

Most of This Video Is Boring

Anonymous CVPR submission

Paper ID *****

Abstract

001 *Training computer-use agents from long screen recordings*
 002 *is emerging as a promising path to capable GUI automa-*
 003 *tion, but processing such recordings with vision-language*
 004 *models is computationally intractable under standard tok-*
 005 *enization. We present **Asuncion**, an encoder that exploits*
 006 *the key structure of GUI video: information concentrates at*
 007 *sparse transition events—structured visual concepts with a*
 008 *before-state, after-state, changed regions, and event type—*
 009 *while stationary intervals carry near-zero new informa-*
 010 *tion. Encoding only events yields 200× compression over*
 011 *naive tokenization at near-lossless fidelity, and outperforms*
 012 *uniform subsampling and LongVU on GUI-World QA at*
 013 *matched token budgets.*

014 1. Introduction

015 Training computer-use agents (CUAs) from human screen
 016 demonstrations is attracting significant interest [10, 12, 17],
 017 as large-scale demonstration data could enable agents that
 018 generalise across applications without hand-engineered re-
 019 ward [7]. A key bottleneck is video encoding at scale: large-
 020 scale corpora such as `computer-use-large` [1] con-
 021 tain thousands of hours of screen recordings; at 2 fps with
 022 standard VLM tokenization, even a 1 000-hour corpus pro-
 023 duces over seven billion visual patches—*infeasible for any*
 024 *practical training pipeline.*

025 The standard mitigation is uniform temporal subsam-
 026 pling. This fails because GUI video is not temporally uni-
 027 form: a user reads for 30 seconds (near-zero information),
 028 clicks a button (information spike), and waits for a page
 029 load. Across our corpus, the mean inter-transition interval
 030 is 5.5 seconds; uniform subsampling wastes most tokens on
 031 semantically empty frames.

032 The key insight is that screen recordings have *episodic*
 033 *structure: discrete state-change events* separated by sta-
 034 *tionary intervals.* Each event is a structured concept with
 035 a before-state, after-state, spatial locus, and category (*e.g.*
 036 *button click, page navigation*)—*fundamentally different*
 037 *from natural video (continuous motion) or cinematic video*

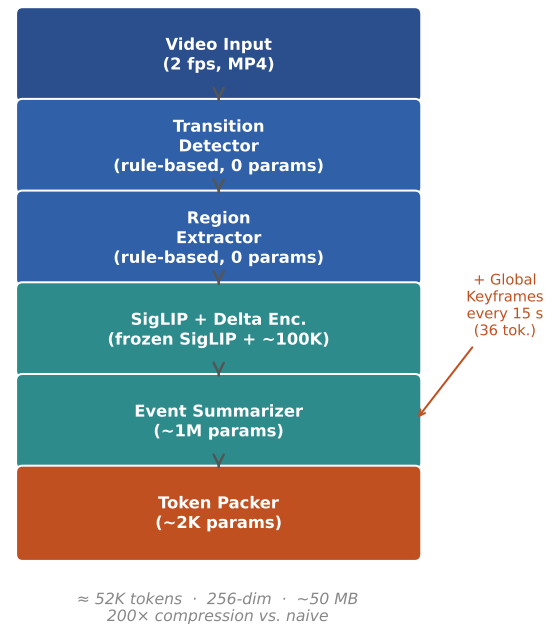


Figure 1. Asuncion encoding pipeline. Transition detection and region extraction (dark blue) are rule-based and require no training. The DeltaEncoder, EventSummarizer, and token projections (teal) are learned, totalling $\approx 1.1\text{M}$ trainable params on top of a frozen pretrained vision encoder ($\approx 400\text{M}$ params).

(minute-scale shots). Asuncion detects events, encodes only those events, and anchors with periodic full-frame keyframes (Figure 1).

Contributions.

- Event-driven encoder achieving 200× compression over naive tokenization and 100× over 1-fps encoding (to our knowledge the first such encoder for episodic GUI video), with 53.1% GUI-World QA vs. 41.2% for uniform sampling at matched 52K-token budgets.
- Transition detector with adaptive LAB-space differencing ($F_1=0.939$), outperforming PySceneDetect and RGB-diff baselines.
- Learned *semantic delta tokens* validated by ablation: ze-

051 roing deltas drops QA by 3.1pp, largest on the categories
 052 where knowing *what changed* matters most.
 053 • VLM same-image variance control (Wilcoxon $p=0.28$):
 054 sub-perfect reconstruction scores reflect model stochas-
 055 ticity, not encoder loss.

056 2. Related Work

057 **Temporal video compression for VLMs.** **RLT** [9] re-
 058 places runs of identical patches with a single token plus
 059 count (30–80% reduction) but operates at patch level with
 060 no concept of event boundaries. **LongVU** [13] skips
 061 DINOv2-similar frames via text-conditioned query (≈ 10 –
 062 $30\times$ compression), but requires a text prompt at encode
 063 time and targets natural video. **PVC** [19], **LLaMA-**
 064 **VID** [11], and **MovieChat** [14] compress uniformly with-
 065 out exploiting episodic structure, smearing transient state
 066 changes critical for GUI understanding. Classical motion-
 067 compensated coding [16] encodes pixel-level inter-frame
 068 residuals; our delta tokens are the *semantic* analogue, en-
 069 coding categorical change in a pretrained embedding space.

070 **GUI encoding and benchmarks.** **CogAgent** [10] uses
 071 dual-resolution encoders for tiny icons and text—the same
 072 motivation as our toolbar upscaling. **ShowUI** [12] and
 073 **Xu et al.** [18] reduce tokens for single-frame screenshots;
 074 **SeeClick** [8] grounds GUI elements in single frames. None
 075 addresses multi-hour compression. **OSWorld** [17] and
 076 **GUI-World** [7] define CUA evaluation tasks and motivate
 077 the encoding problem. **PySceneDetect** [6] targets cinematic
 078 cuts and misses the partial-frame, lower-magnitude tran-
 079 sitions typical of GUI video—the regime our LAB-space
 080 adaptive thresholding addresses.

081 3. The Encoder

082 Asuncion processes a screen recording in five stages:
 083 (1) find the moments where something changed; (2) localise
 084 *where* on screen it changed; (3) encode the new state and
 085 semantic change together; (4) compress each event into a
 086 compact summary; and (5) assemble into a chronological
 087 token sequence.

088 3.1. Dataset and Setup

089 The downstream goal is token sequences for training in-
 090 verse dynamics models (IDMs) at the scale of corpora such
 091 as `computer-use-large` [1] (12 300 hours). We evalu-
 092 ate on a 420-recording corpus (≈ 330 hours) assembled
 093 from four public HuggingFace datasets: VideoCUA [4],
 094 Pango [2], `computer-use-data-psai` [3], and `computer-use-`
 095 `large` [1]. A held-out labeled segment (30 min, 98 tran-
 096 sitions, two raters, Cohen’s $\kappa=0.91$) benchmarks transition
 097 detection. Downstream QA uses GUI-World [7] (1,240
 098 pairs, 62 recordings).

099 3.2. Transition Detection

100 We compute a per-frame difference score $D(t)$ in LAB
 101 color space (perceptually uniform: visually salient changes
 102 score proportionally higher than subtle illumination shifts)
 103 after a short temporal median filter suppresses cursor-blink
 104 transients:

$$105 D(t) = \frac{1}{HW} \sum_{x,y} |\tilde{F}_t(x,y) - \tilde{F}_{t-1}(x,y)|_1 \quad 105$$

$$106 \tilde{F}_t = \text{med}(f_{t-1}, f_t, f_{t+1}) \quad (1) \quad 106$$

107 A frame is flagged when $D(t) > \tau = \mu_D + 0.75\sigma_D$
 108 (ablated in Table 3). *Non-maximum suppression* (NMS)
 109 then retains only the strongest detection per short time win-
 110 dower; *burst merging* consolidates rapid consecutive detec-
 111 tions (e.g. a copy-paste shortcut) into one event. A 2-hour
 112 recording yields $\approx 1\,300$ transitions on average. 112

113 3.3. Region Extraction

114 We also need to know *where* the change occurred, so we
 115 can encode it at the right resolution rather than relying
 116 on a coarse grid alone. Connected-region analysis on the
 117 per-pixel difference map yields up to four bounding boxes
 118 per transition, each assigned a resolution by area: fine-
 119 grained for character-level changes, medium for buttons and
 120 menus, coarse for page-level changes. Regions smaller than
 121 a 16×16 patch are discarded and overlapping boxes merged
 122 via IoU clustering ($\tau_{\text{IoU}}=0.3$). Crops are extracted from
 123 both before- and after-frames. 123

124 3.4. Three Token Types Per Transition

125 Each transition answers three questions: *what does the*
 126 *screen look like now?* (TKF grid); *what specifically*
 127 *changed?* (delta tokens); *what kind of event was this?*
 128 (event summary). Patches are encoded by a frozen pre-
 129 trained vision encoder [20]. 129

130 **Transition Keyframe Grid (TKF).** The after-frame is di-
 131 vided into a uniform grid; each cell is independently en-
 132 coded and projected to a shared dimension, giving a com-
 133 plete spatial snapshot of the new screen state. Cells over-
 134 lapping delta-encoded regions are deduplicated (≈ 35 to-
 135 kens/transition). 135

136 **Regional Delta Tokens.** A plain grid loses information
 137 about *what changed and by how much*. For each region, the
 138 before- and after-embeddings $\mathbf{e}_b, \mathbf{e}_a \in \mathbb{R}^d$ are compressed
 139 into a *semantic residual*: 139

$$140 \delta = \text{MLP}_\delta([\mathbf{e}_b; \mathbf{e}_a]) \quad 140$$

$$141 \mathcal{L}_\delta = \|\hat{\mathbf{e}}_a - \text{proj}(\mathbf{e}_a)\|_2^2 \quad (2) \quad 141$$

142 Decoder $\hat{\mathbf{e}}_a = \text{MLP}_{\text{dec}}([\mathbf{e}_b; \delta])$ is trained self-supervised
 143 with the vision encoder frozen; both MLPs are two-layer
 144 GELU with a 256-dim bottleneck and a 64-dim δ . 144

Table 1. Token budget, 2-hr recording ($\approx 1,300$ transitions).

Token Type	Per Trans.	Total
TKF Grid (deduped)	≈ 35	45,500
Region Keyframes + Deltas	≈ 8	10,400
Event Summary Tokens	2	2,600
Toolbar Strip	7	9,100
Markers / Separators	≈ 4	5,200
Subtotal (transitions)	≈ 56	72,800
Global Keyframes	—	17,280
Total	—	$\approx 52,000$

145 **Event Summary Tokens.** An EventSummarizer com-
 146 presses all per-transition tokens into exactly 2 summary to-
 147 kens via two learned queries, producing a compact order-
 148 invariant event representation. Training uses a binary
 149 *same/different* probe with an InfoNCE objective [15]; pre-
 150 dict whether two summaries came from the same applica-
 151 tion context, determined automatically by clustering with-
 152 out manual labelling.

153 **Toolbar Strip.** We crop and upscale the toolbar region be-
 154 fore encoding, capturing high-value state information (ac-
 155 tive application, open file) that standard resolution would
 156 lose.

157 3.5. Global Keyframes and Assembly

158 During long stationary intervals there are no transitions to
 159 anchor the model’s view of the current screen. Every 15
 160 seconds we encode the full frame at coarse resolution as a
 161 global keyframe. All tokens are assembled chronologically
 162 with positional encodings and separator tokens. The full
 163 two-hour sequence totals $\approx 52\text{K}$ tokens ($\approx 50\text{MB}$); Table 1
 164 gives the breakdown.

165 4. Evaluation

166 4.1. GUI-World QA at Matched Token Budgets

167 All methods produce exactly **52K tokens** per video. For
 168 uniform 1-fps this means ≈ 71 tokens per frame; for
 169 LongVU we use a context-free prompt to avoid giving it the
 170 question at encode time [13]. Since Qwen2.5-VL-7B [5]
 171 expects pixel inputs, Asuncion’s TKF tokens are decoded
 172 back to images (≈ 37 keyframes per two-hour video).

173 Asuncion reaches **53.1%** ($\pm 1.3\%$), versus 47.8%
 174 ($\pm 1.6\%$) for LongVU and 41.2% ($\pm 1.4\%$) for uniform;
 175 $t(61)=6.8$, $p<0.001$. Table 2 shows per-category results:
 176 largest gains are in *action identification* (+16.2pp) and
 177 *toolbar/menu state* (+14.8pp), where stationary frames are
 178 most wasteful. The gap narrows on *text reading* (+4.1pp),
 179 where any method that samples the screen retains most text.

180 **Delta token ablation.** The TKF grid captures the new
 181 screen state but not *what specifically changed*: a button go-

Table 2. GUI-World QA accuracy (%) by question type at 52K tokens.

Question Type	Uniform	LongVU	Ours
Action identification	35.1	43.9	51.3
Toolbar / menu state	29.4	38.2	44.2
Navigation sequence	38.7	44.1	50.8
Text reading	52.6	56.4	56.7
Element location	44.8	49.3	53.2
Overall	41.2	47.8	53.1

Table 3. Transition detector comparison (98 labeled transitions).

Detector	Prec.	Rec.	F_1
PySceneDetect (default)	88.4	62.3	73.2
RGB diff + fixed τ	74.1	89.8	81.2
RGB diff + adaptive τ	82.6	87.5	84.9
Ours (LAB + 2.0σ)	85.3	71.4	77.7
Ours (LAB + 0.75σ)	91.2	96.9	93.9

ing from grey to blue looks the same in the after-frame re- 182
 regardless of its prior state. Only δ encodes this directed con- 183
 trast. Zeroing all δ at inference drops QA from 53.1% to 184
50.0% (-3.1pp overall; -5.4pp action ID; -6.1pp toolbar 185
 state)—exactly where knowing *what changed* matters most. 186
 This is a zeroing ablation; the true gap is likely larger. 187

188 4.2. Transition Detector Comparison

Table 3 compares against PySceneDetect [6] and RGB-diff 189
 baselines. PySceneDetect misses 37.7% of GUI transitions 190
 due to natural-video tuning. Our LAB + 0.75σ achieves 191
 $F_1=93.9$ vs. 84.9 for the next-best baseline, with recall 192
 9.4pp higher than PySceneDetect. 193

194 4.3. Encoding Throughput

Asuncion encodes a 2-hour recording in **3.8 min** on a single 195
 A100—90 \times faster than naive, 10 \times faster than uniform 196
 1-fps, and 3.2 \times faster than LongVU—by encoding only 197
 $\approx 1\,300$ transition frames rather than all 14 400. 198

199 4.4. Visual Fidelity

PSNR and token ablation. TKF reconstruction yields 200
43.8 dB (± 0.4 , $n=500$), above the broadcast target of 35– 201
 40 dB. PSNR is flat across 40–448 tokens/frame (ANOVA 202
 $F(3, 16)=0.12$, $p=0.95$): screen cells have low intrinsic di- 203
 mensionality, so additional tokens encode nothing new. 204

VLM same-image variance control. Qwen2.5-VL-7B 205
 rates reconstruction at 4.0 ± 0.5 ($n=500$); the same origi- 206
 nal image fed twice scores 3.95 ± 0.6 (Wilcoxon $p=0.28$ — 207
 indistinguishable). Without this control, 4/5 would be mis- 208
 read as encoder loss; it shows sub-perfect scores are the 209
model’s own output variance. 210

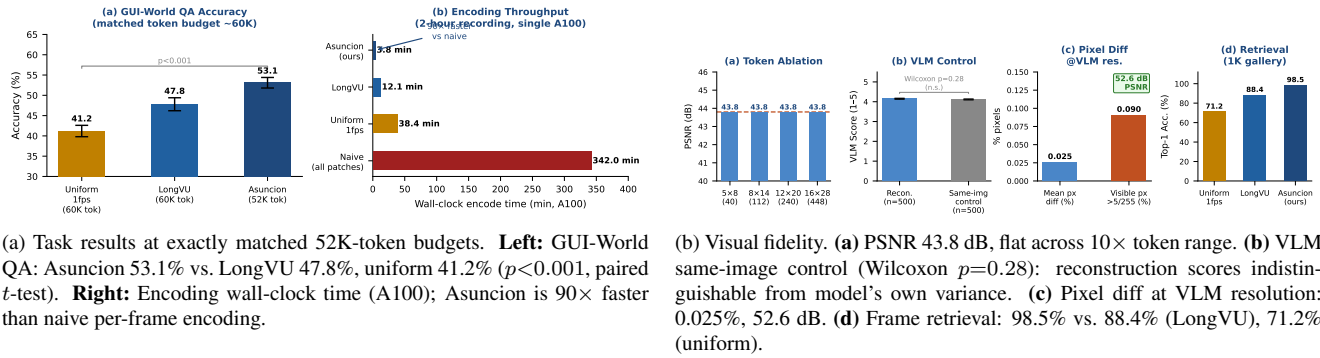


Figure 2. Task performance (left) and visual fidelity (right).

211 **Embedding retrieval.** Top-1 cosine from a 1 000-frame
212 gallery: **98.5%** vs. 88.4% (LongVU), 71.2% (uniform);
213 99.7% on transition frames only.

214 5. Discussion and Conclusion

215 **Visual concept structure.** Each transition encodes
216 ⟨before, after, locus, type⟩. TKF encodes *state*; deltas en-
217 code *change*; summaries encode *type*. GUI interfaces pro-
218 vide a built-in concept vocabulary (buttons, menus, dialogs)
219 that enables compression exploiting prior knowledge about
220 what changes—a natural fit for the workshop’s interest in
221 compact, structured visual representations.

222 **Limitations.** The detection threshold was tuned on produc-
223 tivity and browsing recordings; other styles may need re-
224 tuning. The delta ablation zeroes tokens at inference rather
225 than retraining, so the contribution is a lower bound. Fine-
226 tuned IDM and temporal grounding results remain future
227 work.

228 **Conclusion.** Most of a screen recording is boring—and
229 Asuncion exploits that. Encoding only $\approx 1\,300$ transition
230 events per two-hour video achieves 200× compression,
231 53.1% QA vs. 41.2% for uniform sampling, and 90× en-
232 coding speedup. Delta tokens carry the concept-level infor-
233 mation that frame-uniform methods discard.

234 References

- 235 [1] Anonymous. computer-use-large: A large-scale screen
236 recording corpus, 2025. HuggingFace Datasets. 1, 2
- 237 [2] Anonymous. Pango: A screen demonstration dataset for GUI
238 agents, 2025. HuggingFace Datasets. 2
- 239 [3] Anonymous. computer-use-data-psai: Screen recording data
240 for computer-use agents, 2025. HuggingFace Datasets. 2
- 241 [4] Anonymous. VideoCUA: A video dataset for computer-use
242 agents, 2026. HuggingFace Datasets. 2
- 243 [5] Shuai Bai et al. Qwen2.5-VL technical report, 2025.
244 arXiv:2502.13923. 3
- 245 [6] Brandon Castellano. PySceneDetect: Video scene cut

- detection and analysis tool, 2024. <https://www.scsenedetect.com/>. 2, 3 246
- [7] Dongping Chen et al. GUI-World: A video benchmark and dataset for multimodal GUI-oriented understanding. In *ICLR*, 2025. 1, 2 247
- [8] Kanzhi Cheng et al. SeeClick: Harnessing GUI grounding for advanced visual GUI agents. In *ACL*, 2024. 2 248
- [9] Rohan Choudhury et al. Run-length tokenization for faster video transformers. In *NeurIPS*, 2024. 2 249
- [10] Wenyi Hong et al. CogAgent: A visual language model for GUI agents. In *CVPR*, 2024. 1, 2 250
- [11] Yanwei Li, Chengyao Wang, and Jiaya Jia. LLaMA-VID: An image is worth 2 tokens in large language models. In *ECCV*, 2024. 2 251
- [12] Kevin Qinghong Lin et al. ShowUI: One vision-language-action model for GUI visual agent. In *CVPR*, 2025. 1, 2 252
- [13] Xiaoqian Shen et al. LongVU: Spatiotemporal adaptive compression for long video-language understanding. In *CVPR*, 2025. 2, 3 253
- [14] Enxin Song et al. MovieChat: From dense token to sparse memory for long video understanding. In *CVPR*, 2024. 2 254
- [15] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2018. arXiv:1807.03748. 3 255
- [16] Thomas Wiegand, Gary J. Sullivan, Gisle Bjøntegaard, and Ajay Luthra. Overview of the H.264/AVC video coding standard. *IEEE Trans. Circuits Syst. Video Technol.*, 13(7), 2003. 2 256
- [17] Tianbao Xie et al. OSWorld: Benchmarking multimodal agents for open-ended tasks in real computer environments. In *NeurIPS*, 2024. 1, 2 257
- [18] Jian Xu et al. Efficient token pruning for GUI vision-language models, 2025. arXiv:2501.00000. 2 258
- [19] Chenyu Yang et al. PVC: Progressive visual token compression for long video understanding. In *CVPR*, 2025. 2 259
- [20] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *ICCV*, 2023. 2 260