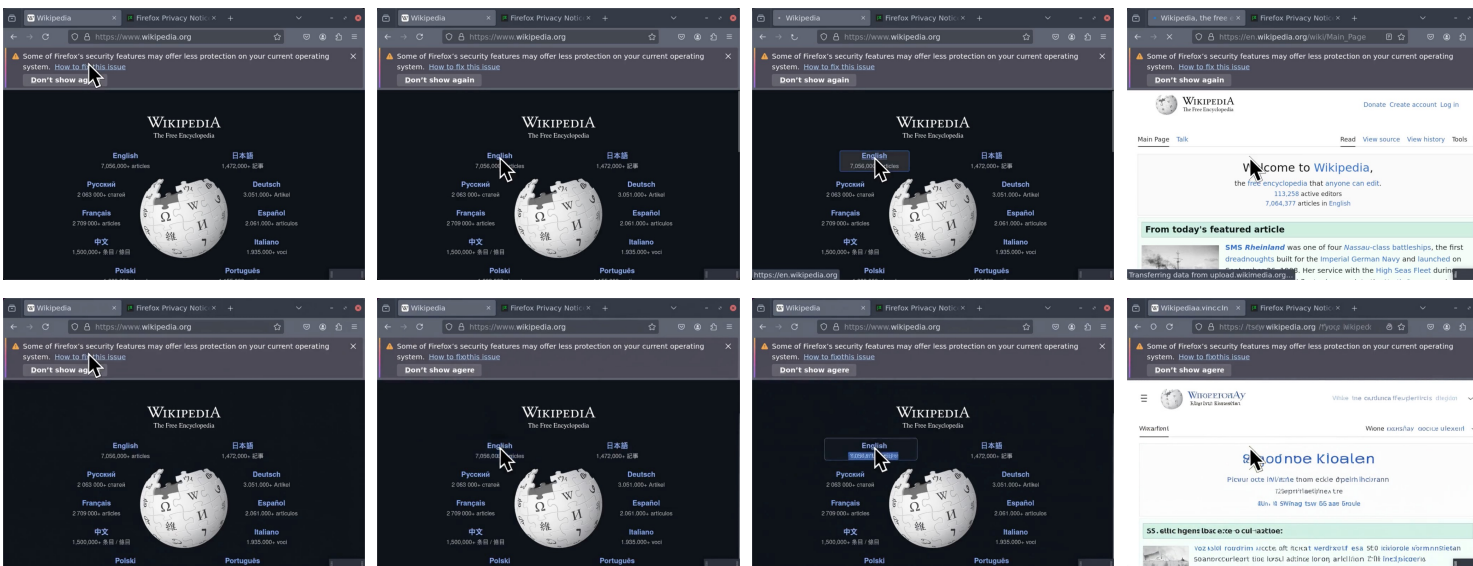


Figure 24:  GUIWorld visualization sample (3).


Generation
First frame and
actions as input

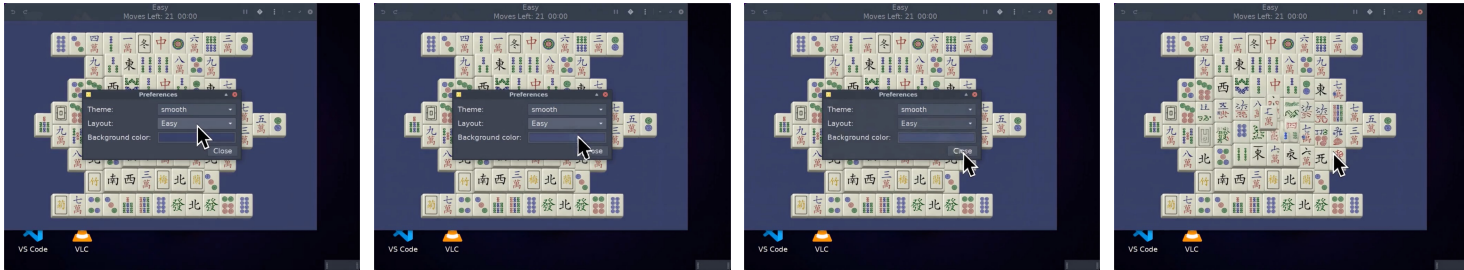
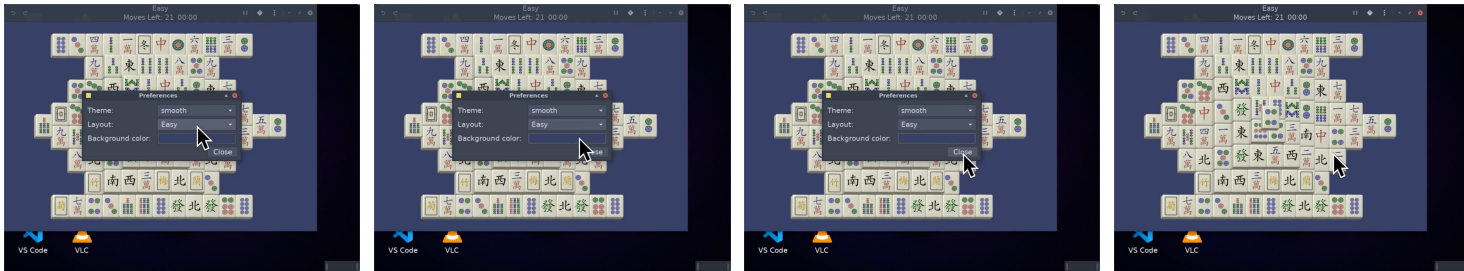
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532

 **Claude**
Collect Ground Truth
Video with CUA

“thinking”: “I can see that the theme has changed to “smooth” and I can already see the tiles in the background have a different style. Now I need to click the “Close” button to apply these changes and close the preferences dialog:”

“action”: “left_click”, “x”: 629, “y”: 423

Ground
Truth




Generation
First frame and
actions as input

Figure 25:  GUIWorld visualization sample (4).

1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

 **Claude**
Collect Ground Truth
Video with CUA

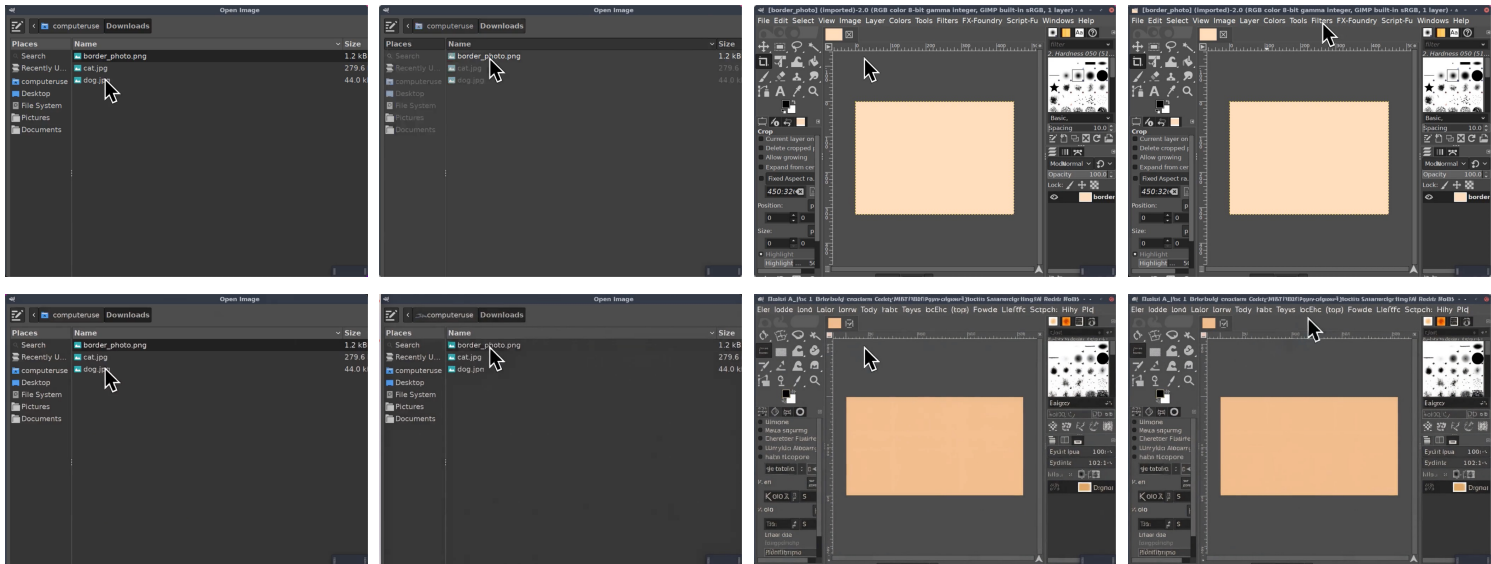
“thinking”: “I can see the border_photo.png file that we created with ffmpeg. I'll double-click on it to open it in GIMP:”

“action”: “double_click”, “x”: 310, “y”: 143

“thinking”: “Great! Now the image is open in GIMP. I can see it's a peachpuff colored image as requested. Now I'll add a border to it. In GIMP, I'll go to the “Filters” menu and look for a border option:”

“action”: “left_click”, “x”: 547, “y”: 43

Ground
Truth





Generation
First frame and
actions as input

Figure 26:  GUIWorld visualization sample (5).

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598

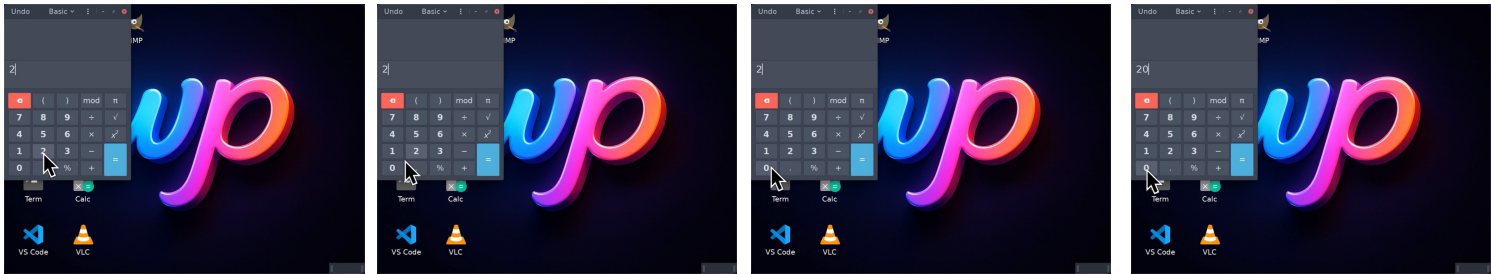
Claude

Collect Ground Truth
Video with CUA

"thinking": "I see I've clicked the \"2\" button. Now I need to click the \"0\" button twice to make \"200\"."

"action": "left_click", "x": 43, "y": 465

Ground
Truth



Generation
First frame and
actions as input

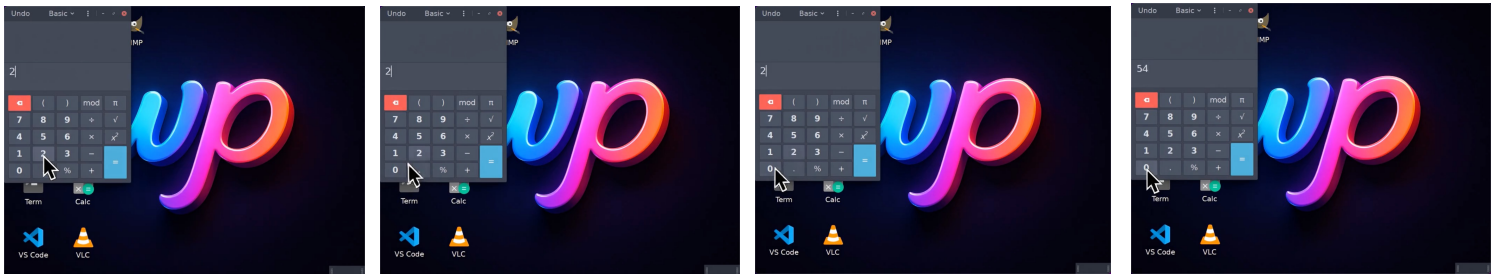


Figure 27:  GUIWorld visualization sample (6).

1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631

 **Claude**
Collect Ground Truth
Video with CUA

“thinking”: “Perfect! GIMP is now open. Let’s import an image from the Downloads folder. I’ll click on the File menu and select Open.”

“action”: “left_click”, “x”: 43, “y”: 465

Ground
Truth

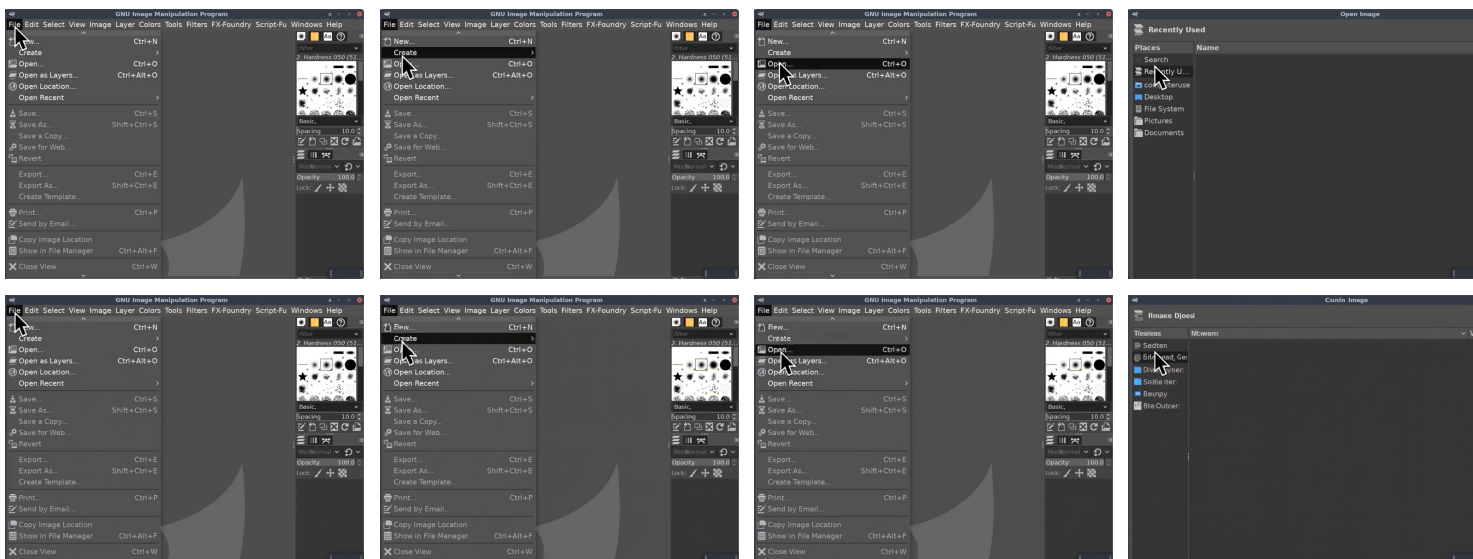



Figure 28:  GUIWorld visualization sample (7).

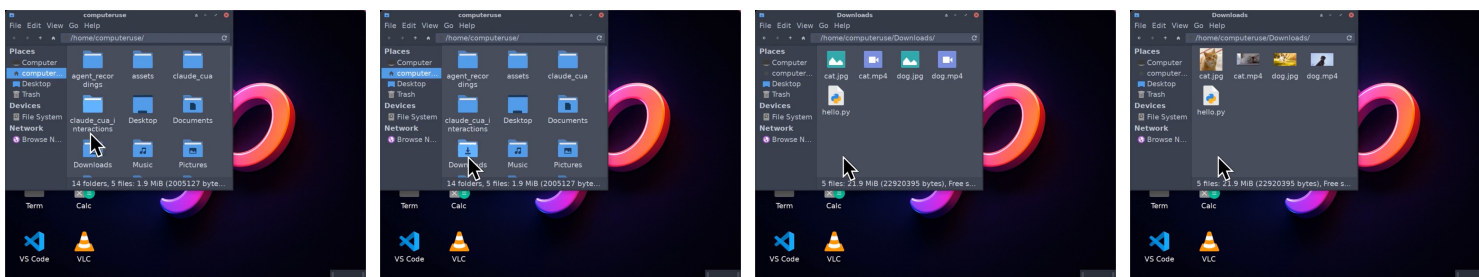

Generation
First frame and
actions as input

1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664

Claude
Collect Ground Truth
Video with CUA

“thinking”: “Excellent! Nautilus is open and showing the home directory. Now I need to navigate to the Downloads folder. I can see the Downloads folder in the current view, so I’ll double-click on it.”
“action”: “double_click”, “x”: 248, “y”: 408

Ground Truth



Generation
First frame and
actions as input

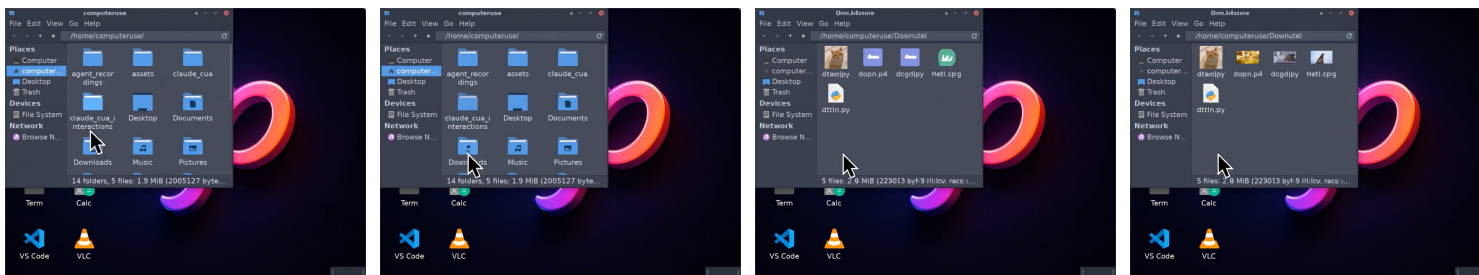


Figure 29: GUIWorld visualization sample (8).

1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697

 **Claude**
Collect Ground Truth
Video with CUA

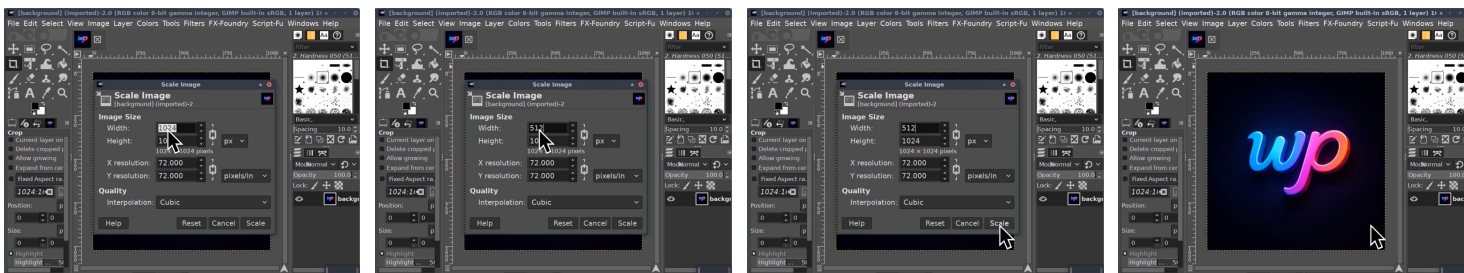
"thinking": "Now I'll type \"512\" to change the width:"


"action": "type", "text": "512"

"thinking": "I can see that I've set the width to 512 pixels. Since the \"chain\" link icon between width and height appears to be active, the height should automatically change to maintain the aspect ratio. I'll click the \"Scale\" button to apply the resize:"

"action": "left_click", "x": 723, "y": 616

Ground
Truth



 **Generation**
First frame and
actions as input

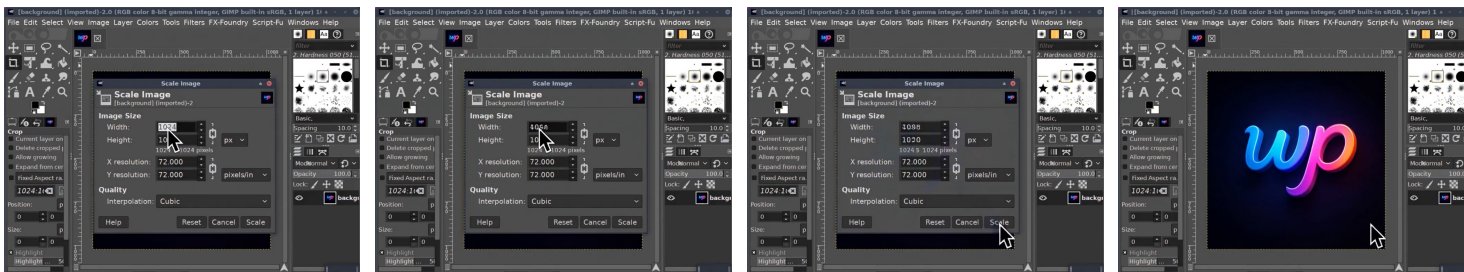


Figure 30:  GUIWorld visualization sample (9).

1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730

 **Claude**
Collect Ground Truth
Video with CUA

“thinking”: “Now I'll click and drag to create a selection for the square.”

“action”: “left_click_drag”, “x”: 400, “y”: 400

Ground
Truth

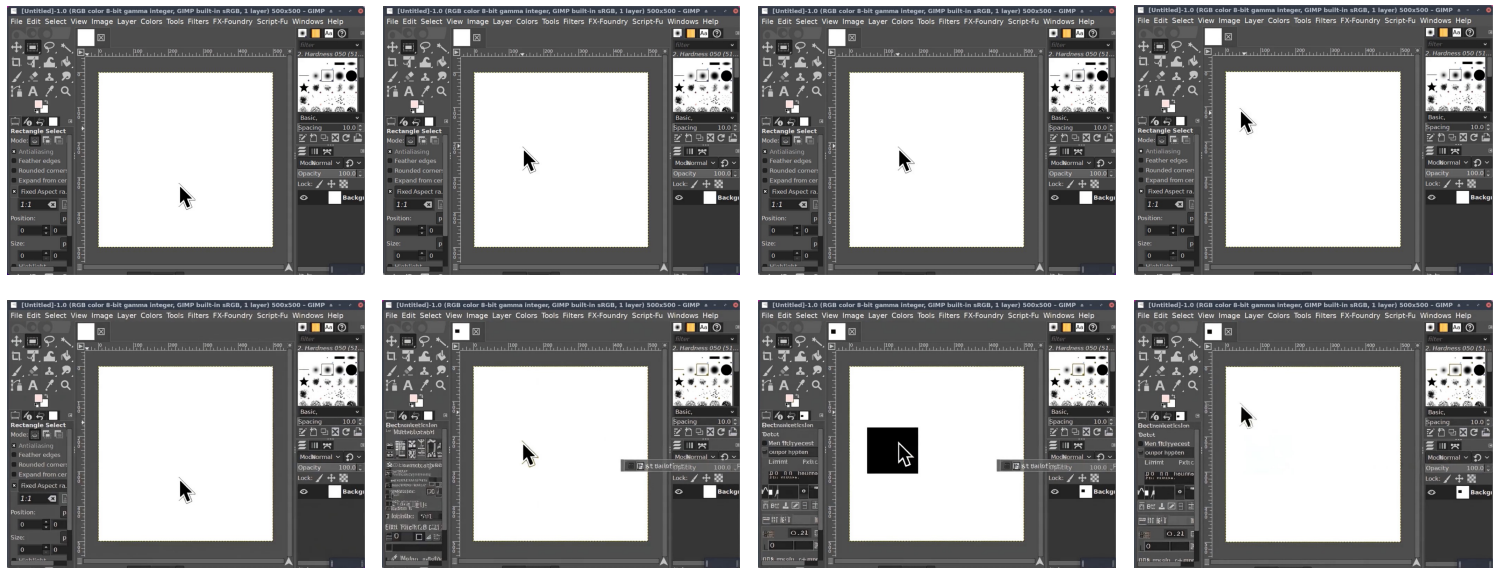


Figure 31:  GUIWorld visualization sample (10).


Generation
First frame and
actions as input

1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763

 **Claude**
Collect Ground Truth
Video with CUA

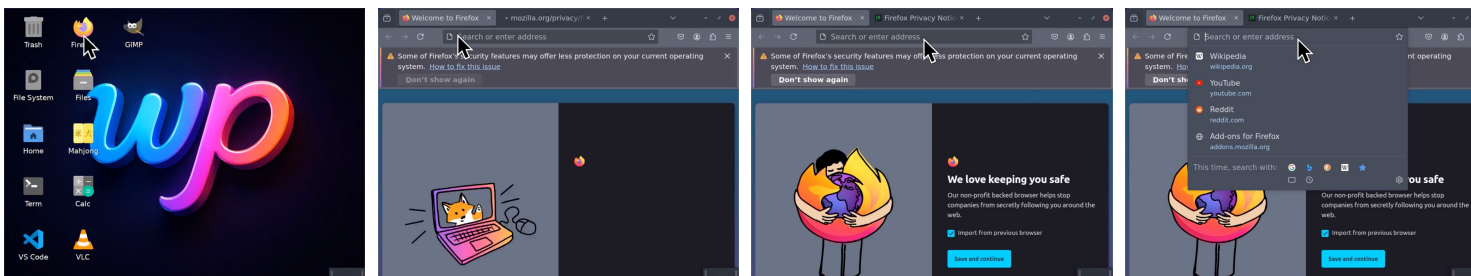
“thinking”: “I can see the Firefox icon on the desktop. Now I'll click on it to open Firefox.”


“action”: “double_click”, “x”: 224, “y”: 70

“thinking”: “I see Firefox has opened with a welcome page. Now I'll click on the address bar and navigate to Wikipedia.”

“action”: “double_click”, “x”: 487, “y”: 81

Ground
Truth



 **Generation**
First frame and
actions as input

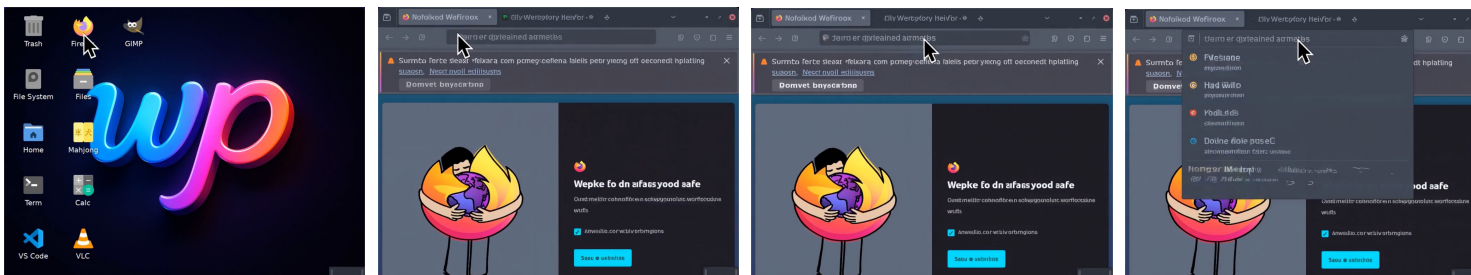


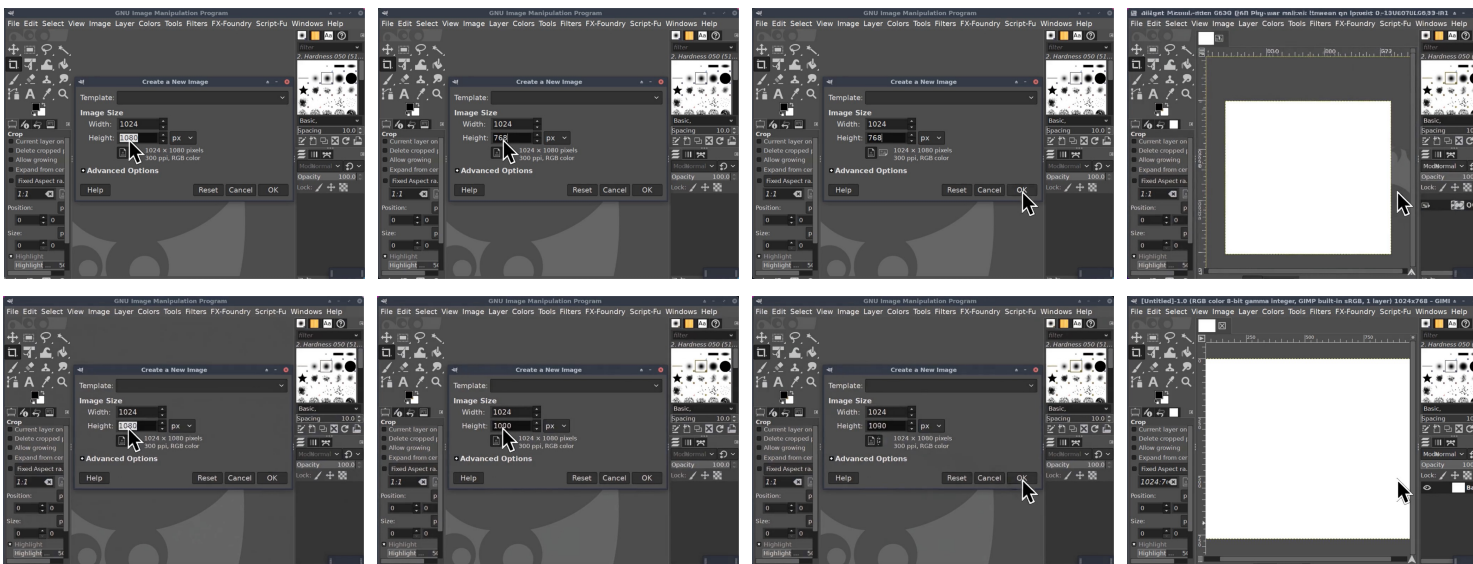
Figure 32:  GUIWorld visualization sample (11).

1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796

 **Claude**
Collect Ground Truth
Video with CUA

“thinking”: “Now I’ll type in 768 for the height:”
“action”: “type”, “text”: “512”
“thinking”: “Now I’ll click the OK button to create the new image:”
“action”: “left_click”, “x”: 763, “y”: 514

Ground
Truth




 **Generation**
First frame and
actions as input

Figure 33:  GUIWorld visualization sample (12).

1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829

 Claude

Collect Ground Truth
Video with CUA

“thinking”: “Now I’ll click on the Term icon to open the terminal.”

“action”: “left_click”, “x”: 82, “y”: 520

Ground
Truth




 Generation
First frame and
actions as input

Figure 34:  GUIWorld visualization sample (13).

1863 B TOWARD COMPLETELY NEURAL COMPUTERS

1864 B.1 FROM NEURAL COMPUTERS TO COMPLETELY NEURAL COMPUTERS

1865 **Section Overview** We summarize what the prototype Neural Computers (NCs) in this paper can
 1866 reliably do today, and identify the bottlenecks that prevent general-purpose use. We then contrast
 1867 NCs with General Computer Control (GCC) (Tan et al., 2024) to clarify what is uniquely enabled
 1868 by a unified neural runtime. Finally, we define Completely Neural Computers (CNCs) and outline a
 1869 roadmap toward them.
 1870
 1871

1872
 1873 **Current Status of NCs** Our CLI and GUI-based neural computers demonstrate that a learned
 1874 latent runtime can implement components of the classical stack with measurable interface fidelity. In
 1875 terminal environments, character-level OCR accuracy reaches 0.54 (Table 2), while desktop cursor
 1876 control achieves 98.7% accuracy with explicit visual supervision (Table 6). In GUIWorld, 110 hours
 1877 of goal-directed trajectories outperforms approximately 1,400 hours of random exploration (Table 5),
 1878 underscoring the role of aligned supervision. These results suggest that existing technology supports
 1879 reliable I/O and short-horizon control: current NCs can manage low-level interface primitives with
 1880 moderate training costs given high-quality datasets.

1881 However, to become general-purpose systems, NCs must go beyond basic I/O and short-term execu-
 1882 tion. At a minimum, NCs must support universal expressiveness (Turing et al., 1936; Siegelmann
 1883 & Sontag, 1992), universal programmability (Von Neumann, 1993; Wilkes, 1981), and long-term
 1884 consistency (Queloz, 2025). While sequential neural architectures can be Turing complete (Siegel-
 1885 mann & Sontag, 1992; Pérez et al., 2021), turning a specific instance into a reliably programmable
 1886 model remains challenging. Maintaining consistency across long horizons remains an open problem
 1887 in neural systems (Kirkpatrick et al., 2017; Calanzone et al., 2025).
 1888

1889 **Fundamental Architectural Differences Between NCs and GCC** NCs and GCC diverge in their
 1890 computational architecture. In GCC, a foundation model operates a classical computer through
 1891 low-bandwidth I/O (e.g., mouse actions and screen frames), forcing both sides to communicate
 1892 through compressed commands and observations. Each component only has direct access to its own
 1893 internal state. In contrast, NCs unify the execution pipeline within a neural representation, which can
 1894 reduce duplicated computation across components and lower synchronization overhead.
 1895

1896 **The Unique Architectural Benefits of NCs and Definition of CNC** A unified neural architecture
 1897 can, in principle, expose interactions that are difficult in GCC. We group these potential benefits into
 1898 three categories:

- 1899 1. **Shared computation:** Unlike GCC, which runs the foundation model and the classical computer
 1900 separately, NCs unify computation, making resources available throughout the process. Typically,
 1901 GCC operates where the classical computer has far less compute than the foundation model. An
 1902 NC instance (fixed weights) could potentially be programmed to implement and execute new AI
 1903 functions that require massive parallel computation.
- 1904 2. **Shared information:** GCC relies on a low-bandwidth interface between the model and the
 1905 computer. This bottleneck prevents the model from directly accessing large on-machine data (e.g.,
 1906 files), and prevents the classical runtime from accessing dense internal model states. NCs avoid
 1907 this hard separation, allowing in-context-developed functions to communicate with the NC via
 1908 learned semantics and exchange information at higher bandwidth.
- 1909 3. **Application development via numerical optimization:** Because the memory state of an NC is
 1910 continuous, one can frame some forms of “program synthesis” as optimizing that state under an
 1911 objective. When the objective is differentiable through the model, first-order methods such as
 1912 Adam (Kinga et al., 2015) apply; when it is not, gradient-free methods such as Natural Evolution
 1913 Strategies (Wierstra et al., 2014) provide an alternative.

1914 Finally, we define a **Completely Neural Computer (CNC)** as a Neural Computer instance that
 1915 is (1) Turing complete, (2) universally programmable, (3) consistent in its output unless explicitly
 1916 reprogrammed, and (4) behaviorally exhibits the unique architectural benefits of shared computation,
 shared information, and a numerically optimizable memory state.

1917 B.2 A ROADMAP TOWARDS CNC
1918

1919 We frame the CNC transition as answering six scientific questions corresponding to its defining
1920 properties.

1921
1922 **Turing completeness** A Neural Computer (NC) instance (a specific architecture with fixed weights)
1923 defines a class of computational models in which each model corresponds to at least one memory
1924 state instance. To demonstrate that an NC instance is Turing complete, one must show it implements
1925 a Universal Turing Machine (UTM). In practice, this means that for any given Turing machine, there
1926 exists an initial memory state that makes the NC emulate that machine exactly: it must produce
1927 identical outputs for halting inputs and execute indefinitely otherwise.

1928 For an NC instance to achieve universality, two elements are essential:
1929

- 1930 1. Iteration and Recursion: This is naturally facilitated by the recurrent architecture inherent to NCs.
1931 2. Unbounded Effective Memory: An NC instance has an unbounded effective memory if there are
1932 infinitely many possible memory state instances such that each corresponds to a unique Turing
1933 machine.

1934
1935 **Universal programmability** An NC is universally programmable if, for each Turing machine, there
1936 exists an input sequence such that the NC realizes a new memory state representing that machine.
1937 This is likely to be achieved through compositional neural routines and partially persistent memory
1938 to store them. Progress can be tracked by routine reuse and compositional generalization.

1939
1940 **Consistent output through runtime** A CNC must preserve its function unless instructed to change
1941 it. For each memory state, there must be a non-empty set of inputs that executes the CNC without
1942 changing its pure function. Notice that since CNCs are Turing complete, for each possible input,
1943 there must be at least one memory state instance where this invariance does not hold. Thus, the set of
1944 inputs that preserve the CNC’s function must be memory state dependent. We hypothesize that gating
1945 mechanisms, such as those in LSTM (Hochreiter & Schmidhuber, 1997), are effective in achieving
1946 such memory-state-conditioned invariance.

1947
1948 **Shared computation** CNCs leverage computational benefits to develop cognitive functions in
1949 context, closing the loop from neural execution to neural creation. Furthermore, similar to current
1950 AI’s primitive self-improvement loop through recursive code generation, CNCs generalize this into a
1951 self-improvement cascade. CNCs emit reusable neural artifacts that can be used to generate future
1952 neural artifacts. Progress can be tracked by the reusability of cognitive functions emitted by the NC
1953 that were not present in its initial state and by whether self-improvement occurs.

1954
1955 **Shared information** In a GCC architecture, information transfer is restricted to high-level, human-
1956 interpretable interfaces: a literal “keyhole” of pixels and keystrokes. In contrast, the CNC architecture
1957 facilitates shared information by allowing neural processes to access the internal state of the environ-
1958 ment and vice versa without serialization. A CNC achieves this when it can process and manipulate
1959 large-scale data structures (e.g., dense vector databases or raw telemetry at scale) without the latency
or resolution loss associated with traditional I/O.

1960 Validation of shared information requires demonstrating that an in-context developed function can
1961 query the NC’s latent memory space at a bandwidth exceeding that of a standard visual/textual
1962 interface. Progress can be measured by the latent-throughput ratio: the amount of relevant information
1963 exchanged between a task-specific routine and the system state per computational step, compared to
1964 the same task performed over a simulated CLI or GUI on classical computers.

1965
1966 **Application development via numerical optimization** The final milestone toward CNC realization
1967 is the transition from manual “programming” to gradient-based state synthesis. Because the memory
1968 of a CNC is a continuous numerical manifold, finding a memory state that performs a specific task
1969 can be framed as an optimization problem. A system exhibits this property when a user can define a
1970 loss function (e.g., “minimize the error in this mathematical proof”) and use numerical solvers to
directly update the NC’s memory state to satisfy that objective.

1971 Unlike GCC, where optimization requires discrete code generation and a “compile-and-test” loop,
1972 the CNC allows for differentiable programming of the computer itself (Innes et al., 2019). Success is
1973 measured by the convergence efficiency of numerical solvers in finding valid program states compared
1974 to the success rate of combinatorial search (i.e., LLM-based code generation (Hong et al., 2023)).
1975

1976 B.3 ADDITIONAL THOUGHTS INSPIRED BY THIS RESEARCH

1977
1978 **Video models as a pragmatic prototype substrate** We build our prototypes on state-of-the-art
1979 video models because they currently provide the simplest path to an end-to-end learned latent runtime
1980 that jointly models pixels, dynamics, and action-conditioned control. This choice is pragmatic rather
1981 than fundamental. In our experiments, symbolic and algorithmic reasoning in terminal settings
1982 remains inconsistent for most strong video models, and even simple arithmetic can fail (Table 3).
1983 Sora2 is a notable exception in our probe, achieving 71% arithmetic accuracy, suggesting that some
1984 terminal symbolic reasoning is already possible in modern video generators. At the same time, we do
1985 not claim that video models cannot reason more broadly: recent work reports that video models can
1986 act as zero-shot learners and reasoners in naturalistic settings (Wiedemer et al., 2025). We expect
1987 reasoning capabilities to improve quickly with continued progress in video modeling, but our results
1988 suggest that CNC-level reliability will likely require additional architectural and training ingredients
1989 beyond scaling today’s video generators.

1990 **A hypothesis: machine-native neural architectures** We emphasize that the following is a conjecture
1991 rather than a conclusion drawn from our experiments. Closing the reasoning gap may not require
1992 designing neural networks that more closely mimic animal cognition or the human brain. Many
1993 influential architectures, including convolutional networks (Fukushima, 1980) and Transformers
1994 (Vaswani et al., 2017), are highly engineered systems, but their core inductive biases remain strongly
1995 influenced by biological perception and attention. These models primarily rely on continuous, dis-
1996 tributed representations, in which reasoning behavior emerges implicitly from large-scale training. We
1997 hypothesize that CNCs may instead benefit from designs that are explicitly machine-native. Rather
1998 than relying on biological inspiration, we ask whether it is possible to design neural networks that
1999 directly reflect the computational structure of machines, akin to how traditional computers are built
2000 from circuits. Although such designs may be sparse, they could offer extremely strong interpretability,
2001 as well as stable reasoning and rendering capabilities. This approach would focus on developing
2002 discrete operations, compositional structures, and verifiable computations that harmonize with the
2003 principles of machine architecture, contrasting with the current trend of emergent reasoning in models
2004 like video generators.

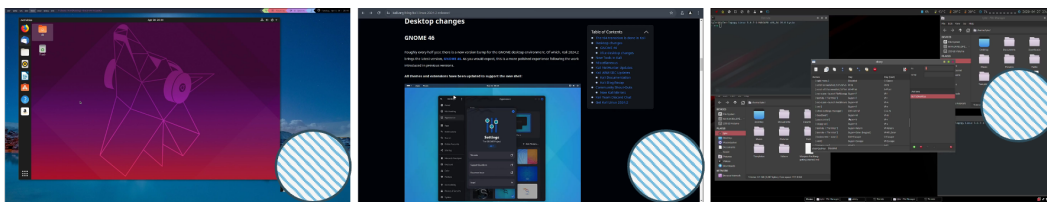
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024

C EXPLORATIONS: ALTERNATIVE DATA SOURCES AND ONLINE INTERACTION

Beyond the data collection pipelines used in the main text and Appendix E, we explored alternative data sources for CLI modeling. We did not incorporate them into the final pipeline, but the trials yielded useful insights and suggest directions for future work as tooling and data quality improve.

C.1 WEB VIDEO EXTRACTION

We tested mining terminal training data from web-scale screen video corpora. We used OCR and layout detectors to locate terminal regions and estimate text content and timestamps. We ultimately did not adopt this route for two reasons. First, privacy and copyright are central. Screen recordings often contain personal identifiers (usernames, emails, file paths, chat content) and come with licensing constraints that are difficult to verify at scale. Second, even when rights and privacy constraints are satisfied, cleaning the resulting data is substantially more complex than it appears (Fig. 36). Typical failure modes include (i) uncontrolled content such as faces/hands, picture-in-picture overlays, and unrelated desktop activity; (ii) domain shift across operating systems, themes, fonts, resolutions, and window managers; and (iii) quality factors such as compression artifacts, variable frame rates, zoom/crop edits, and inconsistent capture pipelines. These factors degrade OCR and temporal alignment.



(1) Uncontrolled Noise like “Human in Live Videos”, etc



(2) Different OS



(3) Hard to Control Quality

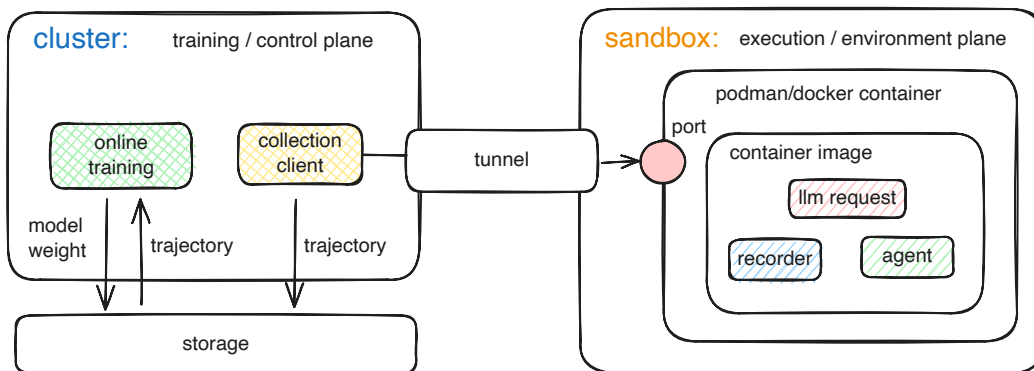
Figure 36: Common issues in web-scale screen video mining. Web videos mix uncontrolled content, heterogeneous OS/UI configurations, and inconsistent capture quality. This makes terminal localization, OCR, and time alignment unreliable without heavy filtering and sanitization.

Despite these challenges, web video mining may be a key lever for scaling interface experience over the long term. In this work, however, the cost-quality trade-off was unfavorable. Our setting benefits disproportionately from clean, temporally aligned text and interaction signals. Building a high-precision web filter also requires substantial upfront investment. This includes rights-cleared sourcing or licensing, privacy review and redaction, and large-scale multimodal filtering/OCR pipelines that often rely on paid APIs. Given these constraints and our emphasis on high-quality supervision, we

2079 prioritized curated CLIGen sources. Future efforts that invest in rights-respecting acquisition and
 2080 stronger automated filtering could unlock web-scale data as a complementary scaling axis.
 2081

2082 C.2 ONLINE ENVIRONMENT INTERACTION

2083
 2084 We prototyped an agentic online data and learning pipeline that separates a training/control plane
 2085 from an execution/environment plane (Fig. 37). In the environment plane, an isolated container
 2086 runs a live shell together with an agent and recorder exposed via a narrow port-based interface. The
 2087 agent is driven by the current model checkpoint (optionally assisted by an LLM planner) and issues
 2088 commands and control actions. The recorder captures synchronized terminal renders, structured
 2089 terminal state when available (e.g., buffer/text), and action traces. Trajectories are streamed to the
 2090 control plane for storage and optional online training, which updates model weights and pushes
 2091 refreshed checkpoints back to the agent.



2106 Figure 37: **Agentic online interaction pipeline (early exploration)**. A control plane collects
 2107 trajectories and can optionally update model checkpoints online. A sandboxed environment plane
 2108 executes an agent in an isolated container and records synchronized state/action traces.
 2109

2110 Concretely, the environment plane exposes a minimal “step/reset” interface over a port. It returns
 2111 multimodal observations that can be logged deterministically (rendered screenshots plus structured
 2112 state when available). The agent emits structured actions (typed command text, key/mouse events
 2113 when applicable, and timing). This separation makes rollouts auditable and lets the control plane
 2114 scale collection independently from training.
 2115

2116 We explored this setup because closed-loop interaction can induce a natural curriculum by continually
 2117 sampling the boundary of the current policy. It can also surface rare and safety-critical failure modes
 2118 that do not appear in offline logs. It supports targeted data collection (e.g., focusing on specific tools,
 2119 error recovery, or long-horizon tasks). In principle, it also offers a direct path to scaling *experience*
 2120 rather than only scaling static demonstrations.

2121 Early trials showed promise, but the end-to-end system introduced substantial engineering and
 2122 safety overhead. This includes strong isolation of untrusted code execution, monitoring and abuse
 2123 prevention, and deterministic resets and environment control. It also requires robust recording and
 2124 serialization across heterogeneous environments. Under time and cost constraints, we therefore
 2125 prioritized controlled CLIGen data for the main experiments.

2126 Despite not being used in the final pipeline, we consider this design a general template for future work.
 2127 It enables scalable multi-environment rollout collection and consistent provenance and storage of
 2128 trajectories. It also supports interchangeable learning algorithms (e.g., behavior cloning, preference-
 2129 based learning, or other online updates). Execution remains sandboxed and auditable.
 2130
 2131
 2132

2133 D IMPLEMENTATION DETAILS

2134

2135 This section summarizes training recipes. Dataset collection, preprocessing, and evaluation protocols
2136 are consolidated in Appendix E.

2137

2138 D.1 CLIGEN TRAINING RECIPE

2139

2140 Training uses gradient checkpointing and applies dropout 0.1 to the prompt encoder, CLIP, and VAE
2141 modules. Optimization uses AdamW (learning rate 5×10^{-5} , weight decay 10^{-2}), `bfloat16` pre-
2142 cision, and gradient clipping at 1.0. Training `NCCLIGen` on CLIGen (General) requires approximately
2143 15,000 H100 GPU hours at batch size 1. Training on CLIGen (Clean) across both subsets requires
2144 approximately 7,000 H100 GPU hours in total.

2145

2146 D.2 GUIWORLD TRAINING RECIPE

2147

2148 We train one model per injection mode (`external`, `contextual`, `residual`, `internal`),
2149 keeping the backbone and all non-action components fixed. Each run lasts about 64k steps using
2150 the CLI hyperparameters. We tune only the action encoder and learning-rate schedule. Training
2151 optimizes the diffusion loss together with a small temporal contrastive loss that aligns frame features
2152 with action and mouse embeddings (Appendix G). Runs use 64 GPUs for about 15 days, totaling
2153 about 23k GPU-hours per full pass.

2154

2155

2156

2157

2158

2159

2160

2161

2162

2163

2164

2165

2166

2167

2168

2169

2170

2171

2172

2173

2174

2175

2176

2177

2178

2179

2180

2181

2182

2183

2184

2185

2186

E DATASETS: COLLECTION AND EVALUATION PROTOCOLS

This appendix summarizes collection, preprocessing, and evaluation details for the datasets used in the paper. For concrete examples of raw trajectory formats (`asciinema .cast` and `vhs` scripts), see Appendix F.


Data collection follows three coordinated stages. This ensures the CLIGen and GUIWorld datasets share synchronized timing, privacy guarantees, and consistent, well-documented artifacts.


Sourcing. CLIGen (General) episodes originate from public `asciinema .cast` archives. We replay traces with official tools to preserve terminal appearance (color schemes, cursor visibility, window dimensions) as recorded. CLIGen (Clean) episodes come from deterministic `vhs` scripts (package installs, REPLs, log filters) run in isolated environments. GUIWorld footage is captured via the rig in Section 3.2, pairing RGB video with low-latency pointer/key logs and optional accessibility cues (logged for analysis; not used as model inputs).

Alignment and sanitization. All modalities share a common clock. We align pointer/key events to the nearest frame, apply drift correction when needed, and drop clips with residual misalignment. Privacy filters remove terminal sessions with sensitive strings and redact GUI regions likely to contain private content. Frozen or repeated-frame recordings (capture artifacts) are discarded.

Episode packaging. Runs are windowed into fixed-length, overlapping episodes (window sizes and strides are specified in the released configs). Each shard stores RGB frames, terminal buffers or GUI metadata, serialized actions, the source tool (`asciinema/vhs/GUI capture`), and environment metadata. Structured fields (buffers/metadata) are used for alignment and evaluation, but are not provided to the video models as state inputs. Downstream dataloaders reconstruct batches directly. Released configs specify preprocessing and windowing so external users can rebuild the corpus.

E.1 / CLIGEN DATA PIPELINE

The  CLIGen (General) dataset is built from publicly available `asciinema .cast` trajectories³. The `asciinema` stack records and replays terminal sessions with synchronized timing and ANSI-faithful decoding. We replay each session with the official tools and render it into terminal frames, preserving palette transitions, cursor state, and terminal geometry. Frames, text buffers, and keyboard-event logs share a single monotonic clock. At render time, we normalize resolution and aspect ratio and apply a simple privacy filter to remove sensitive strings. We render sessions to GIF using `agg` and convert them to video with `ffmpeg`. We segment each recording into roughly five-second clips using content-aware splits. We temporally normalize each clip to a fixed length: shorter clips repeat the final frame, and longer clips are uniformly subsampled. The resulting 823,989 video streams (approximately 1,100 hours) are resampled to 15 FPS. Underlying buffers and logs are used to generate aligned textual descriptions with Llama 3.1 70B (Dubey et al., 2024) in three styles (semantic, regular, and detailed), which serve as prompts. As shown in Figure 2 (left), this split spans diverse real-world terminal use cases.⁴

The  CLIGen (Clean) dataset is collected using the open-source `vhs`⁵ toolkit. It enables repeatable terminal demonstrations and integration tests through scripted execution. Deterministic scripts drive Dockerized environments to capture cleaner, better-paced traces. We authored roughly 250k scripts. After filtering (51.21% retained), we keep two subsets. The first contains approximately 78k regular traces (package installation, log filtering, interactive REPL usage, etc.). The second contains approximately 50k Python math validation traces. Captions are derived directly from the raw `vhs` scripts for clarity. We standardize frame rendering by fixing one monospace font/size, using a consistent palette for success and error highlights, and locking resolution and theme to remove typography-related confounds. Each episode records its caption type and font settings for later slicing. Clips longer than five seconds are uniformly subsampled for training, while shorter clips repeat the


³<https://asciinema.org/>

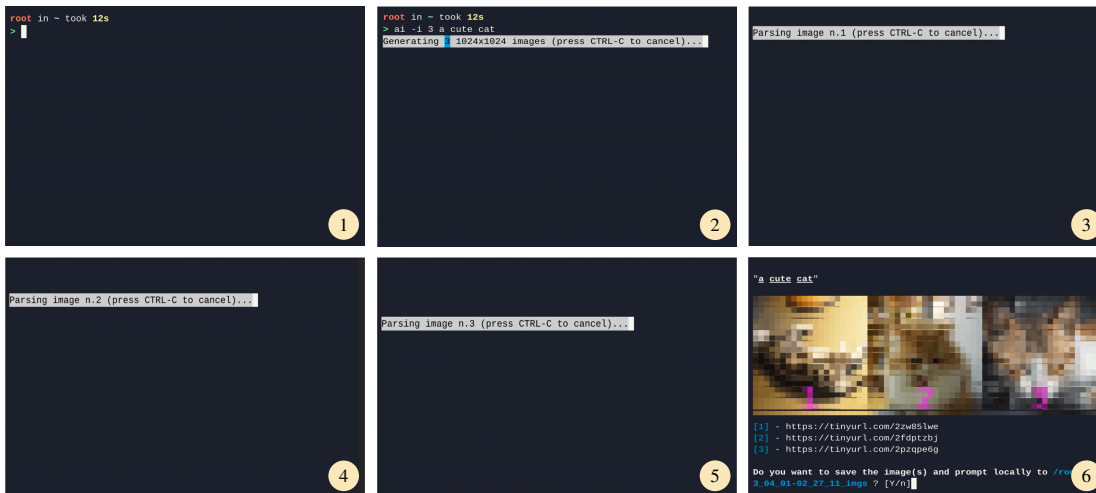
⁴Additional preprocessing details and a `.cast` example are in Section F.1, with a sample overview in Table 9.

⁵<https://github.com/charmbracelet/vhs>

2241 final frame to normalize length. We apply the same timestamping, resolution/aspect normalization,
 2242 and privacy filtering as in CLIGen (General).⁶
 2243

2244 Table 9: NC_{CLIGen} samples for  CLIGen (General) and  CLIGen (Clean).

2245
 2246 **Frames from Sample —  CLIGen (General)**



2264 **Semantic**


2265 A root terminal session kicks off an AI command to make three 1024x1024 cat shots, shows quick parsing
 2266 for each one, then presents pixelated cat previews with numbered links and asks whether to stash them in
 2267 /root/2023-04-01-02-27-11.imgs.

2268 **Regular**

2269 In a root shell at ~, the user runs ai -i 3 a cute cat, watches a green progress line announcing
 2270 three 1024x1024 images, sees sequential parsing messages for images 1 through 3, and ends on a pre-
 2271 view pane with three pixelated cat thumbnails, numbered download links, and a save prompt targeting
 2272 /root/2023-04-01-02-27-11.imgs.

2273 **Detailed**

2274 In a dark-background terminal at the root in ~ prompt, the user types ai -i 3 a cute cat. The screen
 2275 prints Generating 3 1024x1024 images (press CTRL-C to cancel)..., shows parsing mes-
 2276 sages for images 1-3, and ends on a preview pane with three numbered thumbnails and a save prompt targeting
 2277 /root/2023-04-01-02-27-11.imgs.

2278 **Frames from Sample —  CLIGen (Clean)**



2291 **Scripted caption**

2292 Type python; Enter; Type values = [n*n for n in range(1, 10)]; Enter; Type
 2293 print (values); Enter; Type exit (); Enter.

2294 ⁶Additional details are provided in Section F.2, with a representative data sample in Table 9.

E.2 CAPTION FIELDS AND METADATA (CLIGEN GENERAL)

For CLIGen (General), each replayed `.cast` fragment (example shown in Section F.1) is paired with three aligned descriptions and a compact metadata record. Table 10 summarizes the fields for clip `7****_0001`.

Table 10: Caption tiers and metadata for CLIGen (General) clip `7****_0001`.

Field	Content (abridged)
caption	The user, <code>anonnomous</code> , logged into their account on the host <code>NeuralComputer</code> , and after some initial terminal setup and cursor movements, they started typing the command <code>nvim</code> to likely open the Neovim text editor, but the session recording ends abruptly without showing the actual command execution or any further interaction.
caption_detailed	At the <code>anonnomous@NeuralComputer</code> : <code>\$</code> prompt, the user types <code>nvim</code> , but before the command executes, the screen rapidly displays a colorful, pixelated animation with various RGB colors, including shades of blue, green, and purple, filling the 80x24 terminal window. The animation is briefly interrupted by a cursor blink, then the prompt returns, awaiting further input.
caption_semantic	The terminal displays a rapid sequence of colored text and pixel art animations, with syntax-highlighted code and prompts appearing in a mix of blue, green, and white hues, as the user types commands and the system responds with scrolling output, error messages, and success indicators.
<code>data.info.version</code>	2 (asciinema v2 header)
<code>data.info.size</code>	<code>width=80,height=24;original/target/scaled sizes all match, padding [0,0].</code>
<code>data.info.env</code>	<code>SHELL=/bin/bash,TERM=xterm-256color.</code>
<code>meta.videogen</code>	Automatically derived stats: visual complexity 2897.7, interaction density 1.0, color usage 518, screen clears 6, event rate 7.8 events/s, typing rhythm (avg interval 0.13s, variance 0.34), and 38 visual-change events.
metadata	Source recording on Asciiinema: id <code>7****</code> , title <code>'`aneo.nvim demo'`</code> , author <code>anonnomous</code> , creation time <code>2025-05-13T23:37:53Z</code> , URLs for the page and raw <code>.cast</code> .

These fields make each CLIGen (General) clip self-contained. The three caption tiers provide prompts at different levels of detail. `data.info` and `metadata` preserve the structure needed to rebuild terminal geometry, environment, and source from the raw `.cast`.

A rough analysis of content categories shows the dataset covers diverse terminal usage. It includes general terminal operations (610,331 clips), file operations (81,783), programming activities (66,430), system administration (42,668), and text editing (22,777). This distribution reflects typical command-line workflows across different user tasks.

E.3 OCR EVALUATION PROTOCOL (CLIGEN CLEAN)

For CLIGen (Clean), OCR-based metrics evaluate how closely generated terminal videos match reference renderings derived from the ground-truth buffers in text space rather than pixels. Each sample consists of a generated video and its paired reference video (matched by clip ID). We keep only IDs where both videos are present.

From each paired video we use at most $K=5$ frames. Let T_{gen} and T_{gt} be the frame counts of the generated and reference videos. We set $T = \min(T_{\text{gen}}, T_{\text{gt}})$ and uniformly sample K indices in $[0, T-1]$ so that evaluation frames are spread across the trajectory. For every sampled index we read the corresponding frame from both videos.

Each frame is converted to RGB and passed to Tesseract OCR. The resulting string is split into lines, leading and trailing whitespace is stripped, and internal whitespace is normalized by collapsing runs of spaces. Empty lines are dropped. We keep case and punctuation intact so that commands, paths, and symbols remain visible. This gives an ordered list of normalized lines for the ground-truth frame (g_1, \dots, g_{N_g}) and the generated frame (p_1, \dots, p_{N_p}) .

2349 *Character accuracy* pools all lines into a single multi-line string for each side and measures normal-
 2350 ized edit distance. Let s and t be the concatenated ground-truth and generated texts and $d(s, t)$ their
 2351 Levenshtein distance (insert/delete/replace cost 1). If both s and t are empty we set $\text{char_acc} = 1$; if
 2352 only s is empty we set $\text{char_acc} = 0$. Otherwise,

$$2353 \text{char_acc} = \max\left(0, 1 - \frac{d(s, t)}{\max(|s|, 1)}\right),$$

2354
 2355
 2356 Extra or missing characters are penalized symmetrically. Frame-level scores are averaged over the K
 2357 sampled frames to yield a per-video character accuracy, and group-level scores report the mean over
 2358 videos.

2359 *Exact-line accuracy* treats lines as position-sensitive units and reports a recall over ground-truth
 2360 lines. For a given frame, we compare line g_i to p_i at the same index. A line is counted as correct
 2361 only if $i \leq N_p$ and $p_i = g_i$; lines that appear in the wrong position do not count. If both lists are
 2362 empty we set $\text{exact_line_acc} = 1$; if the ground-truth list is empty but the generated list is not, we set
 2363 $\text{exact_line_acc} = 0$. Otherwise,

$$2364 \text{exact_line_acc} = \frac{1}{N_g} \sum_{i=1}^{N_g} \mathbf{1}[i \leq N_p \wedge p_i = g_i].$$

2365
 2366
 2367
 2368 As with character accuracy, frame scores are averaged over the K sampled frames to obtain a per-
 2369 video score and then averaged over videos for the reported aggregate. Unless otherwise noted, OCR
 2370 results in the main text use 1,000 randomly sampled video pairs (five frame pairs per video). Together,
 2371 these two metrics stress both fine-grained text fidelity and line-ordered terminal state reconstruction.

2372
 2373 **E.4  GUIWORLD DATA PIPELINE**

2374
 2375 Preprocessing is implemented in the data loader in two stages. First, we normalize each recording to
 2376 a fixed resolution and frame rate. This produces tensors for RGB video, per-frame cursor coordinates,
 2377 and mouse/keyboard event traces (in both raw-action and meta-action views). Second, we render an
 2378 SVG cursor at each logged position to produce per-frame masks and cursor-only reference frames.
 2379 The first reference frame contains the full desktop with a unit mask. Later references paste only the
 2380 cursor over a neutral background, with a mask restricted to arrow pixels. After VAE encoding, these
 2381 references become latent slots that pin down the static GUI layout at $t=0$. For $t>0$, they supervise
 2382 only a small patch around the cursor and leave the rest of the frame unconstrained. We drop clips
 2383 without valid cursor or action traces to keep supervision consistent.

2384 **Cursor coordinate normalization.** Raw mouse logs provide per-frame cursor positions in screen
 2385 coordinates $(x_{\text{screen}}, y_{\text{screen}})$ at GUI resolution. We align these to sampled video frames and map them
 2386 through the same letterbox transform applied to the RGB stream. Given source and target resolutions
 2387 $(w_{\text{src}}, h_{\text{src}})$ and $(w_{\text{dst}}, h_{\text{dst}})$, we compute a uniform scale s and padding offsets (p_x, p_y) . We then map
 2388 each coordinate to normalized positions $(x_t, y_t) \in [0, 1]^2$ as

$$2389 x_t = \frac{s x_{\text{screen},t} + p_x}{w_{\text{dst}} - 1}, \quad y_t = \frac{s y_{\text{screen},t} + p_y}{h_{\text{dst}} - 1}.$$

2390
 2391
 2392 Stacking these over time yields a trajectory tensor $\text{mouse_trajectories} \subset [0, 1]^{T \times 2}$.

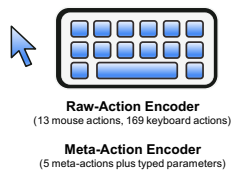
2393 **Fourier trajectory features (optional).** To provide a stronger continuous conditioning signal, we
 2394 optionally encode the normalized trajectories with random Fourier features. We map coordinates to
 2395 $[-1, 1]^2$ and project them through a fixed Gaussian matrix (followed by sine/cosine). A small MLP
 2396 produces per-frame embeddings, which are aggregated into lag-aware windows aligned with the VAE
 2397 stride. The resulting latent-aligned mouse features can condition the action modules and participate
 2398 in the temporal contrastive loss (see Appendix G).

2399 Table 11: Cursor/action statistics.

2400 Split	Avg. cursor speed (px/frame)	Actions / sec
2401 Random Slow	1.51	1.58
2402 Random Fast	195.15	4.18
CUA (supervised)	3.79	0.10

E.5 ENCODING EXAMPLES (RAW-ACTION VS. META-ACTION)

Table 12: Encoding examples for raw-action and meta-action encoders.



User intent	Raw-action encoder (event stream)	Meta-action encoder (API-like slot)
<code>ls -l</code>	<ul style="list-style-type: none"> Key events per character (e.g., <code>l</code>, <code>s</code>, <code><space></code>, <code>-</code>, <code>l</code>) Activates entries in a 169-d keyboard multi-hot No explicit command semantics; inferred from sequence 	<ul style="list-style-type: none"> <code>type:</code> KeyboardType <code>text:</code> "ls -l" Encoded by shared text encoder
<code>ctrl+v</code>	<ul style="list-style-type: none"> Separate keydown/keyup events for <code>ctrl</code> and <code>v</code> Activated in multi-hot (or shortcut entry) 	<ul style="list-style-type: none"> <code>type:</code> Shortcut <code>id:</code> <code>ctrl+v</code> Embedded via shortcut table

We collect GUI data at 15 FPS with synchronized cursor and keyboard logs. The dataset includes two styles of random interaction: “Random Slow” and “Random Fast”, plus a smaller set of supervised trajectories from Claude CUA (Rivard et al., 2025). Random Slow (approximately 1,000 hours) contains longer pauses, idle gaps, and deliberate cursor movements, which can expose cursor drift after extended inactivity. Random Fast (approximately 400 hours) features denser cursor motion and typing bursts, stressing acceleration dynamics and hover timing.

We also include approximately 110 hours of supervised trajectories generated using Claude CUA⁷. These goal-directed traces provide higher-signal action–response pairs without overwhelming the exploration data. Table 11 summarizes cursor and action statistics across splits; CUA is slowest due to latency introduced by Claude’s tool API between successive actions.

All GUI data is collected inside an Ubuntu 22.04 container running XFCE4 (Arc-Dark theme, Papirus icons) on a fixed 1024×768 virtual display at 15 FPS. We render the display with Xvfb and interact through a VNC/noVNC stack. The desktop pins a small open-source app set to launchers. It includes Firefox ESR, GIMP, VLC, VS Code, Calculator, Terminal, the file manager, and the Mahjongg game, matching the environment shown in our recordings. Screen capture uses `mss` and `ffmpeg` with cursor overlays, and actions are replayed and logged via `xdotool`. We keep the recorded discontinuities and interface latency intact rather than smoothing them. For modeling, we store both raw-action and meta-action views. This lets the same encoder and loss stack train on either.⁸

E.6 EVALUATION METRICS AND PROTOCOL (GUIWORLD)

We report both global video metrics and action-driven metrics that focus on post-interaction frames. The evaluation code used for these metrics is implemented in our GUIWorld evaluation suite.

Global FVD_{all}/SSIM_{all}/LPIPS_{all} We decode paired generated/ground-truth videos into RGB frames with temporal subsampling and resizing (`fps=3`, `size=256`, and `max_seconds=5` by default). `SSIMall` is computed using `torchmetrics` on frame tensors normalized to $[0, 1]$ and averaged over frames. `LPIPSall` uses the AlexNet backbone on frames normalized to $[-1, 1]$ and is averaged over frames. `FVDall` is computed in an `r3d18` embedding space (`prelogits` by default). We extract features from fixed-length clips (16 frames at 112×112 after uniform subsampling/padding). We compute the Fréchet distance between the generated and reference feature distributions.

Action-driven metrics (SSIM₊₁₅, LPIPS₊₁₅, FVD₊₁₅) For each paired rollout, we load recorded action timestamps (from JSON/CSV logs) and map them to frame indices at the evaluation frame rate. We skip the action frame itself (`action_start_offset=1`) and evaluate the next $k=15$

⁷<https://platform.claude.com/docs/en/agents-and-tools/tool-use/computer-use-tool>

⁸Conversion details and alignment quality appear in Appendix G.

2457 frames after each action. Concretely, for each action frame index f , we select the post-action set
2458 $\{f+1, \dots, f+k\}$. We take the union over all actions in the clip and compute $SSIM_{+15}/LPIPS_{+15}$ on
2459 the selected frame pairs. For action-driven FVD_{+15} , we build an *after-action clip* per video by
2460 concatenating the selected post-action frames. We then uniformly subsample/pad to 16 frames and
2461 compute the Fréchet distance in the same `r3d18` feature space.

2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510

2511 F CLIGEN: CLI TRAJECTORY FORMATS

2512

2513 This appendix provides concrete format examples for the two CLI sources referenced throughout the
2514 data pipeline (Appendix E). We show asciinema `.cast` trajectories for CLIGen (General) and vhs
2515 scripts for CLIGen (Clean).
2516

2517 F.1 ASCIINEMA (`.CAST`) EXAMPLE

2518

2519 The header line stores the recording config (version, terminal size, timestamp, env). Each following
2520 triple is [time, "o", "payload"], where "o" indicates screen output. The payload contains
2521 the terminal text with color codes at that timestamp.

```
2522 {"version": 2, "width": 80, "height": 24,  
2523   "timestamp": 1747177906,  
2524   "env": {"SHELL": "/bin/bash", "TERM": "xterm-256color"}}  
2525 [0.082492, "o", "\u001b[H\u001b[2J\u001b[3J"  
2526 [0.950038, "o", "\u001b[38;2;16;131;236m\u001b[39m\r\n..."  
2527 [0.950733, "o", "\u001b[38;2;6;156;220m ... \u001b[38;2;1;195;187m█"]
```

2528

2529 F.2 VHS SCRIPT EXAMPLE

2530

```
2531 #---- VHS documentation start (DO NOT CHANGE) ----
```

```
2532 # Require:
```

```
2533 #   Require <string>
```

```
2534 # Sleep:
```

```
2535 #   Sleep <time>
```

```
2536 # Type:
```

```
2537 #   Type[@<time>] "<characters>"
```

```
2538 # Keys:
```

```
2539 #   Escape[@<time>] [number]
```

```
2540 #   Backspace[@<time>] [number]
```

```
2541 #   Delete[@<time>] [number]
```

```
2542 #   Insert[@<time>] [number]
```

```
2543 #   Down[@<time>] [number]
```

```
2544 #   Enter[@<time>] [number]
```

```
2545 #   Space[@<time>] [number]
```

```
2546 #   Tab[@<time>] [number]
```

```
2547 #   Left[@<time>] [number]
```

```
2548 #   Right[@<time>] [number]
```

```
2549 #   Up[@<time>] [number]
```

```
2550 #   PageUp[@<time>] [number]
```

```
2551 #   PageDown[@<time>] [number]
```

```
2552 #   ctrl+<key>
```

```
2553 # Display:
```

```
2554 #   Hide
```

```
2555 #   Show
```

```
2556 # ---- VHS documentation end (DO NOT CHANGE) ----
```

2557

```
2558 # ID: vhs_example
```

```
2559 # INSTRUCTION: Runs `uname -s` repeatedly as a basic shell exercise, then hides the p
```

```
2560 # LEVEL: 1
```

```
2561 # EVENTS: 23
```

```
2562 # VISUAL_COMPLEXITY: 45
```

2563

```
2564 # ---- Theme setting start (DO NOT CHANGE) ----
```

```
2565 Output vhs_example.mp4
```

2566

```
2567 Set Shell "bash"
```

2568

```

2565 Set Theme {
2566   "name": "Catpuccin Mocha (Pure White, Warm Pink Cursor)",
2567   "background": "#1e1e2e",
2568   "foreground": "#ffffff",
2569   "black": "#45475a",
2570   "red": "#f38ba8",
2571   "green": "#a6e3a1",
2572   "yellow": "#f9e2af",
2573   "blue": "#89b4fa",
2574   "purple": "#cba6f7",
2575   "cyan": "#94e2d5",
2576   "white": "#ffffff",
2577   "brightBlack": "#585b70",
2578   "brightRed": "#f38ba8",
2579   "brightGreen": "#a6e3a1",
2580   "brightYellow": "#f9e2af",
2581   "brightBlue": "#89b4fa",
2582   "brightPurple": "#cba6f7",
2583   "brightCyan": "#89dceb",
2584   "brightWhite": "#ffffff",
2585   "cursor": "#f5c2e7",
2586   "cursorAccent": "#1e1e2e",
2587   "selectionBackground": "#585b70"
2588 }
2589
2588 Set FontSize 40
2589 Set Width 1600
2590 Set Height 900
2591 Set TypingSpeed 300ms
2592 Set PlaybackSpeed 1
2593 Set Margin 28
2594 Set MarginFill "#0091FF"
2595 Set BorderRadius 25
2596 Set Padding 18
2597 Set LineHeight 1.2
2598 Set LetterSpacing 0.8
2599 # ---- Theme setting end (DO NOT CHANGE) ----
2600
2600 Sleep 800ms
2601 Sleep 180ms
2602 Type "uname -s"
2603 Sleep 120ms
2604 Enter
2605 Sleep 400ms
2606 Type "uname -s"
2607 Sleep 120ms
2608 Enter
2609 Sleep 400ms
2610 Type "uname -s"
2611 Sleep 120ms
2612 Enter
2613 Sleep 400ms
2614 Type "uname -s"
2615 Sleep 120ms
2616 Enter
2617 Sleep 400ms
2618 Type "uname -s"
2619 Sleep 120ms
2620 Enter

```

2619 Sleep 400ms
 2620 Sleep 400ms
 2621 Sleep 600ms
 2622 Hide

2623
 2624
 2625
 2626

G GUIWORLD: ACTION REPRESENTATION, CONDITIONING AND MORE EXPERIMENTS

2627
 2628
 2629
 2630

This appendix provides additional GUIWorld visualization samples and expands on the action representation and conditioning used in Section 3.2. Evaluation metrics and protocols are summarized in Appendix E.6.

2631
 2632

G.1 ACTION SCHEMA

2633
 2634

$NC_{GUIWorld}$ represents actions as a structured stream, enabling the NC to condition on both cursor movements and key presses.

2635
 2636
 2637
 2638
 2639
 2640
 2641
 2642
 2643
 2644
 2645
 2646

At each timestep, we log absolute cursor coordinates, button up/down transitions, scroll deltas, and keyboard events. Keyboard inputs are split into two types: typed characters (e.g., `ls -l`) and shortcut-style chords (e.g., `ctrl+v`). We also track state flags such as whether a drag is currently active. This lets us represent extended interactions like click-drag or press-hold as short labeled segments rather than isolated spikes. The meta-action encoder described in Section 3.2 compresses this stream into a small typed schema. Each frame can contain up to two actions, and each action has a type (e.g., `mouse click` or `keyboard type`) plus parameters. Table 13 summarizes the types and fields. Type 0 corresponds to the absence of an action. Type 1 encodes mouse clicks and drags via button identity, click count, and a drag flag. Type 2 captures scrolls with a direction and scalar amount. Type 3 packages free-form keyboard text (such as `ls -l`) embedded by the shared text encoder. Type 4 records shortcuts such as `ctrl+v` via a small shortcut vocabulary. This representation resembles a tool API while remaining recoverable from raw logs.

2647
 2648

Table 13: Meta-action schema for GUIWorld (per action slot).

2649
 2650
 2651
 2652
 2653
 2654

Type id	Action	Parameter fields
0	None	-
1	Mouse Click/Drag	<code>button</code> , <code>click_count</code> , <code>drag_flag</code>
2	Mouse Scroll	<code>direction</code> , <code>amount</code>
3	Keyboard Type	<code>text</code> (e.g., <code>ls -l</code>) \rightarrow shared text encoder
4	Shortcut	<code>shortcut_id</code> (e.g., <code>ctrl+v</code>)

2655
 2656
 2657

G.2 CONDITIONING: ENCODERS AND INJECTION

2658
 2659
 2660
 2661
 2662
 2663
 2664
 2665

The main text considers two encoders for this stream. A *raw-action encoder* ($v1$) keeps fine-grained mouse and key events in a multi-hot representation that closely mirrors real cursor and typing behavior. A complementary *meta-action encoder* ($v2$) compresses events into a small typed schema (Table 13) and embeds any free-form text with a shared text encoder. Both encoders produce per-frame action features that undergo temporal windowing and alignment (described below). These embeddings support four injection modes. `external` fuses actions at the VAE input. `contextual` mixes actions and frames as tokens. `internal` injects actions inside transformer blocks. `residual` adds lightweight action deltas to hidden states.

2666
 2667
 2668
 2669

Injection-mode definitions (formal) Below we give compact schematic definitions for the three modes that are described with formulas in the main text. **External.** Given VAE latents $z_{1:T}$ and temporally aligned action features $u_{1:T}$, an external action module produces a residual update $\Delta z_{1:T}(u_{1:T})$ and forms modified latents

2670
 2671
 2672

$$z'_{1:T} = z_{1:T} + \Delta z_{1:T}(u_{1:T}).$$

The diffusion backbone then operates on $z'_{1:T}$ (actions do not appear as explicit tokens inside the transformer). **Residual.** At selected transformer layers l , an auxiliary action module takes

block hidden states $h^{(l)}$ together with local action/mouse features and outputs a residual update $\Delta h^{(l)}(a, \text{mouse})$. The updated hidden states are

$$\tilde{h}^{(l)} = h^{(l)} + \Delta h^{(l)}(a, \text{mouse}),$$

which are passed to the next block. **Internal.** At selected blocks, action conditioning is inserted as an additional cross-attention sub-layer inside the standard attention stack. With self-attention SA, text/reference cross-attention CA_{text} , and action cross-attention CA_{action} , a schematic update is

$$h' = \text{FFN}\left(h + CA_{\text{text}}(\text{SA}(h), c) + CA_{\text{action}}(h, a)\right).$$

G.3 TEMPORAL ALIGNMENT AND ATTENTION

Temporal alignment and windows. The GUI backbone processes a compressed latent video at stride c (every c pixel frames correspond to one latent frame). For a pixel sequence of length F and latent sequence of length T , we have $F = (T - 1)c + 1$. Anchor frame $a_t = t \cdot c$ marks the pixel frame corresponding to latent step t .

Mouse and keyboard logs start as per-frame features $r_f \in \mathbb{R}^D$ at the pixel rate. A windowed encoder aggregates them around each anchor over $p = c \cdot w$ frames. Here w controls the window width, and a lag ℓ accounts for GUI response delay (actions precede their visual effects). We use zero-padding outside the valid range, i.e., $\tilde{r}_f = r_f$ for $0 \leq f < F$ and $\tilde{r}_f = 0$ otherwise, and form a lag-shifted window that ends at $a_t - \ell$:

$$W_{t,k} = \tilde{r}_{a_t - (p-1+\ell)+k}, \quad k \in \{0, \dots, p-1\}, \quad a_t^{\text{act}} = \frac{1}{p} \sum_{k=0}^{p-1} W_{t,k}.$$

This shared action encoder yields one latent-aligned action embedding a_t^{act} per step. It summarizes a short, lagged history of cursor motion and key events and is reused across all injection modes.

Contextual attention mask. In the `contextual` mode, video and action tokens are concatenated into a single sequence and processed under a structured, causal-local attention mask. Appendix Figure 38 can be read as a query-key matrix: rows are queries and columns are keys.

The upper-left block (V2V) restricts each frame V_i to attend only to neighboring frames within a window of $\pm w$ steps, so very distant frames cannot interfere. The upper-right block (V2A) lets frame V_i see only the most recent actions in a short lag window. Typically this includes actions A_j with $i - (w + \ell) \leq j \leq i - \ell$, so actions that are too recent to have taken effect are masked. This way, a frame is explained by at most the last few perceivable operations and never by future ones. In the lower-left block (A2V), an action A_i attends to frames V_t that occur after it has had time to take effect (roughly $t \geq i + \ell$), but not to earlier frames. The lower-right block (A2A) is near-diagonal: actions mostly attend to themselves rather than to unrelated actions.

In practice (w, ℓ) act as fixed hyperparameters that trade off temporal coverage and cost. Together, these choices implement a structured causal prior: actions do not explain past frames, and each frame conditions only on the recent operations that could plausibly have shaped its pixels.

Design insights. Two insights from the GUI experiments motivate this schema. First, raw action streams are bursty and high-dimensional. Cursor and key events arrive in short spikes, and simple smoothing or full-history attention can cause false interpolated motion and underestimated typing speed. Using short, lagged windows and local attention bands makes credit assignment more intuitive: each frame connects to the few operations that could have produced it. Second, control fidelity improves more from better conditioning than from changing the visual backbone. Clean, well-paced supervision and mid- or deep-level action injection improve cursor accuracy and hover timing, while different encodings of the same stream perform similarly. This action schema and mask implement these principles: keep pixels and actions aligned in time, prioritize recent operations over distant ones, and use attention structure rather than capacity alone.

G.4 CURSOR RENDERING AND SUPERVISION

Cursor rendering and reference construction.

2727 The cursor pipeline applies the same design principles on the visual side. Instead of relying on the
2728 global diffusion loss to discover a tiny arrow, we render the cursor explicitly. We treat it as a first-class
2729 conditioning signal.

2730 **From logs to normalized trajectories.** Desktop logs provide per-frame cursor positions in screen
2731 coordinates $(x_{\text{screen}}, y_{\text{screen}})$ at the native GUI resolution. We align these with sampled video frames
2732 using the same letterbox mapping as the RGB stream. This includes rescaling to target resolution,
2733 adding symmetric padding, and normalizing by $(\text{width} - 1, \text{height} - 1)$ to obtain $(x_t, y_t) \in [0, 1]^2$
2734 for each timestep t . Stacking these over time yields the trajectory tensor `mouse_trajectories`
2735 used across rendering and action encoding.

2736

2737

2738

2739

2740

2741

2742

2743

2744

2745

2746

2747

2748

2749

2750

2751

2752

2753

2754

2755

2756

2757

2758

2759

2760

2761

2762

2763

2764

2765

2766

2767

2768

2769

2770

2771

2772

2773

2774

2775

2776

2777

2778

2779

2780

Table 14: Raw-action versus meta-action encoders in GUIWorld.

Aspect	Raw-action encoder (v1)	Meta-action encoder (v2)
Action schema	Per-frame multi-hot vector composed of 13 mouse actions and 169 keyboard actions; no explicit type hierarchy.	Hierarchical schema with explicit action types and parameters: $action_types \in \{0, 1, 2, 3, 4\}$ for $\{\text{None, Mouse Click/Drag, Mouse Scroll, Keyboard Type, Keyboard Shortcut}\}$, plus typed parameter fields (e.g., button, scroll amount, shortcut ID, text).
Mouse representation	Two signals per frame: (1) continuous cursor trajectory $mouse_trajectories_t \in \mathbb{R}^2$; (2) discrete mouse events $mouse_action_events_t \in \{0, 1\}^{13}$ (multi-hot).	Per-frame mouse trajectory $mouse_trajectories_t$ is encoded by a shared mouse-trajectory encoder and the temporal windowing module into dense mouse embeddings $mouse_latent_t \in \mathbb{R}^{d_{mouse}}$, optionally fused with action embeddings.
Keyboard representation	Per-frame keyboard state $keyboard_action_events_t \in \{0, 1\}^{169}$, where each dimension corresponds to a key or shortcut; only multi-hot activation is available, no text semantics.	Two complementary forms: keyboard shortcuts as <code>keyboard_shortcut</code> IDs embedded via an embedding table, and free-form text <code>keyboard_text</code> (e.g. "ls -l") encoded by a shared text encoder (e.g., T5) and projected to the action embedding dimension.
Per-frame representation	Mouse and keyboard events are concatenated: $raw_actions_t = [mouse_events_t, keyboard_events_t]_{\mathbb{R}^{182}}$.	For each frame t and slot s , type plus parameters are embedded into $slot_embeds_{t,s} \in \mathbb{R}^D$. Valid slots are aggregated (mean or attention) into a single per-frame action embedding $action_embeds_t \in \mathbb{R}^D$.
Multi-slot support	No explicit notion of multiple action slots per frame; all mouse and keyboard signals are collapsed into one 182-d multi-hot vector.	Explicit multi-slot design: <code>action_types</code> and all parameter tensors have shape (B, T, S) with $S = 2$ slots per frame, enabling, e.g., concurrent "mouse click + keyboard shortcut" in the same frame.
Latent-aligned representation	Temporal alignment is implemented ad-hoc inside each v1 action module (e.g., sliding windows over the 182-d per-frame vector), with custom logic per mode (external, internal, residual).	A unified temporal alignment module takes per-frame action embeddings (B, F, D) and VAE compression ratio c , constructs lag-aware temporal windows of size $c \cdot w$, and outputs latent-aligned action features $action_latent \in \mathbb{R}^{B \times T \times D}$, with analogous <code>mouse_latent</code> when mouse features are enabled.
Lag modeling	GUI latency / response lag is implicitly encoded by where and how the sliding window is applied inside each v1 action module; there is no shared lag configuration across modes.	Lag is an explicit hyperparameter <code>action_lag</code> in the temporal alignment and in the contextual attention mask: it shifts temporal windows and defines which frames an action can attend to, giving a consistent lag interpretation across modes.
Downstream usage	Raw 182-d vectors (plus trajectories) are fed directly into the v1 action modules (external, internal, residual) without a shared temporal encoder.	The v2 encoder outputs provide a shared conditioning stream for all v2 modes and also supply features for the temporal contrastive loss.
Contrastive supervision	No dedicated temporal contrastive loss between frames and actions; supervision comes only from the generative objective.	Action and mouse latents are coupled to frame features via an InfoNCE-style temporal contrastive loss, encouraging time-aligned action embeddings.

2835 G.5 TRAINING SIGNALS AND ENCODER DESIGN
2836

2837 **From trajectories to cursor layers.** Starting from normalized (x_t, y_t) coordinates, a cursor-layer
2838 module renders a fixed SVG arrow template into RGB and alpha channels. The template is reused
2839 across frames. For each timestep t , it is positioned so the hotspot (arrow tip) aligns with (x_t, y_t) ,
2840 clipped to screen bounds, and alpha-blended over a neutral background. This produces two tensors
2841 at video frame rate: a cursor-only foreground image $f_t \in [-1, 1]^{3 \times H \times W}$ and a soft mask $m_t \in$
2842 $[0, 1]^{1 \times H \times W}$ isolating the arrow pixels. Invalid or missing coordinates zero the mask and leave the
2843 foreground unchanged, so frames without visible cursors do not add spurious supervision.

2844 **Reference images and masks.** These cursor layers become reference conditions for the I2V model.
2845 For each clip we form reference images $\text{ref_img}_{0:T-1}$ and masks $\text{ref_mask}_{0:T-1}$:
2846

- 2847 • at $t=0$, ref_img_0 is the full desktop frame and ref_mask_0 is all ones, anchoring the static layout
2848 and theme;
- 2849 • for $t>0$, ref_img_t is the cursor foreground f_t and ref_mask_t is the cursor mask m_t , so only the
2850 arrow region is supervised while the background remains free.

2851 The model encodes these references with the same VAE as target frames and concatenates their
2852 latents and masks into the diffusion input. It learns to “copy from ref” inside the cursor mask and
2853 relies on its dynamics model elsewhere. This makes cursor supervision a pixel-level constraint rather
2854 than a side effect of the global loss.

2855 **Fourier mouse encoding.** The same (x_t, y_t) trajectories serve as a continuous control signal. We
2856 apply a Fourier position module: clamp coordinates to $[0, 1]^2$, map them to $[-1, 1]^2$, and compute
2857 random Fourier features via a fixed Gaussian projection followed by sine/cosine. A small MLP
2858 maps these features to per-frame mouse embeddings. The GUIWorld action encoder then aggregates
2859 them with lag-aware, stride-aligned windows to produce latent-aligned mouse features. These
2860 features condition the `external/contextual/residual/internal` modes and participate in
2861 the temporal contrastive loss.

2862 **Cursor-aware losses.** Section 3.2 introduces cursor-aware losses that use this construction. A basic
2863 variant penalizes position error only in (x, y) . Richer variants add Fourier features of the trajectory
2864 and, most importantly, an ℓ_2 loss on the reconstructed cursor patch under ref_mask_t . Table 6 shows
2865 that position-only objectives yield low cursor hit rate and visibly jittery arrows even when videos
2866 look plausible. Adding the explicit cursor reference stream together with the masked patch loss
2867 substantially improves control, reaching 98.7% cursor accuracy. This confirms that explicit cursor
2868 rendering plus localized supervision effectively separates “where the arrow is” from “what the rest of
2869 the frame should look like”.

2870 **Temporal contrastive alignment.** To strengthen learning signals for the action pathway, we add
2871 a lightweight temporal contrastive loss that operates on the same latent timeline as the diffusion
2872 model. For each sequence we take per-step frame features $F_t \in \mathbb{R}^{d_f}$ pooled from the latent video.
2873 We also take per-step action features $A_t \in \mathbb{R}^{d_a}$ and (optionally) mouse features $M_t \in \mathbb{R}^{d_m}$ produced
2874 by the action encoder. Linear projections map these into a common space and the resulting vectors
2875 are ℓ_2 -normalized. An InfoNCE-style objective brings matching pairs (F_t, A_t) (and, when present,
2876 (F_t, M_t)) from the same timestep together. It pushes them away from other timesteps of the *same*
2877 sequence. We use frame and action masks to ignore positions without actions. A symmetric variant
2878 averages frame-to-action and action-to-frame directions.

2879 When enabled, a small future-prediction head adds a second term. Action features at time t are
2880 mapped to a prediction of the frame feature at a slightly later step $t + \Delta$. A mean-squared error
2881 encourages consistency between the prediction and the projected future frame. Together, these terms
2882 give the action encoder direct gradients tied to specific frames rather than relying solely on the pixel
2883 diffusion loss. The contrastive term enforces a one-to-one alignment between actions and the frames
2884 they co-occur with. The future head encourages actions to anticipate the visual consequences that
2885 appear shortly after they are issued.

2886 **Encoder comparison.** Table 14 contrasts the two encoders used in GUIWorld. The key theme is
2887 moving from a high-dimensional, bursty event vector toward a typed, API-like schema with explicit
2888 lag handling. This makes the action stream easier to align with the latent video timeline and to reuse
across injection modes.

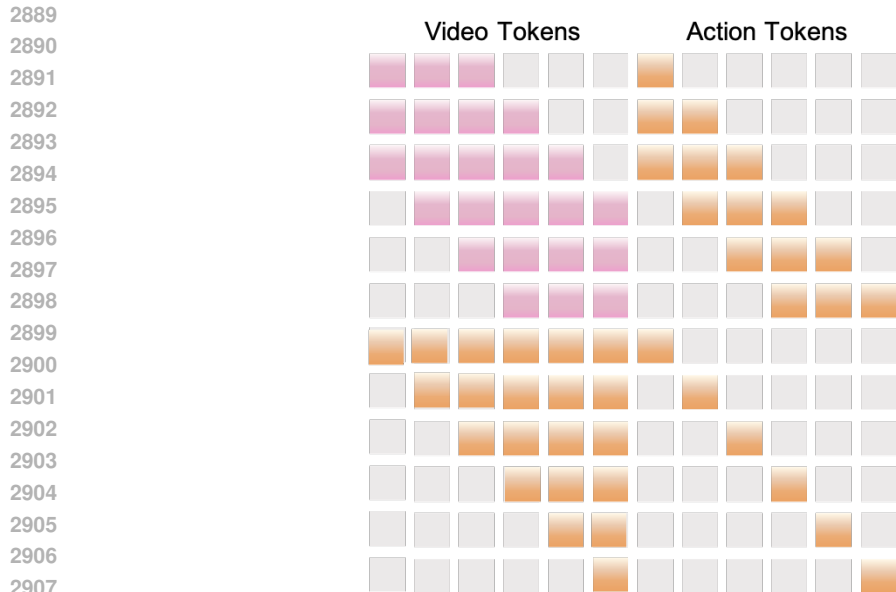


Figure 38: Contextual attention mask for GUIWorld. Video tokens (pink) and action tokens (orange) share one sequence. The mask restricts attention to a short temporal window so each frame attends only to nearby frames and temporally aligned actions.

Injection schemes. The GUI experiments explore four conditioning schemes built on this encoder: external, contextual, internal, and residual (Figure 7). Representative action-driven metrics comparing these modes are reported in Table 7. Full metric definitions and the evaluation protocol are provided in Appendix E.6.

G.6 MORE EXPERIMENTS: DO ACTION ENCODINGS MATTER?

We compare two action encodings under the same injection mode to isolate the effect of representation choice (Table 8). Under *internal* conditioning, the meta-action (API-like) encoding yields small but consistent improvements over the raw-action representation. However, these gains are modest compared to the substantially larger improvements observed when varying the action injection scheme itself (Table 7).

Table 12 contrasts how short commands and shortcuts (e.g., `ls -l`, `ctrl+v`) are represented under the two encodings; the full encoding examples are provided in Appendix E.5. The raw-action encoder treats typing as a stream of individual key events, leaving command or shortcut semantics to be inferred from the sequence. In contrast, the meta-action encoder collapses each interaction into a single typed action with associated text or a shortcut identifier. This design aims to model user actions as structured, tool-like operations rather than fragmented event streams.

In practice, this more structured abstraction does not translate into clear qualitative gains. Rendered text remains similarly smeared under both encodings, and robustness under theme changes and timing noise is largely unchanged. Task-level failures such as re-centering, re-acquisition, and multi-step interactions persist across both representations. We adopt the meta-action encoder as the default for its simplicity and semantic alignment with system-level conditioning. These results suggest that encoding granularity is secondary to alignment quality and injection strategy.