

# ObjectAlign: Neuro-Symbolic Object Consistency Verification and Correction

Anonymous CVPRW submission

Paper ID 17

## Abstract

001 *Video editing and synthesis often introduce object inconsis-*  
002 *tencies, such as frame flicker and identity drift that degrade*  
003 *perceptual quality. To address these issues, we introduce*  
004 *ObjectAlign, a novel framework that seamlessly blends per-*  
005 *ceptual metrics with symbolic reasoning to detect, verify, and*  
006 *correct object-level and temporal inconsistencies in edited*  
007 *video sequences. The novel contributions of ObjectAlign*  
008 *are as follows: First, we propose learnable thresholds for*  
009 *metrics characterizing **object consistency** (i.e. CLIP-based*  
010 *semantic similarity, LPIPS perceptual distance, histogram*  
011 *correlation, and SAM-derived object-mask IoU). Second,*  
012 *we introduce a neuro-symbolic verifier that combines two*  
013 *components: (a) a formal, SMT-based check that operates*  
014 *on masked object embeddings to provably guarantee that ob-*  
015 *ject identity does not drift, and (b) a temporal fidelity check*  
016 *that uses a probabilistic model checker to verify the video’s*  
017 *formal representation against a temporal logic specification*  
018 *( $\Phi$ ). A frame transition is subsequently deemed “consistent”*  
019 *based on a single logical assertion that requires satisfying*  
020 *both the learned metric thresholds and this unified neuro-*  
021 *symbolic constraint, ensuring both low-level stability and*  
022 *high-level temporal correctness. Finally, for each contigu-*  
023 *ous block of flagged frames, we propose a neural network*  
024 *based interpolation for adaptive frame repair, dynamically*  
025 *choosing the interpolation depth based on the number of*  
026 *frames to be corrected. By decoupling consistency verifica-*  
027 *tion from the heavy generative diffusion process, ObjectAlign*  
028 *operates with **minimal computational overhead**, making it*  
029 *highly suitable for resource-constrained edge deployments.*  
030 *Our results show up to 1.4 point improvement in CLIP Score*  
031 *and up to 6.1 point improvement in warp error compared to*  
032 *SOTA baselines on the DAVIS and Pexels video datasets.*

## 033 1. Introduction

034 Recent advances in artificial intelligence have significantly  
035 enhanced the quality, realism, and efficiency of synthetic im-  
036 age and video generation models [14, 18, 35, 44]. These im-  
037 provements have broadened applications in content creation,

real-time video editing, and interactive media [10, 27, 33].  
Despite these strides, a critical yet often overlooked chal-  
lenge persists, namely *maintaining consistent object repre-*  
*sentation* across different video frames. This is important  
since subtle inconsistencies, including semantic drift, visual  
flickering, or transient artifacts, frequently arise during video  
synthesis and editing, diminishing the visual coherence and  
perceptual realism [14].

Current diffusion-based editing methods [14, 32, 37, 42,  
53], predominantly use extended attention mechanisms to  
propagate information across frames to maintain tempo-  
ral coherence. However, extending attention across mul-  
tiple frames significantly increases the computational cost  
and memory requirements, often becoming prohibitively  
expensive [14, 37] and unsuitable for deployment on memory-  
constrained edge devices. Moreover, these approaches do  
not provide formal guarantees for consistency, leaving room  
for errors that degrade the video quality.

To overcome these limitations, there is an emerging need  
for robust verification methods capable of *provably* ensuring  
consistency between frames. Unlike methods relying solely  
on perceptual metrics which may still miss subtle inconsis-  
tencies or offer no formal assurances, a provable guarantee,  
such as that provided by a Satisfiability Modulo Theories  
(SMT) solver [11], can offer a mathematically-grounded as-  
sertion that specified consistency constraints (e.g., bounds on  
semantic feature drift) are met. This is crucial for detecting  
errors that evade heuristic checks.

In this paper, we propose **ObjectAlign**, a neuro-symbolic  
framework that can rigorously verify and adaptively repair  
object-level and temporal inconsistencies in edited video  
sequences. Our approach bridges perceptual metrics with  
symbolic verification techniques, ensuring both practical  
performance and formal consistency guarantees. To this end,  
we introduce three key contributions:

- First, we propose a new methodology that integrates multi-  
perceptual and semantic metrics, including CLIP-based  
semantic similarity, LPIPS perceptual distance, color his-  
togram correlation, and segmentation-mask IoU, into a  
unified, *learnable threshold-based classifier* for identify-  
ing object-level inconsistencies. By learning thresholds

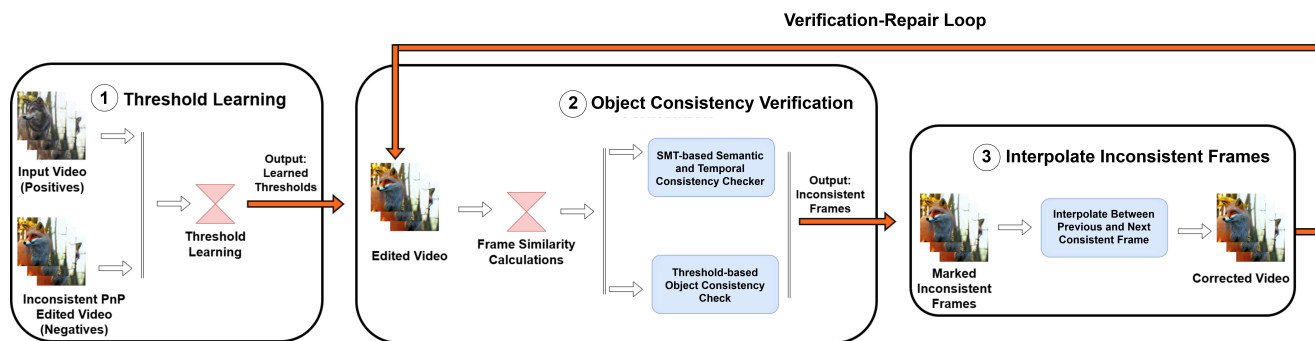


Figure 1. **Overview of ObjectAlign.** ① We first learn per-metric consistency thresholds from “positive” original video clips and “negative” inconsistently edited clips. ② Next, for each consecutive frame pair in a newly edited video, we compute semantic and perceptual similarities and apply both the learned threshold checks and an SMT-based object consistency check on the embeddings to flag inconsistent transitions. ③ Finally, each contiguous block of flagged frames is repaired by adaptively interpolating between the nearest preceding and succeeding consistent keyframes, with the interpolation depth chosen according to the segment length. The corrected frames can then be re-verified in a closed loop until no inconsistencies remain.

- 079 directly from data, our approach offers both flexibility and  
 080 interpretability in inconsistency detection.  
 081 • Second, we introduce a formal verification method to *prov-*  
 082 *ably verify semantic and temporal consistency.* Specif-  
 083 ically, we embed object features as constraints within a  
 084 symbolic reasoning framework, enforcing per-dimension  
 085 semantic bounds on masked CLIP embeddings. We also  
 086 ensure temporal fidelity verification through a probabilistic  
 087 model checker to verify the video satisfies a given tempo-  
 088 ral logic specification ( $\Phi$ ). This ensures a mathematically  
 089 grounded guarantee of semantic and temporal consistency  
 090 within defined thresholds.  
 091 • Finally, we develop an *adaptive interpolation strategy* for  
 092 correcting flagged inconsistencies. Our repair mechanism  
 093 dynamically adjusts the interpolation depth based on the  
 094 number of contiguous inconsistent frames identified, re-  
 095 constructing corrupted frames from adjacent consistent  
 096 keyframes, thus preserving a smooth temporal coherence.

097 Indeed, as shown in Figure 1, ObjectAlign effectively inte-  
 098 grates learnable perceptual metrics, formal semantic ver-  
 099 ification, and adaptive interpolation-based correction into a  
 100 unified end-to-end pipeline. Our evaluation demonstrates  
 101 that ObjectAlign reduces perceptual flickering and semantic  
 102 drift, decreasing the warp error [30] from 107.4 to 101.3  
 103 compared to Plug and Play Diffusion (PnP) [53] on clips  
 104 from the DAVIS [45] and Pexels [1] video datasets.

105 The remainder of this paper is structured as follows: Sec-  
 106 tion 2 discusses related work in video synthesis and formal  
 107 verification techniques. Section 3 provides necessary back-  
 108 ground on diffusion models, perceptual metrics, and SMT  
 109 solvers. Section 4 describes the ObjectAlign methodology  
 110 and technical innovations in detail. Experimental results  
 111 and ablations are presented in Section 5. Finally, Section 6  
 112 summarizes our main contributions.

## 2. Related work 113

### 2.1. Video Editing and Object Consistency 114

Recent works have explored training-free frameworks for  
 115 improving or stylizing text-to-video generation by leverag-  
 116 ing pre-trained text-to-image (T2I) models to edit video  
 117 frames [21, 23, 54, 58]. Approaches such as SDEdit [42],  
 118 InstructPix2Pix [4], and ControlNet [62] provide general-  
 119 purpose image editing capabilities that have been adapted  
 120 for video by applying them frame-by-frame or with addi-  
 121 tional guidance. Several methods enhance video genera-  
 122 tion through refined text prompts [24, 36], or by combining  
 123 text and image modalities for editing [61]. Plug-and-Play  
 124 Diffusion [53] and Free2Guide [24] further enable flexible,  
 125 training-free editing. Dreamix [43] and Tune-A-Video [55]  
 126 demonstrate the use of video diffusion models and spatio-  
 127 temporal tuning for improved consistency and style transfer.  
 128 Real-time editing approaches such as StreamDiffusion and  
 129 StreamV2V [27, 33] enable efficient video editing. 130

A key limitation of these approaches lies in their dif-  
 131 ficulty to maintain temporal coherence and object consis-  
 132 tency across frames. Methods like TokenFlow [14], Rerender  
 133 [58], and VideoP2P [34] address this by identifying key  
 134 frames [5] and propagating features across frames. Other  
 135 approaches, such as Ground-A-Video [21], FateZero [46],  
 136 and Ada-VE [37], rectify cross-frame attention or integrate  
 137 motion cues to improve consistency. Despite these advances,  
 138 cross-frame attention remains computationally expensive  
 139 and does not provide formal guarantees of consistency [14]. 140

### 2.2. Neuro-Symbolic Verification 141

Neuro-symbolic methods aim at integrating the advancement  
 142 of neural networks with the rigor of symbolic reasoning  
 143 [9, 16]. Neuro-symbolic methods use symbolic reasoning  
 144 to provide formal guarantees in various domains. Specifically,  
 145 in image and video synthesis, formal verification approaches  
 146

147 such as SMT [11] and temporal logic [7, 50] can rigorously  
148 validate the consistency and semantic correctness of the  
149 generated content.

150 Recently, neuro-symbolic verification has been explored for video searching, editing, and evaluation tasks  
151 [7, 50]. Video classification employs graph-based relational modeling [13, 52], while event detection lever-  
152 ages spatiotemporal pattern recognition in video streams [31, 40, 57]. Neuro-symbolic frameworks enhance video  
153 question-answering [6, 59], with applications extending to robotic action planning [15, 28, 51] and safety verification  
154 in autonomous driving systems [22, 41]. These methods either construct graph structures [39, 56, 60], use latent-space  
155 representations as symbolic representations [3, 29, 49], or use formal language methods [2] to design specifications.  
156

157 In contrast with this prior work, our ObjectAlign uniquely  
158 combines learnable perceptual metrics with symbolic constraints using the SMT solving, thus providing formal guar-  
159 antees for object consistency in video editing. Additionally, we complement our verification with adaptive interpolation  
160 to repair inconsistencies dynamically. ObjectAlign is the first work to explore object consistency correction by post-  
161 processing inconsistent frames identified by learnable perceptual metrics or formal verification.  
162

## 171 3. Preliminaries

### 172 3.1. Latent Diffusion Models

173 Diffusion models [12, 17] are generative models comprised of two main stochastic phases: (a) a *forward process* that  
174 progressively adds noise to data, and (b) a *reverse process* that learns to remove this noise to generate data.  
175

176 The **forward process** is typically formulated as a fixed Markov chain that gradually introduces Gaussian noise to  
177 an initial data sample  $\mathbf{z}_0$  over  $T$  discrete time steps. If  $\mathbf{z}_0$  is a sample from the true data distribution  $p_{\text{data}}$  (e.g., a  
178 clean image), this process yields a sequence of increasingly noisy samples  $\mathbf{z}_1, \dots, \mathbf{z}_T$ . The final sample  $\mathbf{z}_T$  is ideally distributed  
179 according to a normal distribution,  $\mathcal{N}(0, \mathbf{I})$ , where  $\mathbf{I}$  is the identity covariance matrix. The transition at each step  
180  $t$  is defined as:  
181

$$182 q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{\alpha_t} \mathbf{z}_{t-1}, (1 - \alpha_t) \mathbf{I}), \quad (1)$$

183 where  $\mathbf{z}_{t-1}$  is the data sample at the previous time step,  $\mathbf{z}_t$   
184 is the sample at the current time step, and  $\alpha_t$  is a parameter derived from a predefined noise schedule (e.g.,  $\alpha_t = 1 -$   
185  $\beta_t \in (0, 1)$ , where  $\beta_t \in (0, 1)$  are small positive constants representing variance schedules).  
186

187 The **reverse process** aims to reverse this noising procedure. It starts with a sample  $\mathbf{z}_T \sim \mathcal{N}(0, \mathbf{I})$  and iteratively  
188 denoises it to produce a sample  $\mathbf{z}_0$  that resembles data from the true distribution  $p_{\text{data}}$ . This process is also a Markov  
189

190 chain, parameterized by a neural network with parameters  $\theta$ . The model is trained to predict the conditional probability  
191 distribution  $p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t)$  for each step  $t$ :  
192

$$193 p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t) = \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{z}_t, t), \sigma_t^2 \mathbf{I}), \quad (2)$$

194 where  $\boldsymbol{\mu}_\theta(\mathbf{z}_t, t)$  is the mean of the Gaussian distribution for  $\mathbf{z}_{t-1}$ , predicted by the neural network conditioned on the  
195 noisy sample  $\mathbf{z}_t$  and the time step  $t$ . The term  $\sigma_t^2$  represents the variance at time step  $t$ , which is often predefined or  
196 learned as part of the noise schedule. The neural network is trained to make  $p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t)$  accurately approximate the  
197 true posterior  $q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{z}_0)$ .  
198

### 199 3.2. Metrics for Object Consistency Verification

200 To assess object consistency between video frames, ObjectAlign employs a combination of perceptual metrics and  
201 formal verification techniques. Let  $f_i, f_j \in \mathcal{I}$  be two video frames, where  $\mathcal{I}$  denotes the entire space of all video frames.  
202

203 **Perceptual Consistency Metrics.** We utilize four established metrics to capture different aspects of visual and semantic similarity:  
204

- 205 • **Learned Perceptual Image Patch Similarity (LPIPS):** This metric quantifies low-level perceptual similarity. The  
206 LPIPS distance between any two frames  $f_i, f_j$  is given by:  
207

$$208 \text{LPIPS}(f_i, f_j) = \|\phi(f_i) - \phi(f_j)\|_2, \quad (3)$$

209 where  $\phi: \mathcal{I} \rightarrow \mathbb{R}^d$  is a deep feature extractor [63] and  $\|x\|_2$  denotes the Euclidean norm. Smaller LPIPS values  
210 indicate that the frames are more similar at a patch-perceptual level.  
211

- 212 • **CLIP-based Semantic Similarity:** To measure high-level semantic alignment, we use the cosine similarity between  
213 image embeddings from a Contrastive Language-Image Pre-training (CLIP) model [47]:  
214

$$215 \text{Sim}_{\text{CLIP}}(f_i, f_j) = \frac{\langle e(f_i), e(f_j) \rangle}{\|e(f_i)\| \cdot \|e(f_j)\|}, \quad (4)$$

216 where  $e: \mathcal{I} \rightarrow \mathbb{R}^k$  is the CLIP image encoder,  $\langle x, y \rangle$  denote the standard dot-product and  $\|x\|_2$  denotes the Euclidean norm. Values closer to 1 signify stronger semantic  
217 correspondence between the frames.  
218

- 219 • **Histogram Correlation:** To check for significant color shifts between frames, we compute the correlation between their color histograms. Let  $h(f) \in \mathbb{R}^c$  be the flattened and normalized color histogram vector for frame  
220  $f$ , where  $c$  represents the dimensionality. The histogram correlation is:  
221

$$222 \text{Sim}_{\text{Hist}}(f_i, f_j) = \frac{h(f_i)^\top h(f_j)}{\|h(f_i)\|_2 \cdot \|h(f_j)\|_2}, \quad (5)$$

223 where  $(\cdot)^\top$  denotes transpose. Values closer to 1 indicate higher similarity in the overall color distributions.  
224

241 • **Mask IoU:** For object-level geometric consistency, we  
 242 compute the IoU of foreground object masks. Let  
 243  $M: \mathcal{I} \rightarrow \{0, 1\}^{H \times W}$  be the binary foreground mask ob-  
 244 tained for a frame (e.g., via the Segment Anything Model  
 245 (SAM) [26]), where  $H$  and  $W$  are the frame height and  
 246 width. The IoU is:

$$247 \quad \text{IoU}(M(f_i), M(f_j)) = \frac{|M(f_i) \cap M(f_j)|}{|M(f_i) \cup M(f_j)|}. \quad (6)$$

248 This value, ranging from 0 to 1, quantifies the spatial  
 249 overlap of the primary objects.

250 These perceptual metrics provide complementary empiri-  
 251 cal checks on object consistency, covering low-level appear-  
 252 ance, high-level semantic content, color distribution, and  
 253 object geometry, respectively. However, taken alone, they  
 254 cannot inherently provide formal guarantees of coherence.

255 **Formal Verification with SMT Solvers.** To address the  
 256 limitations of purely metric-based approaches and introduce  
 257 rigorous consistency checks, ObjectAlign provides formal  
 258 verification using SMT solvers. SMT solvers can determine  
 259 the satisfiability of logical formulas with respect to back-  
 260 ground theories, enabling us to enforce provable bounds on  
 261 specific features. In our context, we use an SMT solver  
 262 to enforce *semantic stability* by asserting bounds on the  
 263 object drift. We use object masks  $M_i$  to compute sepa-  
 264 rate embeddings for the foreground object  $e(f_i, M_i)$  and the  
 265 background  $e(f_i, \neg M_i)$ . We then use an SMT solver to for-  
 266 mally verify a *conjunctive* formula that ensures both **object**  
 267 **identity stability** and **background stability**:

$$268 \quad (\forall j \mid e_j(f_i, M_i) - e_j(f_{i+1}, M_{i+1}) \mid \leq \epsilon_s) \wedge \\ (\forall j \mid e_j(f_i, \neg M_i) - e_j(f_{i+1}, \neg M_{i+1}) \mid \leq \epsilon_{bg}) \quad (7)$$

269 where  $\epsilon_s$  and  $\epsilon_{bg}$  are semantic drift tolerances. An SMT  
 270 solver checks if this set of constraints is satisfiable; if it is,  
 271 then we have a formal guarantee that no individual semantic  
 272 feature dimension has drifted beyond the specified tolerances  
 273  $\epsilon_s$  and  $\epsilon_{bg}$ . ObjectAlign leverages this neuro-symbolic veri-  
 274 fication to complement the aforementioned learned perceptual  
 275 metrics, providing a more robust and reliable consistency  
 276 assessment than using the perceptual metrics alone.

## 277 4. Proposed Methodology

278 ObjectAlign consists of three stages executed in a closed  
 279 verification–repair loop (Fig. 1): ① *metric-based scoring*  
 280 *with learned thresholds*, ② *neuro-symbolic consistency veri-*  
 281 *fication*, and ③ *adaptive frame repair via neural interpola-*  
 282 *tion*. The loop repeats until every neighbouring frame pair  
 283 satisfies *all* consistency constraints.

## 4.1. Inconsistency identification (Step ① in Fig. 1)

### 4.1.1. Metric Based Consistency Scoring

284 **Feature vector.** For two consecutive frames  $f_i, f_{i+1}$ , we  
 285 extract (a) cosine similarity of CLIP embeddings  $S_{\text{cos}}$ , (b)  
 286 color–histogram correlation  $S_{\text{hist}}$ , (c) mask–IoU  $S_{\text{iou}}$ , and  
 287 (d) perceptual distance  $D_{\text{lpiips}}$ . We invert LPIPS so that  
 288 larger values denote higher consistency, i.e.  $\tilde{S}_{\text{lpiips}} = -D_{\text{lpiips}}$ .  
 289 Hence the feature vector is  $\mathbf{s}_i = [S_{\text{cos}}, S_{\text{hist}}, S_{\text{iou}}, \tilde{S}_{\text{lpiips}}]^\top$ .  
 290 These specific metrics are chosen for their complementary  
 291 strengths in assessing frame-to-frame object consistency:  
 292 CLIP [47] similarity captures high-level semantic content  
 293 alignment, LPIPS [63] evaluates low-level perceptual appear-  
 294 ance, color histogram correlation checks for drastic color  
 295 shifts, and mask IoU quantifies object-level geometric over-  
 296 lap and spatial stability, thereby providing a comprehensive  
 297 empirical check as noted in Section 3.

298 **Learnable thresholds.** We treat each dimension of the

299 feature vector  $\mathbf{s}_i$  independently and learn a threshold vector  
 300  $\boldsymbol{\tau} = [\tau_{\text{cos}}, \tau_{\text{hist}}, \tau_{\text{iou}}, \tau_{\text{lpiips}}]^\top$  from a small *positive* set  $\mathcal{P}$  (ad-  
 301 jacent frames from the unedited video) and a *negative* set  $\mathcal{N}$   
 302 (pairs of original vs. edited frames, considered inconsistent).  
 303 For each frame pair  $i$ , we compute the element-wise differ-  
 304 ence vector  $\Delta_i$  between its feature vector  $\mathbf{s}_i$  and the learned  
 305 threshold vector  $\boldsymbol{\tau}$ :  $\Delta_i = \mathbf{s}_i - \boldsymbol{\tau}$ . The probability that a pair  
 306 is consistent for a single threshold  $k$  ( $P_k(i)$ ) is modeled by  
 307 the sigmoid function:  
 308

$$309 \quad P_k(i) = \sigma(\lambda \Delta_k), \quad \sigma(z) = \frac{1}{1 + e^{-z}}, \quad (8)$$

310 where  $\lambda$  is a sharpness constant. The four thresholds in  
 311  $\boldsymbol{\tau}$  are simultaneously optimized by minimizing the binary  
 312 cross-entropy loss:  
 313

$$314 \quad \mathcal{L}_{\text{BCE}} = - \frac{1}{|\mathcal{P}| + |\mathcal{N}|} \sum_{i \in \mathcal{P} \cup \mathcal{N}} \left[ y_i \log(P_k(i)) \right. \\ \left. + (1 - y_i) \log(1 - P_k(i)) \right] \quad (9)$$

315 where  $y_i = 1$  for  $i \in \mathcal{P}$  and  $y_i = 0$  for  $i \in \mathcal{N}$ . Optimiza-  
 316 tion is performed using Adam [25].

### 4.1.2. Neuro-Symbolic Verification (Step ② in Fig. 1)

317 While the metric classifier is effective in practice for captur-  
 318 ing perceptual inconsistencies, it offers no *formal* guarantee  
 319 against all forms of object drift, particularly subtle semantic  
 320 shifts that may fall within learned perceptual thresholds but  
 321 still represent a logical inconsistency. The SMT-based veri-  
 322 fication step (see ② in Fig 1) addresses this by combining  
 323 low-level feature stability with high-level temporal fidelity.  
 324

325 Given that the scalar perceptual metrics (e.g.,  $S_{\text{hist}}$ ) are  
 326 directly evaluated against their learned thresholds (Eq. (12)),  
 327 SMT verification is reserved for the high-dimensional CLIP

328 embeddings to enforce semantic stability. We therefore im-  
 329 pose an SMT constraint on the **masked CLIP embeddings**  
 330 (introduced in Sec. 3.2). Specifically, we verify the stability  
 331 of both the foreground object  $e(f, M)$  and the background  
 332  $e(f, \neg M)$  independently, defining this semantic stability con-  
 333 straint as  $C_{\text{neuro}}$ :

$$334 \quad C_{\text{neuro}} \equiv (\forall j \mid e_j(f_i, M_i) - e_j(f_{i+1}, M_{i+1}) \mid \leq \epsilon_s) \wedge \\ (\forall j \mid e_j(f_i, \neg M_i) - e_j(f_{i+1}, \neg M_{i+1}) \mid \leq \epsilon_{bg}) \quad (10)$$

335 We complement this stability check with a high-level  
 336 temporal fidelity metric [8, 50]. This component calculates  
 337 a *satisfaction probability* by verifying the video’s formal  
 338 representation (automaton  $\mathcal{A}_\nu$ ) against the text prompt’s  
 339 temporal logic specification ( $\Phi$ ) using a probabilistic model  
 340 checker function,  $\Psi$ . A video is considered formally verified  
 341 *only if* it satisfies both the low-level stability constraints  
 342 ( $C_{\text{neuro}}$ ) and the high-level temporal requirements. We define  
 343 this unified neuro-symbolic constraint,  $\mathcal{P}_{\text{formal}}$ , as the logical  
 344 conjunction of these two conditions:

$$345 \quad \mathcal{P}_{\text{formal}} \equiv C_{\text{neuro}} \wedge (\Psi(\mathcal{A}_\nu, \Phi) \geq \tau) \quad (11)$$

346 Here,  $\mathcal{P}_{\text{formal}}$  is satisfied if and only if the SMT solver finds  
 347 the frame-to-frame drift constraints  $C_{\text{neuro}}$  (the first conjunct)  
 348 satisfiable for all frames, *and* the probabilistic model checker  
 349 finds that the temporal fidelity  $\Psi(\mathcal{A}_\nu, \Phi)$  meets or exceeds a  
 350 specified probability threshold  $\tau$  (the second conjunct).

### 351 4.1.3. Joint Consistency Criterion

352 A transition is declared *consistent* ( $C(i) = 1$ ) if and only  
 353 if (iff) *all* thresholds are satisfied simultaneously *and* the  
 354 formal constraints are met:

$$355 \quad C(i) = (S_{\text{cos}} \geq \tau_{\text{cos}}) \wedge (S_{\text{hist}} \geq \tau_{\text{hist}}) \wedge \\ (S_{\text{iou}} \geq \tau_{\text{iou}}) \wedge (D_{\text{lipips}} \leq \tau_{\text{lipips}}) \wedge \\ (P_{\text{formal}}(i) = 1) \quad (12)$$

356 All indices with  $C(i) = 0$  form the inconsistent set  $\mathcal{I}$ .

357 **Formal Consistency Guarantees.** The joint consistency  
 358 criterion  $C(i)$  defined in Eq. (12) combines the learned  
 359 threshold checks with a formal SMT constraint. Our joint  
 360 consistency criterion assures that a frame-pair declared con-  
 361 sistent by our pipeline ( $C(i) = 1$ ) exhibits a bounded drift  
 362 according to every metric and the formal semantic check  
 363 included in our criteria. These formal bounds underpin Ob-  
 364 jectAlign’s robustness in improving video consistency.

## 365 4.2. Adaptive Frame Repair (Step ③ in Fig. 1)

366 The adaptive frame repair stage (see ③ in Fig 1) funda-  
 367 mentally relies on the presence or eventual emergence of  
 368 consistent anchor frames surrounding any block of identified

---

### Algorithm 1 OBJECTALIGN verification–repair loop

---

```

1: Input: edited video  $V = \{f_0, \dots, f_{T-1}\}$   $\triangleright T$  is the total
   number of frames in video  $V$ .
2: Learn thresholds  $\tau$  on positive and negative set  $\triangleright$  Sec. 4.1
3: repeat
4:    $\mathcal{I} \leftarrow \emptyset$ 
5:   for  $i = 0$  to  $T - 2$  do  $\triangleright$  Iterate over all  $T - 1$  frame
   transitions  $(f_i, f_{i+1})$ 
6:     compute  $P_{\text{metric}}(i), P_{\text{formal}}(i)$ 
7:     if  $C(i)$  (Eq. 12) is false then
8:        $\mathcal{I} \leftarrow \mathcal{I} \cup \{i\}$ 
9:     for each contiguous run  $[i_s, i_e] \subseteq \mathcal{I}$  do  $\triangleright i_s, i_e$ : start/end
   indices of a run of inconsistent transitions
10:       $k \leftarrow i_e - i_s + 1$   $\triangleright$  Number of frames  $f_{i_s}, \dots, f_{i_e}$  to
   repair
11:      replace  $f_{i_s}, \dots, f_{i_e}$  with corrected frames  $\triangleright$  Sec. 4.2
12: until  $\mathcal{I} = \emptyset$ 
13: return verified & corrected video  $V$ 

```

---

inconsistencies. This principle can be conceptualized using  
 Linear Temporal Logic (LTL) [38].

Let  $AP_{IB}$  be an atomic proposition that is true when  
 a contiguous block of frames is currently identified as an  
 ‘*InconsistentBlock*’ requiring repair. Let  $AP_{CAB}$  be true if  
 a ‘*ConsistentAnchorBefore*’ (i.e., a suitable frame  $f_{i_s-1}$ )  
 exists or is established, and  $AP_{CAA}$  be true if a ‘*Consis-*  
*tentAnchorAfter*’ (i.e., a suitable frame  $f_{i_e+1}$ ) exists or is  
 established. The iterative verification–repair loop of Objec-  
 tAlign (Algorithm 1) operates under the premise that the  
 video sequence will eventually satisfy the property:

$$380 \quad \square(AP_{IB} \implies (\diamond AP_{CAB} \wedge \diamond AP_{CAA})) \quad (13)$$

381 This LTL formula asserts that it is always ( $\square$ ) the case  
 382 that if an inconsistent block requiring repair ( $AP_{IB}$ ) exists,  
 383 then eventually ( $\diamond$ ) a consistent anchor frame will be found  
 384 or established before it ( $AP_{CAB}$ ), and eventually ( $\diamond$ ) a con-  
 385 sistent anchor frame will be found or established after it  
 386 ( $AP_{CAA}$ ), thus enabling interpolation. Our framework aims  
 387 to progressively achieve this state, allowing for repair even  
 388 when initial edits contain extended inconsistent segments.

389 Given a contiguous sequence of frames marked as incon-  
 390 sistent, we apply adaptive neural network based interpolation  
 391 using RIFE [19]. Let  $[i_s, i_e] \subseteq \mathcal{I}$  represent a maximal run  
 392 of  $k = i_e - i_s + 1$  inconsistent frames. To reconstruct  
 393 these frames, we first identify the closest consistent frames  
 394 immediately preceding and following this sequence:  $f_{i_s-1}$   
 395 and  $f_{i_e+1}$  (whose existence is anticipated by the property  
 396 in Eq. (13)). We then dynamically select the interpolation  
 397 depth ( $\gamma$ ) as a function of the number of frames needing  
 398 repair ( $k$ ), defined by:  $\gamma = \lceil \log_2(k + 1) \rceil$ . This adaptive in-  
 399 terpolation depth ensures that longer runs of inconsistencies  
 400 are addressed with deeper interpolation, generating sufficient  
 401 intermediate frames to preserve smooth and coherent motion.

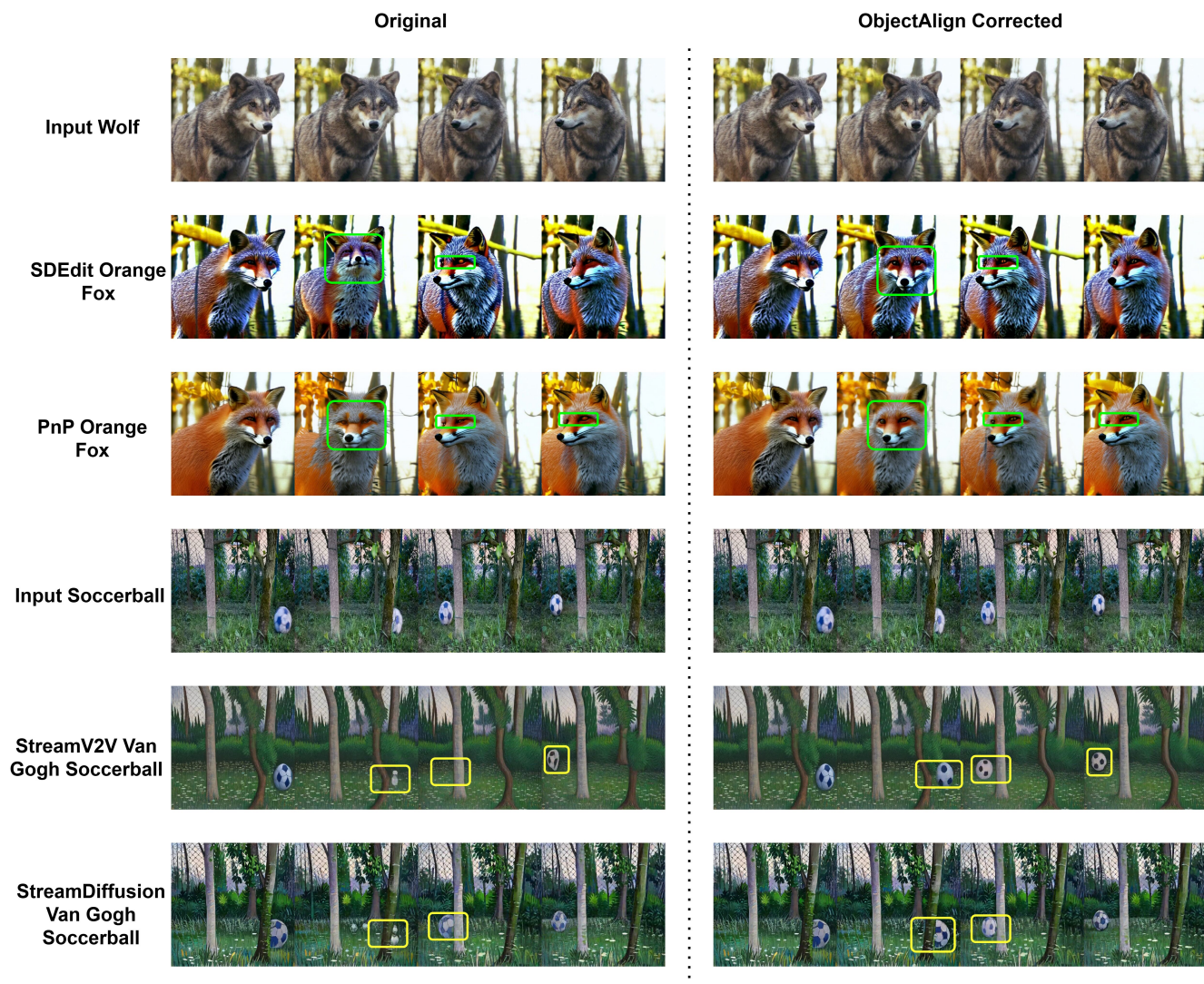


Figure 2. **Qualitative comparison of ObjectAlign corrections across different editing pipelines.** Before ObjectAlign correction (left), both SDEdit and PnP in the "Orange Fox" edits incorrectly alter the wolf's shape and color across consecutive frames (*highlighted in green boxes*). Similarly, StreamV2V and StreamDiffusion in the "Van Gogh Soccerball" edits cause the soccerball to intermittently disappear and reappear (*highlighted in yellow boxes*). These inconsistencies are accompanied by noticeable color and style drift, perceptual flickering, and identity misalignment. After applying ObjectAlign (right), these issues are effectively mitigated, resulting in greater semantic and temporal consistency.

402 The repaired frames replace the original inconsistent frames,  
403 producing an updated, more consistent video sequence.

404 The ObjectAlign pipeline re-runs Sections 4.1 to 4.2 on  
405  $V$  until  $\mathcal{I} = \emptyset$  as shown in Algorithm 1. ObjectAlign is able  
406 to provide both *empirical* quality and *formal* consistency  
407 guarantees through its neuro-symbolic verification pipeline.

## 408 5. Experimental Results

### 409 5.1. Experimental Setup

410 Our experiments utilize a dataset of 95 text-to-video prompt  
411 pairs obtained from the DAVIS [45] and Pexels [1] datasets.  
412 The prompts include both manually designed descriptions

and those inspired from [33, 53]. The videos cover a diverse  
413 array of subjects including animals and humans. The  
414 prompts encompass various scenarios involving object edits,  
415 style changes, and structural modifications. Our primary  
416 baselines for comparison are PnP [53], SDEdit [42],  
417 StreamV2V [33], and StreamDiffusion [27] without our  
418 object consistency corrections. Evaluations focus on object-  
419 level consistency and semantic fidelity using established  
420 metrics such as CLIP [47] scores and Warp Error [30]. All  
421 experiments were performed on one Nvidia 3090 GPU.  
422

Table 1. **ObjectAlign Improvements.** Comparison of video quality and consistency metrics *before* (Original) and *after* applying ObjectAlign (Edited) to videos processed by different base editing methods (PnP, SDEdit, StreamV2V, and StreamDiffusion). Scores are shown for different edit categories, with (+Improvement) indicating the improvement attributed to ObjectAlign. Higher scores are better for CLIP Score and VBench [20] motion smoothness, subject consistency, and background consistency metrics; lower is better for Warp Error.

Edit Category	PnP [53]		SDEdit [42]		StreamV2V [33]		StreamDiffusion [27]		
	Orig	Edited	Orig	Edited	Orig	Edited	Orig	Edited	
<b>Object Edits</b>	CLIP Score ↑	97.0	98.1 (+1.1)	96.7	98.1 (+1.4)	97.1	97.9 (+0.8)	95.2	96.0 (+0.8)
	Warp Error ↓	107.4	101.3 (+6.1)	105.5	100.8 (+4.7)	100.5	98.8 (+1.7)	108.7	103.5 (+5.2)
	Motion Smoothness ↑	0.917	0.930 (+0.013)	0.903	0.925 (+0.022)	0.916	0.935 (+0.019)	0.887	0.901 (+0.014)
	Subject Consistency ↑	0.913	0.925 (+0.012)	0.900	0.915 (+0.015)	0.920	0.931 (+0.011)	0.884	0.899 (+0.015)
	Background Consistency ↑	0.921	0.925 (+0.004)	0.904	0.917 (+0.013)	0.917	0.924 (+0.007)	0.892	0.908 (+0.016)
<b>Style Edits</b>	CLIP Score ↑	97.6	98.2 (+0.6)	97.3	98.0 (+0.7)	97.5	97.9 (+0.4)	95.8	96.4 (+0.6)
	Warp Error ↓	106.3	101.6 (+4.7)	105.3	100.5 (+4.8)	99.5	98.6 (+0.9)	107.2	103.3 (+3.9)
	Motion Smoothness ↑	0.938	0.973 (+0.035)	0.930	0.958 (+0.028)	0.933	0.940 (+0.007)	0.920	0.936 (+0.016)
	Subject Consistency ↑	0.905	0.937 (+0.032)	0.904	0.925 (+0.021)	0.912	0.928 (+0.016)	0.896	0.908 (+0.012)
	Background Consistency ↑	0.913	0.932 (+0.019)	0.903	0.920 (+0.017)	0.915	0.928 (+0.013)	0.900	0.913 (+0.013)
<b>Overall Average</b>	CLIP Score ↑	97.3	98.2 (+0.9)	97.0	98.1 (+1.1)	97.3	97.9 (+0.6)	95.5	96.2 (+0.7)
	Warp Error ↓	106.9	101.5 (+5.4)	105.4	100.7 (+4.7)	100.0	98.7 (+1.3)	108.0	103.4 (+4.6)
	Motion Smoothness ↑	0.928	0.952 (+0.024)	0.917	0.942 (+0.025)	0.925	0.938 (+0.013)	0.904	0.919 (+0.015)
	Subject Consistency ↑	0.909	0.931 (+0.022)	0.902	0.920 (+0.018)	0.916	0.930 (+0.014)	0.890	0.904 (+0.014)
	Background Consistency ↑	0.917	0.929 (+0.012)	0.905	0.919 (+0.014)	0.916	0.926 (+0.010)	0.896	0.911 (+0.015)

## 423 5.2. Qualitative Results

424 Figure 2 presents visual comparisons between ObjectAlign  
425 and the baseline PnP [53] and SDEdit [42] methods on mul-  
426 tiple challenging scenarios. We observe that ObjectAlign  
427 significantly reduces perceptual flicker, artifact generation,  
428 and object drift. In particular, ObjectAlign effectively main-  
429 tains stable object identities across frames, producing results  
430 noticeably smoother and more temporally coherent than base-  
431 line methods.

## 432 5.3. Quantitative Results

433 **CLIP Score.** We measure semantic consistency using the  
434 CLIP similarity score [47], defined as the cosine similarity  
435 of CLIP embeddings between consecutive frames. Higher  
436 scores reflect greater semantic stability (↑). When Objec-  
437 tAlign is applied to correct the outputs of various base  
438 editing methods, it consistently enhances semantic stabil-  
439 ity. For instance, drawing from the "Overall Average" re-  
440 sults in Table 1, applying ObjectAlign to videos edited by  
441 PnP improves the CLIP Score from an original 97.3 to  
442 98.2. For SDEdit, the score increases from 97.0 to 98.1;  
443 for StreamV2V [33], it improves from 97.3 to 97.9; and for  
444 StreamDiffusion [27], the score is enhanced from an original  
445 95.5 to 96.2. This demonstrates that ObjectAlign effectively  
446 preserves or improves semantic content preservation across  
447 frames when applied to a range of editing techniques.

448 **Warp Error.** Temporal coherence is evaluated via Warp  
449 Error [30], which computes pixel-wise discrepancies after  
450 warping edited frames by the original video’s optical flow.  
451 Lower Warp Error indicates greater temporal consistency (↓).  
452 When applied to various base editing methods, ObjectAlign  
453 consistently reduces their Warp Error, thereby enhancing  
454 temporal coherence. For instance, as detailed in Table 1  
455 (Overall Average section), ObjectAlign improves the Warp  
456 Error for PnP from an original score of 106.9 down to 101.5.  
457 Similarly, for SDEdit, the error is reduced from 105.4 to  
458 100.7; for StreamV2V, from 100.0 to 98.7; and for StreamD-  
459 iffusion, from 108.0 to 103.4. These results confirm that our  
460 method produces more temporally consistent videos when  
461 used to correct the outputs of established editing techniques.

462 **VBench Perceptual Metrics.** To further assess video qual-  
463 ity across diverse perceptual dimensions, we employ met-  
464 rics from the VBench benchmark [20], specifically Motion  
465 Smoothness, Subject Consistency, and Background Consis-  
466 tency. For these metrics, higher scores are preferable (↑).  
467 As shown in Table 1, ObjectAlign consistently improves  
468 these scores when applied to the outputs of different editing  
469 methods across both object and style edit categories. For  
470 instance, when ObjectAlign is applied to videos edited us-  
471 ing PnP, the Motion Smoothness score increases from 0.928  
472 to 0.952, and Subject Consistency improves from 0.909

473 to **0.931**. Similarly, for a baseline like StreamDiffusion,  
474 ObjectAlign enhances Motion Smoothness from 0.904 to  
475 **0.919** and Subject Consistency from 0.890 to **0.904**. These  
476 examples, representative of the broader findings in Table 1,  
477 indicate enhanced visual quality in terms of smoother motion,  
478 more stable subject appearance, and more coherent  
479 backgrounds. Further details on the VBench benchmark can  
480 be found in the original VBench documentation [20].

481 **Runtime Efficiency.** The runtime overhead introduced by  
482 ObjectAlign is minimal, requiring approximately 3% additional  
483 computation time compared to baseline PnP [53]  
484 editing, and 4% additional runtime compared to SDEdit  
485 [42]. The runtime overhead is primarily due to adaptive  
486 interpolation and SMT-based verification. The average processing  
487 time per frame remains acceptable for practical use,  
488 enabling ObjectAlign integration into existing image and  
489 video-editing workflows even in compute-limited edge environments.  
490 Overall, ObjectAlign achieves a superior balance  
491 between semantic consistency, temporal coherence, and computational  
492 efficiency compared to existing methods.

## 493 5.4. Ablation Studies

494 **Ablation Study on Diffusion-Based Inpainting for Frame  
495 Repair.** To evaluate the efficacy of our adaptive interpolation  
496 for frame repair (Section 4.2), we conducted an ablation  
497 study comparing it against an alternative approach using a  
498 pre-trained Stable Diffusion (SD) inpainting pipeline [48].  
499 For this experiment, inconsistent frame outputs of PnP [53]  
500 were targeted for repair. Segmentation masks obtained via  
501 SAM [26] from a consistent reference frame guided the in-  
502 painting region, and textual prompts were provided. As  
503 shown in Table 2, the SD inpainting method yielded minimal  
504 beneficial impact on key metrics such as CLIP Score and  
505 Warp Error when applied to repair inconsistent PnP outputs,  
506 improving them by only 0.1 points. In contrast, ObjectAlign’s  
507 interpolation demonstrates substantial improvements  
508 on CLIP Score [47] and warp error [30] for the same PnP  
509 outputs. These findings support our choice of targeted interpolation  
510 from consistent anchor frames for frame repair.

Table 2. **Ablation Study: Frame Repair Methods for PnP Outputs.** Comparison of ObjectAlign’s RIFE-based interpolation against Stable Diffusion (SD) inpainting for repairing frames from the PnP baseline. "Original" refers to PnP output before repair. SD Inpainting refers to scores after the inpainting based repair method. Improvements (+Value) are relative to PnP (Original). The ablation study is performed over 18 edited video sequences.

Metric	Original	SD Inpainting	Interpolation (ObjectAlign)
CLIP Score ↑	97.0	97.1 (+0.1)	97.7 (+0.7)
Warp Error ↓	106.4	106.6 (-0.2)	100.3 (+6.1)

**Ablation Study on Verification Checks.** We ablate the impact of individual consistency verification checks—semantic similarity (CLIP cosine), perceptual similarity (LPIPS), histogram correlation, object-mask overlap (IoU), and the SMT-based semantic embedding constraint—on identifying inconsistent frames. This study is performed over 35 edited video sequences of lengths between 40 and 280 frames each. As detailed in Table 3, we observe that the IoU-based object-mask consistency check is most frequently triggered (22.3% of total frames), reflecting its sensitivity to spatial discrepancies in segmentation masks. The SMT-based embedding constraint triggers second most often (16.6%), underscoring the benefit of formal semantic bounds. The perceptual LPIPS check triggers third (15.5%), highlighting its effectiveness at detecting subtle visual artifacts. Histogram correlation and CLIP-based semantic similarity check flag inconsistencies less frequently (8.7% and 7.1%, respectively), indicating that color and global semantic shifts are comparatively rarer. Overall, the combined use of complementary verification checks ensures robust detection of diverse inconsistency types, each targeting different aspects of perceptual, spatial, and semantic coherence.

Table 3. **Ablation on individual consistency verification checks.** We report the percentage of frames flagged as inconsistent by each verification check over all sequences. Higher percentage indicates greater sensitivity of the check in detecting inconsistencies.

Verification Check	IoU	SMT	LPIPS	Histogram	CLIP Cosine
Percentage flagged (%)	22.3	16.6	15.5	8.7	7.1

## 513 6. Conclusion

514 In this paper, we have introduced ObjectAlign, a neuro-  
515 symbolic framework designed to detect, formally verify,  
516 and adaptively correct object-level inconsistencies in edited  
517 video sequences. Our approach integrates learnable perceptual  
518 metrics, neuro-symbolic verification, and adaptive neural  
519 network based interpolation to ensure semantic fidelity,  
520 temporal fidelity, and visual coherence.

521 Experimental evaluations demonstrate ObjectAlign’s ability  
522 to substantially reduce semantic drift and visual artifacts,  
523 achieving superior performance in both perceptual consistency  
524 and temporal coherence, compared to existing baseline  
525 methods. Furthermore, ablation studies confirm the importance  
526 of each component in our design, highlighting the effectiveness  
527 of combining learnable consistency thresholds, symbolic reasoning,  
528 and adaptive interpolation. ObjectAlign thus represents an important  
529 step towards provably consistent and visually stable video editing.  
530

551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
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  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606

## References

- [1] <https://www.pexels.com>. 2, 6
- [2] Christel Baier and Joost-Pieter Katoen. *Principles of Model Checking*. The MIT Press, 2008. 3
- [3] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *International Conference on Machine Learning*, pages 813–824. PMLR, 2021. 3
- [4] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2022. 2
- [5] Duygu Ceylan, Chun-Hao Paul Huang, and Niloy Jyoti Mitra. Pix2video: Video editing using image diffusion. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 23149–23160, 2023. 2
- [6] Zhenfang Chen, Kexin Yi, Yunzhu Li, Mingyu Ding, Antonio Torralba, Joshua B Tenenbaum, and Chuang Gan. Comphy: Compositional physical reasoning of objects and events from videos. *arXiv preprint arXiv:2205.01089*, 2022. 3
- [7] Minkyu Choi, Harsh Goel, Mohammad Omama, Yunhao Yang, Sahil Shah, and Sandeep Chinchali. Towards neuro-symbolic video understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, Milan, Italy, 2024. Springer. 3
- [8] Minkyu Choi, S P Sharan, Harsh Goel, Sahil Shah, and Sandeep Chinchali. We’ll fix it in post: Improving text-to-video generation with neuro-symbolic feedback, 2025. 5
- [9] Brandon C. Colelough and William Regli. Neuro-symbolic ai in 2024: A systematic review, 2025. 2
- [10] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *arXiv preprint arXiv:2209.04747*, 2022. 1
- [11] Leonardo De Moura and Nikolaj Bjørner. Satisfiability modulo theories: introduction and applications. *Commun. ACM*, 54(9):69–77, 2011. 1, 3
- [12] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 3
- [13] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *IEEE/CVF International Conference on Computer Vision*, pages 6201–6210. IEEE, 2019. 3
- [14] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373*, 2023. 1, 2
- [15] Mohammadhosein Hasanbeigi, Yiannis Kantaros, Alessandro Abate, Daniel Kroening, George J Pappas, and Insup Lee. Reinforcement learning for temporal logic control synthesis with probabilistic satisfaction guarantees. In *2019 IEEE 58th conference on decision and control (CDC)*, pages 5338–5343. IEEE, 2019. 3
- [16] Pascal Hitzler, Aaron Eberhart, Monireh Ebrahimi, Md Kamruzzaman Sarker, and Lu Zhou. Neuro-symbolic approaches in artificial intelligence. *National Science Review*, 9(6):nwac035, 2022. 2
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3
- [18] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv:2204.03458*, 2022. 1
- [19] Zhewei Huang, Tianyuan Zhang, Wen Heng, Boxin Shi, and Shuchang Zhou. Real-time intermediate flow estimation for video frame interpolation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022. 5
- [20] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 7, 8
- [21] Hyeonho Jeong and Jong Chul Ye. Ground-a-video: Zero-shot grounded video editing using text-to-image diffusion models. *ArXiv*, abs/2310.01107, 2023. 2
- [22] Susmit Jha, Vasumathi Raman, Dorsa Sadigh, and Sanjit A Seshia. Safe autonomy under perception uncertainty using chance-constrained temporal logic. *Journal of Automated Reasoning*, 60:43–62, 2018. 3
- [23] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators, 2023. 2
- [24] Jaemin Kim, Bryan S Kim, and Jong Chul Ye. Free<sup>2</sup>guide: Gradient-free path integral control for enhancing text-to-video generation with large vision-language models, 2024. 2
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 4
- [26] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 4, 8
- [27] Akio Kodaira, Chenfeng Xu, Toshiaki Hazama, Takanori Yoshimoto, Kohei Ohno, Shogo Mitsuohori, Soichi Sugano, Hanying Cho, Zhijian Liu, and Kurt Keutzer. Streamdiffusion: A pipeline-level solution for real-time interactive generation. *arXiv preprint arXiv:2312.12491*, 2023. 1, 2, 6, 7
- [28] Hadas Kress-Gazit, Georgios E Fainekos, and George J Pappas. Temporal-logic-based reactive mission and motion planning. *IEEE transactions on robotics*, 25(6):1370–1381, 2009. 3
- [29] Aliaksandr Krushchanka, Vladimir Golovko, Egor Mikhno, Mikhail Kovalev, Vadim Zahariev, and Aleksandr Zagorskij. A neural-symbolic approach to computer vision. In *International Conference on Open Semantic Technologies for Intelligent Systems*, pages 282–309. Springer, 2021. 3
- [30] Wei-Sheng Lai, Jia-Bin Huang, Oliver Wang, Eli Shechtman, Ersin Yumer, and Ming-Hsuan Yang. Learning blind video temporal consistency. In *Proceedings of the European conference on computer vision (ECCV)*, pages 170–185, 2018. 2, 6, 7, 8

- 665 [31] Nanjun Li, Faliang Chang, and Chunsheng Liu. Human-  
666 related anomalous event detection via spatial-temporal graph  
667 convolutional autoencoder with embedded long short-term  
668 memory network. *Neurocomputing*, 490:482–494, 2022. 3
- 669 [32] Feng Liang, Bichen Wu, Jialiang Wang, Licheng Yu, Kun-  
670 peng Li, Yanan Zhao, Ishan Misra, Jia-Bin Huang, Peizhao  
671 Zhang, Peter Vajda, et al. Flowvid: Taming imperfect optical  
672 flows for consistent video-to-video synthesis. *arXiv preprint*  
673 *arXiv:2312.17681*, 2023. 1
- 674 [33] Feng Liang, Akio Kodaira, Chenfeng Xu, Masayoshi  
675 Tomizuka, Kurt Keutzer, and Diana Marculescu. Looking  
676 backward: Streaming video-to-video translation with feature  
677 banks. *arXiv preprint arXiv:2405.15757*, 2024. 1, 2, 6, 7
- 678 [34] Shaoteng Liu, Yuechen Zhang, Wenbo Li, Zhe Lin, and Jiaya  
679 Jia. Video-p2p: Video editing with cross-attention control.  
680 *Proceedings of the IEEE/CVF Conference on Computer Vi-*  
681 *sion and Pattern Recognition*, pages 8599–8608, 2024. 2
- 682 [35] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao,  
683 Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jian-  
684 feng Gao, et al. Sora: A review on background, technology,  
685 limitations, and opportunities of large vision models. *arXiv*  
686 *preprint arXiv:2402.17177*, 2024. 1
- 687 [36] Yang Luo, Xuanlei Zhao, Mengzhao Chen, Kaipeng Zhang,  
688 Wenqi Shao, Kai Wang, Zhangyang Wang, and Yang You.  
689 Enhance-a-video: Better generated video for free, 2025. 2
- 690 [37] Tanvir Mahmud, Mustafa Munir, Radu Marculescu, and Di-  
691 ana Marculescu. Ada-ve: Training-free consistent video edit-  
692 ing using adaptive motion prior. In *2025 IEEE/CVF Win-*  
693 *ter Conference on Applications of Computer Vision (WACV)*,  
694 pages 940–949. IEEE, 2025. 1, 2
- 695 [38] Zohar Manna and Amir Pnueli. *The Temporal Logic of Reac-*  
696 *tive and Concurrent Systems: Specification*. Springer-Verlag,  
697 1992. 5
- 698 [39] Effrosyni Mavroudi, Benjamín Béjar Haro, and René Vidal.  
699 Representation learning on visual-symbolic graphs for video  
700 understanding. In *European Conference on Computer Vision*,  
701 pages 71–90. Springer, 2020. 3
- 702 [40] Gérard G. Medioni, Isaac Cohen, François Brémond, Som-  
703 boon Hongeng, and Ramakant Nevatia. Event detection and  
704 analysis from video streams. *IEEE Trans. Pattern Anal. Mach.*  
705 *Intell.*, 23(8):873–889, 2001. 3
- 706 [41] Noushin Mehdipour, Matthias Althoff, Radboud Duintjer  
707 Tebbens, and Calin Belta. Formal methods to comply with  
708 rules of the road in autonomous driving: State of the art and  
709 grand challenges. *Automatica*, 152:110692, 2023. 3
- 710 [42] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun  
711 Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: Guided image  
712 synthesis and editing with stochastic differential equations. In  
713 *International Conference on Learning Representations*, 2022.  
714 1, 2, 6, 7, 8
- 715 [43] Eyal Molad, Eliahu Horwitz, Dani Valevski, Alex Rav Acha,  
716 Yossi Matias, Yael Pritch, Yaniv Leviathan, and Yedid Hoshen.  
717 Dreamix: Video diffusion models are general video editors,  
718 2023. 2
- 719 [44] Mustafa Munir, Saloni Modi, Geffen Cooper, Huntae Kim,  
720 and Radu Marculescu. Three decades of low power: From  
721 watts to wisdom. *IEEE Access*, 12:19447–19458, 2024. 1
- [45] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Ar-  
belález, Alex Sorkine-Hornung, and Luc Van Gool. The 2017  
davis challenge on video object segmentation. *arXiv preprint*  
*arXiv:1704.00675*, 2017. 2, 6
- [46] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei,  
Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fus-  
ing attentions for zero-shot text-based video editing. *2023*  
*IEEE/CVF International Conference on Computer Vision*  
*(ICCV)*, pages 15886–15896, 2023. 2
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya  
Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,  
Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning  
transferable visual models from natural language supervi-  
sion. In *International conference on machine learning*, pages  
8748–8763. PMLR, 2021. 3, 4, 6, 7, 8
- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz,  
Patrick Esser, and Björn Ommer. High-resolution image  
synthesis with latent diffusion models. In *Proceedings of*  
*the IEEE/CVF Conference on Computer Vision and Pattern*  
*Recognition*, pages 10684–10695, 2022. 8
- [49] Soumalya Sarkar, Kin Gwn Lore, and Soumik Sarkar. Early  
detection of combustion instability by neural-symbolic analy-  
sis on hi-speed video. In *Proceedings of the NIPS Workshop*  
*on Cognitive Computation: Integrating Neural and Symbolic*  
*Approaches co-located with the 29th Annual Conference on*  
*Neural Information Processing Systems*, Montreal, Canada,  
2015. CEUR-WS.org. 3
- [50] SP Sharan, Minkyu Choi, Sahil Shah, Harsh Goel, Moham-  
mad Omama, and Sandeep Chinchali. Neuro-symbolic eval-  
uation of text-to-video models using formal verification. In  
*Proceedings of the Computer Vision and Pattern Recognition*  
*Conference*, pages 8395–8405, 2025. 3, 5
- [51] Yasser Shoukry, Pierluigi Nuzzo, Ayca Balkan, Indranil Saha,  
Alberto L Sangiovanni-Vincentelli, Sanjit A Seshia, George J  
Pappas, and Paulo Tabuada. Linear temporal logic motion  
planning for teams of underactuated robots using satisfiability  
modulo convex programming. In *2017 IEEE 56th annual*  
*conference on decision and control (CDC)*, pages 1132–1137.  
IEEE, 2017. 3
- [52] Du Tran, Heng Wang, Matt Feiszli, and Lorenzo Torresani.  
Video classification with channel-separated convolutional net-  
works. In *IEEE/CVF International Conference on Computer*  
*Vision*, pages 5551–5560. IEEE, 2019. 3
- [53] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel.  
Plug-and-play diffusion features for text-driven image-to-  
image translation. *Proceedings of the IEEE/CVF Conference*  
*on Computer Vision and Pattern Recognition*, 2023. 1, 2, 6,  
7, 8
- [54] Wen Wang, Yan Jiang, Kangyang Xie, Zide Liu, Hao Chen,  
Yue Cao, Xinlong Wang, and Chunhua Shen. Zero-shot video  
editing using off-the-shelf image diffusion models, 2024. 2
- [55] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian  
Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and  
Mike Zheng Shou. Tune-a-video: One-shot tuning of image  
diffusion models for text-to-video generation. *arXiv preprint*  
*arXiv:2212.11565*, 2022. 2
- [56] Yu Xiong, Qingqiu Huang, Lingfeng Guo, Hang Zhou, Bolei  
Zhou, and Dahua Lin. A graph-based framework to bridge

- 780 movies and synopses. In *IEEE/CVF International Conference*  
781 *on Computer Vision*, pages 4591–4600. IEEE, 2019. 3
- 782 [57] Zhongwen Xu, Yi Yang, and Alexander G. Hauptmann. A  
783 discriminative CNN video representation for event detection.  
784 In *IEEE Conference on Computer Vision and Pattern Recog-*  
785 *nition*, pages 1798–1807, Boston, MA, USA, 2015. IEEE  
786 Computer Society. 3
- 787 [58] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy.  
788 Rerender a video: Zero-shot text-guided video-to-video trans-  
789 lation. *SIGGRAPH Asia 2023 Conference Papers*, 2023. 2
- 790 [59] Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Push-  
791 meet Kohli, and Josh Tenenbaum. Neural-symbolic VQA:  
792 disentangling reasoning from vision and language understand-  
793 ing. In *Advances in Neural Information Processing Systems*,  
794 pages 1039–1050, 2018. 3
- 795 [60] Dongran Yu, Bo Yang, Qianhao Wei, Anchen Li, and Shirui  
796 Pan. A probabilistic graphical model based on neural-  
797 symbolic reasoning for visual relationship detection. In  
798 *IEEE/CVF Conference on Computer Vision and Pattern*  
799 *Recognition*, pages 10599–10608, New Orleans, LA, USA,  
800 2022. IEEE. 3
- 801 [61] David Junhao Zhang, Dongxu Li, Hung Le, Mike Zheng  
802 Shou, Caiming Xiong, and Doyen Sahoo. Moonshot: Towards  
803 controllable video generation and editing with multimodal  
804 conditions. *arXiv preprint arXiv:2401.01827*, 2024. 2
- 805 [62] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding  
806 conditional control to text-to-image diffusion models. In  
807 *Proceedings of the IEEE/CVF international conference on*  
808 *computer vision*, pages 3836–3847, 2023. 2
- 809 [63] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman,  
810 and Oliver Wang. The unreasonable effectiveness of deep  
811 features as a perceptual metric. In *Proceedings of the IEEE*  
812 *conference on computer vision and pattern recognition*, pages  
813 586–595, 2018. 3, 4