
Frame-Level Captions for Long Video Generation with Complex Multi Scenes

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Generating long videos that can show complex stories, like movie scenes from
2 scripts, has great promise and offers much more than short clips. However, current
3 methods that use autoregression with diffusion models often struggle because
4 their step-by-step process naturally leads to a serious error accumulation (drift).
5 Also, many existing ways to make long videos focus on single, continuous scenes,
6 making them less useful for stories with many events and changes. This paper
7 introduces a new approach to solve these problems. First, we propose a novel
8 way to annotate datasets at the **frame-level**, providing detailed text guidance
9 needed for making complex, multi-scene long videos. This detailed guidance
10 works with a **Frame-Level Attention Mechanism** to make sure text and video
11 match precisely. In inference, we develop **Parallel Multi-Window Denoising**, a
12 new method that handles a long video as multiple overlapping windows. These
13 windows are processed in parallel, and the noise prediction in overlapping areas
14 is averaged, which allows bidirectional information interaction and introduces no
15 error accumulation. A key feature is that each part (frame) within these windows
16 can be guided by its own distinct text prompt. Our training uses **Diffusion Forcing**
17 to provide the model with the ability to handle time flexibly. We tested our approach
18 on difficult VBench 2.0 benchmarks ("Complex Plots" and "Complex Landscapes")
19 based on the WanX2.1-T2V-1.3B model. The results show our method is better at
20 following instructions in complex, changing scenes and creates high-quality long
21 videos. We plan to share our dataset annotation methods and trained models with
22 the research community.

23 1 Introduction

24 The ability to create long video sequences from text instructions opens exciting doors for rich,
25 evolving stories, such as turning scripts into videos, producing short films, or showing complex
26 processes. Unlike short clips, long videos provide the needed duration for multiple connected scenes,
27 detailed character interactions, and consistent plotlines that follow complex user requests [17, 42, 32,
28 23, 15, 20, 22]. However, creating high-quality, consistent, and accurate long videos from text is still
29 a major challenge for current generative models.

30 A primary difficulty lies in the common autoregressive (step-by-step) methods used with diffusion
31 models to make longer videos. Their sequential way of working naturally leads to errors accumulation
32 over time. This shows up as lower visual quality, the video drifting away from the original text's
33 meaning, and a loss of consistency as the video gets longer, seriously weakening the quality of
34 extended generations [14, 18, 37, 36]. Furthermore, much current research on long videos deals with
35 single, continuous scenes or slowly changing environments. This limited focus reduces the usefulness
36 of long video generation for dynamic stories with many events, which is a key goal for creative

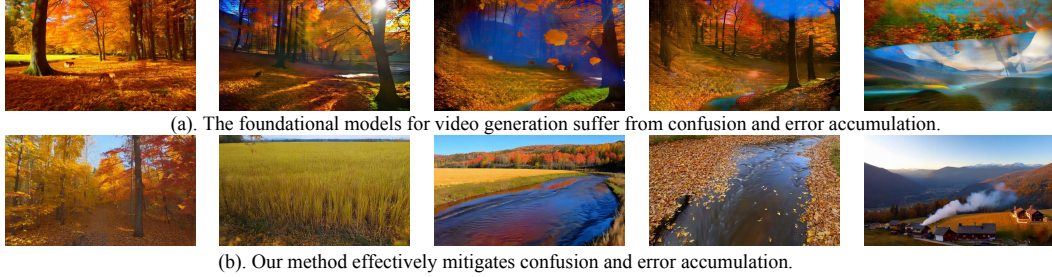


Figure 1: Illustration of issues of error accumulation and semantic confusion in the first row, and videos generated by ours are shown in the second row.

37 applications. Standard methods using single, general (global) video descriptions struggle with small,
 38 quick changes, leading to timing issues [26, 18, 19]. While some approaches try to make multi-shot
 39 videos by first detecting scene cuts with tools like PySceneDetect and then processing these shots, as
 40 in Presto [30] and Long Context Tuning (LCT) [12], these methods can be complicated, risk losing
 41 information, and depend heavily on good shot detection and captioning. They often still describe
 42 video at a "shot-text" level, which doesn't fully capture smooth, continuous changes.

43 To solve these basic problems, our work offers a new way of thinking, focused on highly detailed,
 44 frame-level text guidance and a novel non-sequential method for creating the video. Our first main
 45 contribution is an innovative **frame-level dataset annotation methodology**. We move beyond
 46 general or shot-level captions to provide very detailed text descriptions for each conceptual part (or
 47 latent segment) of a video. This rich information about meaning is essential for guiding models to
 48 understand and create the complex, changing details needed for stories with many scenes and detailed
 49 prompts. This directly addresses the limits of less detailed global or shot-level captions. This detailed
 50 annotation is designed to work closely with our **Frame-Level Attention Mechanism**, which clearly
 51 links each video segment's visual features to its specific text description, improving content accuracy
 52 and consistency over time (Section 4.2).

53 To properly use such detailed and dynamic text prompts, models need to be trained to handle time
 54 flexibly. We achieve this using **Diffusion Forcing** (Section 4.3), a training strategy that shows
 55 the model video segments being denoised at different rates. This prepares it to manage varied
 56 timing patterns and allows for strong, adaptable inference. Building on these training improvements,
 57 we introduce our second major innovation: **Parallel Multi-Window Denoising (PMWD)**, a new
 58 inference method designed to create very long videos that are highly consistent (Section 4.4). PMWD
 59 divides the target long video into multiple overlapping sections (windows), usually matching the
 60 model's training window size. Importantly, unlike step-by-step methods, all these windows are
 61 processed *at the same time (in parallel)* during each step of the diffusion denoising process. The data
 62 in the overlapping areas between windows is then averaged. This averaging not only ensures smooth
 63 connections but also allows information to flow in both directions, meaning later parts of the video
 64 can help refine earlier ones. A special feature of PMWD is that each conceptual frame, even within
 65 these simultaneously processed windows, can be guided by its own distinct, frame-level prompt.

66 We test our approach thoroughly using highly challenging benchmarks, specifically the "Complex
 67 Plots" and "Complex Landscapes" prompt categories from VBench 2.0 [45]. We use the state-of-
 68 the-art open-source WanX2.1-T2V-1.3B model [32] as our base. Our experiments show that our
 69 combined frame-level approach is much better at following prompts when creating very long videos
 70 with multiple scenes, different characters, and complex actions.

71 In summary, our main contributions are:

- 72 • **Scalable Frame-Level Dataset Methodology:** We introduce an efficient and scalable approach for
 73 constructing datasets with dense, frame-by-frame textual annotations. This enables highly granular
 74 video-text alignment crucial for generating complex, multi-scene narratives without relying on
 75 traditional shot detection.
- 76 • **Frame-Level Attention for Precise Guidance:** We propose a novel attention mechanism that
 77 directly couples each video segment's visual features with its unique frame-level prompt. This
 78 significantly enhances semantic fidelity, content accuracy, and temporal consistency in generated
 79 videos.

- 80 • **Parallel Multi-Window Denoising for Coherent Long Video Generation:** We develop PMWD,
81 a inference strategy that processes a long video as multiple overlapping windows, denoised
82 simultaneously in parallel. Guided by distinct frame-level prompts and leveraging overlap averaging
83 for bidirectional context, PMWD effectively avoids the error accumulation common in sequential
84 methods. This is enabled by training strategies like Diffusion Forcing that provide temporal
85 flexibility.
- 86 • **State-of-the-Art Performance on Complex Videos:** Through comprehensive evaluations on
87 challenging VBench 2.0 benchmarks ("Complex Plots" and "Complex Landscapes") using the
88 WanX2.1-T2V-1.3B model, we demonstrate our integrated approach’s superior ability to follow
89 intricate prompts in multi-element long videos, achieving high-fidelity results with minimal error
90 accumulation.

91 2 Related Work

92 **Video Generation Dataset.** Large-scale video datasets have driven advancements in video generation,
93 but many existing datasets like YouCook2 [47], VATEX [38], and ActivityNet [5] were not designed
94 for this purpose and lack fine-grained annotations. Similarly, large-scale datasets like YTTemporal-
95 180M [43] and HD-VILA-100M [40] suffer from low-quality captions generated through speech
96 recognition, limiting their utility for high-quality video generation. Datasets like Panda-70M [8]
97 offer extensive data but rely on simplistic global descriptions, which hinder the model’s ability to
98 capture fine temporal details. Newer datasets, such as Koala-36M [35] and LongTake-HD [41],
99 provide more detailed annotations but still rely on segment-level or shot-based annotations, limiting
100 long-duration video generation. In contrast, our method introduces a frame-level captioning approach,
101 where each frame is independently annotated with a description that maintains contextual relevance
102 to the preceding and succeeding frames. This ensures better alignment between visual content and
103 text while preserving the temporal continuity and motion dynamics, ultimately improving the overall
104 quality of long-form video generation.

105 **Long Video Generation.** Video generation has evolved from simple single-shot models to more com-
106 plex long-form and multi-scene models. Early methods relied on GANs [9, 29, 31, 39], constrained
107 by single-domain datasets. Diffusion models [3, 11, 2, 1, 10] introduced temporal layers, enabling
108 motion modeling. DiT-based architectures [4, 25, 28, 17, 42, 32, 23, 15, 20] have achieved tremen-
109 dous success in scaling diffusion transformers, significantly enhancing video quality. However, these
110 models were limited to generating short clips. FreeNoise [27] and StreamingT2V [13] extended video
111 sequences using auto-regressive methods and temporal attention mechanisms. Gen-L-Video [34]
112 processes videos as sequences of overlapping short clips and employs a temporal co-denoising
113 technique, wherein multiple predictions for each individual frame are averaged. Despite these
114 advancements, challenges in content diversity and temporal consistency persisted. The Diffusion
115 Forcing [6] paradigm addressed these issues by combining diffusion’s high-quality generation with
116 auto-regressive models for sequence extension.

117 In multi-scene video generation, models like Mask2DiT [26], LCT [12], VideoStudio [21],
118 SKYREELS-V2 [7], MovieDreamer [44], StoryAnchors [33], and VGoT [46] focused on scene-level
119 consistency but struggled with temporal coherence across scenes. Recent methods, including Sto-
120 ryDiffusion [48] and MEVG [24], employed attention mechanisms to enhance visual and dynamic
121 consistency. Our approach uses frame-level attention for dynamic scene extension without fixed
122 scene durations, improving flexibility and coherence in long-form videos. Combined with Diffusion
123 Forcing [6], our method ensures smooth scene transitions, extended video lengths, and maintains
124 both visual richness and temporal consistency.

125 3 Frame-Level Dataset

126 Previous video-text datasets such as Panda70M [8] and Koala-36M [35] provide only global-level
127 captions, resulting in coarse supervision that cannot reflect detailed visual changes within videos.
128 LongTake-HD [41] offers shot-level sub-captions but still depends on explicit shot boundaries, making
129 it difficult to model continuous motion and intra-shot dynamics. In contrast, our dataset uses frame-
130 level uniform sampling and annotation, enabling dense and temporally continuous supervision. This
131 design captures both subtle and significant changes without being limited by artificial segmentation,
132 supporting more precise alignment between video frames and text descriptions. Overall, our dataset

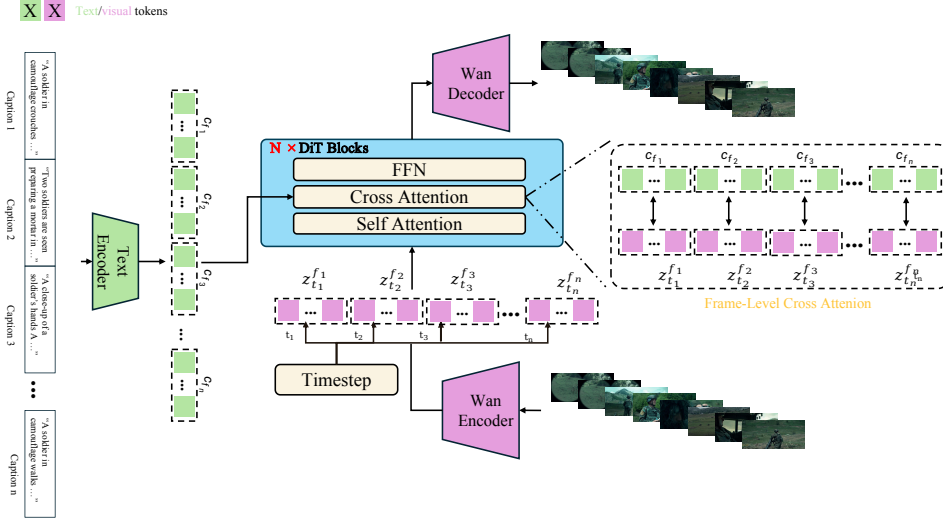


Figure 2: **Overview of the proposed frame-level training method.** Frame-Level Cross-Attention links the visual data of each video segment (latent token) directly to its own specific text description.

133 offers finer-grained, structurally consistent, and temporally faithful video-text supervision, facilitating
 134 improved learning of dynamic visual content.

135 **Large-scale frame-level video dataset construction.** We present a frame-level video dataset
 136 comprising 700,000 high-quality clips, designed to enhance fine-grained text-video alignment and
 137 provide dense semantic supervision for diffusion-based video generation models. The dataset
 138 systematically balances visual diversity, temporal continuity, and annotation precision, and can be
 139 further improved with larger scale in the future. We collect raw videos longer than 10 minutes from
 140 multiple platforms, remove near-duplicate content using perceptual hashing, and discard the first and
 141 last 10% of frames to ensure the sampled content is dynamic and semantically meaningful. Each
 142 processed video is evenly divided into four segments, from which an 8-second continuous clip is
 143 extracted. At a frame rate of 24 fps, one frame is sampled every 8 frames, resulting in approximately
 144 24 frames per clip. This design balances temporal context and computational efficiency, and ensures
 145 compatibility with mainstream video VAE tokenization schemes, enabling precise one-to-one frame-
 146 token supervision.

147 **Adaptive frame-level annotation.** We utilize multimodal large language models to generate frame-
 148 level captions, automatically choosing between shared or independent descriptions based on the
 149 degree of visual change. Identical captions are assigned to frames with minimal differences, while
 150 significant changes trigger independent frame-level descriptions, achieving unified and adaptive
 151 semantic alignment. To enhance consistency and richness, we design structured annotation prompts
 152 that require each frame description to cover main subjects, actions, environment, shot size, and
 153 camera angle. Outputs strictly follow JSON format with no redundant commentary, ensuring precise
 154 semantic and structural alignment at the frame level. We provide the full frame-level annotation
 155 prompt, specifying formatting and content requirements to facilitate reproducibility and further
 156 research; the template is too long to show here, see appendix. During inference, we use gemini pro
 157 2.5 to convert a user input from short/detailed caption to a frame-level detailed caption. More details
 158 can be found in appendix.

159 4 Method

160 Our work introduces a set of interconnected improvements for DiT-based video diffusion models,
 161 aimed at generating complex, long videos. We focus on enhancing how models understand text
 162 instructions for each video part, how they handle timing and changes, and how they create long,
 163 consistent videos during inference.

164 4.1 Overall Framework

165 Current Diffusion Transformer (DiT) based models are skilled at creating high-quality short videos.
166 They typically use a VAE (Variational Autoencoder) to compress videos into compact latent data, and
167 a DiT then generates video from these latents. However, when generating long videos with detailed
168 stories and dynamic action, these models face several basic problems:

- 169 • **Imprecise Content Control:** Using a single text description (caption) for the entire video often
170 leads to unclear or mixed-up details for different parts of the video, making it hard to accurately
171 control specific events or elements across various scenes.
- 172 • **Limited Handling of Timing:** Standard models usually denoise all parts of the video at the same
173 rate in each step. This restricts their ability to show varied motion speeds, different pacing in
174 scenes, or sudden changes effectively.
- 175 • **Difficulty with Long Video Coherence:** When creating long videos by generating segments one
176 after another from models trained on short clips, errors tend to build up. This can make the video
177 lose consistency over time.

178 To address these key challenges, we propose three main contributions:

- 179 1. A **Frame-Level Cross-Attention** mechanism (Section 4.2) for precise, localized text-based control
180 over the content of each video segment during training.
- 181 2. A **Diffusion Forcing** training strategy (Section 4.3) to teach models how to handle varied timing
182 by exposing them to video segments denoised at different rates.
- 183 3. A **Synchronized Multi-Window Denoising (PMWD)** inference method (Section 4.4), designed
184 to generate long, coherent videos by significantly reducing the build-up of errors.

185 While these principles can be applied to many DiT-based video models, we demonstrate our methods
186 by adapting and fine-tuning the WanX2.1-T2V-1.3B [32] framework, a well-known open-source
187 model, as our base.

188 4.2 Frame-Level Cross Attention

189 To accurately control video content in line with detailed narratives, we introduce Frame-Level Cross-
190 Attention. This method links the visual data of each video segment (latent token) directly to its own
191 specific text description. Simultaneously, the DiT’s standard self-attention mechanism continues to
192 capture overall temporal relationships, ensuring smooth motion. This approach provides both precise
193 local content guidance and global video coherence.

194 Our process starts by assigning an independent text description to each conceptual "frame" (latent
195 unit) of a video. When the original video is converted into latent data by the VAE, each resulting
196 latent token z_f is directly paired with the embedding of its corresponding frame-level caption, c_f .
197 This creates a detailed, one-to-one mapping between text and video segments over time, offering
198 exact guidance for generation. We modify the DiT’s cross-attention mechanism so that each latent
199 token z_f attends exclusively to its paired caption embedding c_f , rather than to a single caption shared
200 by the entire video. Formally, this is:

$$\text{CrossAttention}(q_f, c_f) = \text{Softmax} \left(\frac{q_f W_q (c_f W_k)^T}{\sqrt{d}} \right) (c_f W_v), \quad (1)$$

201 where q_f is the query projected from z_f , and W_q, W_k, W_v are learnable matrices. This targeted
202 attention mechanism reduces the unclear meaning that can arise from global captions, greatly
203 improving text-to-video alignment and allowing for precise control over dynamic content within each
204 segment.

205 4.3 Diffusion Forcing for Temporal Flexibility

206 Creating long videos with dynamic action and varied pacing requires the model to handle time flexibly.
207 Standard diffusion models are often limited in this area because they apply the same noise level to
208 all video segments at each step of the denoising process, restricting their ability to generate diverse
209 visual qualities, dynamic changes, or quick scene transitions.

210 To give models this needed flexibility, we use a **Diffusion Forcing (DF)** strategy during training.
211 This technique assigns an independent noise level to each video segment (latent token) in a training
212 sequence. Specifically, we pick a reference segment, set its target noise removal stage (timestep), and
213 then determine the noise stages for other segments in relation to it: preceding segments get "cleaner"
214 (earlier) timesteps, and subsequent segments get "noisier" (later) timesteps. This approach maintains
215 temporal smoothness while training the model to manage different denoising states simultaneously
216 within one sequence.

217 This training approach makes the model highly adaptable at inference time. By adjusting a "step-size"
218 parameter—which controls the allowed difference in noise schedules between adjacent segments—we
219 can smoothly shift the generation style. We can opt for fully synchronized diffusion (small step-
220 size, for high consistency) or for more dynamic, evolving outputs (large step-size, resembling
221 autoregressive generation). This adaptability allows the model to produce either smooth, consistent
222 videos or to progressively unfold complex scene transitions and actions as guided by the text
223 prompts. Furthermore, already partially denoised historical segments can serve as stable conditions
224 for generating later segments, aiding long-range consistency without forcing all segments to share the
225 same noise level at the same time.

226 4.4 Flexible Inference Modes for Long Video Generation

227 The temporal flexibility gained from Diffusion Forcing during training allows for various inference
228 methods to generate videos much longer than the training segments (the "*train short, test long*"
229 approach).

230 **Sequential Sliding Window Approaches.** Common methods for long video generation use a
231 sequential sliding window. These include simple autoregressive techniques, where a new segment
232 of M latents is generated based on $N - M$ previous latents from an N -latent window (often re-
233 noising the context), and more advanced methods like FIFO-Diffusion [16], which uses a queue with
234 diagonally progressing noise levels for better temporal consistency. However, a core problem with all
235 such step-by-step sequential methods is the unavoidable build-up of errors, which reduces quality
236 and long-range consistency in very long videos.

237 **Parallel Multi-Window Denoising (PMWD).** To effectively overcome this error accumulation
238 problem, we introduce PMWD. This novel inference strategy takes full advantage of our frame-level
239 prompt system to generate long videos more as a complete whole, rather than piece by piece. For a
240 target long video of L latents, we view it as K overlapping windows, each the length of a training
241 segment. Crucially, all K windows are processed *in parallel* (at the same time) during each step
242 of the diffusion denoising process. Every latent, whether new or historical context, is guided by its
243 own dedicated frame-level prompt. This parallel, parallel method for the entire sequence inherently
244 avoids the cascading error build-up seen in typical autoregressive techniques. Latents located in the
245 overlapping regions between adjacent windows are averaged after each denoising step. This averaging,
246 along with the parallel processing, allows information to flow in both directions (bidirectionally)
247 between an earlier and a later window. Unlike methods where only the past influences the future,
248 PMWD allows upcoming video segments to help refine earlier ones. This is especially useful for
249 creating natural-looking scene changes and maintaining consistency in stories with multiple scenes.

250 5 Experimental Results

251 5.1 Experimental Setup

252 We fully fine-tune the open-source WanX-2.1-T2V-1.3B model with Diffusion Forcing technique
253 on resolution 81x480x832 for 100,000 iterations using our internal dataset (detailed in Section 3) of
254 dense frame-level annotations. Training occurred on H-series GPUs with a global batch size of 64.

255 5.2 Evaluation Dataset

256 We evaluate the model’s capability to generate complex videos by utilizing prompts from the VBench
257 2.0 benchmark, specifically focusing on the **Complex Plots** and **Complex Landscapes**.

258 **Complex Plots** assess the model’s ability to construct coherent and consistent multi-scene narratives
259 based on prompts describing multi-stage events. These prompts often involve extended descriptions

260 (150+ words) outlining a sequence of actions or a story with multiple acts, challenging the model to
 261 maintain plot consistency and logical flow throughout the generated video.

262 **Complex Landscapes** evaluate whether the model can faithfully translate long-form landscape
 263 descriptions (150+ words) into video, including multiple scene transitions dictated by camera move-
 264 ments. These prompts test the model’s understanding of spatial relationships and its ability to render
 265 dynamic changes in the environment as described in the text.

266 5.3 Evaluation Metrics

267 We evaluate video quality using metrics for overall video-text alignment and also propose a new
 268 metric for the issue of semantic confusion in multi scenes generation. Let P_g be the global prompt,
 269 V the generated video, $\{P_1, \dots, P_F\}$ the sequence of F frame-level prompts, and $\{V_1, \dots, V_F\}$ the
 270 corresponding sequence of generated frames.

271 **Standard VBench Evaluation.** To provide a comprehensive assessment of fundamental video
 272 quality aspects, particularly for the complex scenarios presented by our chosen VBench 2.0 prompt
 273 categories (Complex Plots and Complex Landscapes), we incorporate a curated subset of established
 274 metrics from the VBench benchmark. This evaluation focuses on key indicators such as: aesthetic
 275 quality, image quality, and motion smoothness. These selected metrics offer standardized measures
 276 of the perceptual quality and spatio-temporal coherence of the generated videos.

277 **Video-Level Video-Text Similarity.** This standard metric evaluates overall coherence between P_g
 278 and V . $\Phi_V(V)$ represents overall video features (uniformly sample 8 frames as input of ViClip).

$$\mathcal{S}_{\text{global}} = \text{Sim}(\Phi_T(P_g), \Phi_V(V)) \quad (2)$$

279 where we use a pre-trained vision-language model (e.g., ViCLIP) for text embeddings $\Phi_T(\cdot)$ and
 280 video/latent ‘frame’ embeddings $\Phi_V(\cdot)$, with $\text{Sim}(\cdot, \cdot)$ denoting cosine similarity.

281 **Confusion Degree (CD).** Despite the widespread use of $\mathcal{S}_{\text{global}}$, this global metric may assign
 282 favorable scores even when content from different scenes are inappropriately combined. To pinpoint
 283 such temporal and semantic inaccuracies, we introduce the Confusion Degree (CD). A high CD score
 284 reveals a model’s difficulty in maintaining a clear, sequential narrative, often resulting in a muddled
 285 or incoherent visual story. We first define two fundamental frame-level similarity metrics as follows:

$$\begin{aligned} S_{TT}(P_i, P_j) &= \text{Sim}(\Phi_T(P_i), \Phi_T(P_j)) \\ S_{TF}(P_i, V_j) &= \text{Sim}(\Phi_T(P_i), \Phi_V(V_j)) \end{aligned} \quad (3)$$

286 , where $S_{TT}(P_i, P_j)$ represents **frame-level text-text similarity** and $S_{TF}(P_i, V_j)$ represents **frame-**
 287 **level text-frame similarity**. Then $\tilde{S}_{TT}(P_i, P_j) = S_{TT}(P_i, P_j)/S_{TT}(P_i, P_i)$ and $\tilde{S}_{TF}(P_i, V_j) =$
 288 $S_{TF}(P_i, V_j)/S_{TF}(P_i, V_i)$ are applied as normalization function to ensure $\tilde{S}_{TT}(P_i, P_i) = 1$ and
 289 $\tilde{S}_{TF}(P_i, V_i) = 1$.

290 The confusion degree of a text P_i in the generated video V is defined as:

$$\text{CD}(P_i) = \sum_{j \in \{1, \dots, F\}} \max(0, \tilde{S}_{TF}(P_i, V_j) - \tilde{S}_{TT}(P_i, P_j)) \quad (4)$$

291 where $\tilde{S}_{TF}(P_i, V_j) - \tilde{S}_{TT}(P_i, P_j)$ indicates that the content of P_i is more aligned with frame V_j
 292 than its inherent semantic relationship with P_j would suggest, thereby signaling confusion. Then the
 293 confusion of a video V is defined as

$$\text{CD} = \frac{1}{F} \sum_{i=1}^F \text{CD}_i \quad (5)$$

294 , representing the average confusion degree across all frames. Lower CD values indicate superior
 295 narrative consistency and reduced semantic confusion throughout the video.

296 5.4 Comparison and Discussion

297 **Analysis of Video Generation under Complex Prompts.** Tab. 1 provides a comparative analysis of
 298 models trained and inferred using either global video-level or granular frame-level prompts. When

Method	Video Length	Prompts Type	Confusion Degree↓	Video-level Text-Video Consistency↑	Frame-level Text-Video Consistency↑	Motion Smoothness↑	Aesthetic Quality↑	Image Quality↑
DF + Video-level Prompt	5s	Complex Plot	0.2952 ± 0.0461	0.2100 ± 0.0410	0.1635 ± 0.0282	98.43	59.20	67.92
DF + Video-level Prompt	30s	Complex Plot	0.2962 ± 0.0487	0.2053 ± 0.0368	0.1518 ± 0.0258	98.63	52.02	58.07
DF + Frame-level Prompt	30s	Complex Plot	0.1385 ± 0.0498	0.2196 ± 0.0309	0.2054 ± 0.0231	98.53	55.04	61.56
DF + Video-level Prompt	5s	Complex Landscape	0.2745 ± 0.0412	0.2101 ± 0.0341	0.1831 ± 0.0227	98.70	61.32	59.61
DF + Video-level Prompt	30s	Complex Landscape	0.2806 ± 0.0474	0.2066 ± 0.0351	0.1723 ± 0.0230	98.58	52.63	51.02
DF + Frame-level Prompt	30s	Complex Landscape	0.1528 ± 0.0479	0.2195 ± 0.0326	0.2139 ± 0.0167	98.99	56.31	55.98

Table 1: Comparing video-level versus frame-level prompting for complex narrative videos. While global Video-Level Text-Video Consistency can yield misleadingly high scores despite internal scene blending or semantic confusion, metrics like Confusion Degree and frame-level consistency more effectively expose these flaws, highlighting the superior prompt adherence of frame-level strategies. Further analysis in Section 5.4.

Method	Prompt Level	Confusion Degree ↓	Video-level ↑ Text-Video Consistency	Frame-level ↑ Text-Video Consistency	Motion ↑ Smoothness	Aesthetic ↑ Quality	Image ↑ Quality
First-In-First-Out (FIFO)	Video	0.2962 ± 0.0487	0.2053 ± 0.0368	0.1518 ± 0.0258	98.63	52.02	58.07
First-In-First-Out (FIFO)	Frame	0.2416 ± 0.0514	0.2100 ± 0.0370	0.1660 ± 0.0266	98.81	51.32	59.25
Sequential Sliding Window	Frame	0.1773 ± 0.0550	0.2134 ± 0.0312	0.1842 ± 0.0227	98.80	52.04	60.11
Parallel Multi-Window Denoising	Frame	0.1385 ± 0.0498	0.2196 ± 0.0309	0.2054 ± 0.0231	98.53	55.04	61.56

Table 2: Inference method comparison for 30s complex plot videos. Parallel Multi-Window Denoising (PMWD) achieves lower error accumulation (improved aesthetic/image quality) and better prompt adherence (reduced Confusion Degree, higher text-video consistency) versus causal methods (FIFO, Sliding Window). Detailed analysis in Section 5.4.

299 generating short videos (e.g., 5 seconds) conditioned on a single **video-level prompt**, the model
300 operates closer to an ideal scenario without temporal error accumulation. However, such prompts
301 often lead to a high Confusion Degree (CD), as the model struggles to render extensive semantic
302 information within a condensed timeframe, resulting in blended or muddled content.

303 Conversely, employing **frame-level prompts** demonstrates a marked improvement in prompt adher-
304 ence, evidenced by lower CD scores alongside high frame-level consistency metrics. This enhanced
305 ability to follow detailed, segmented instructions makes the frame-level prompting strategy more
306 reliable and effective for generating coherent multi-scene long videos. Furthermore, metrics such
307 as aesthetic and image quality serve as indirect indicators of error accumulation; significant degra-
308 dation in these scores over time typically reflects compounding errors. This accumulation is an
309 inherent consequence of the causal nature of sequential generation processes, a fundamental issue
310 that even precise frame-level semantic guidance cannot resolve on its own when operating within
311 such autoregressive frameworks.

312 **Comparative Analysis of Long Video Inference Strategies.** We further analyze the efficacy of
313 different inference strategies for extending video generation beyond training lengths, comparing
314 our proposed Parallel Multi-Window Denoising (PMWD) with established sequential methods like
315 FIFO-Diffusion and naive sliding windows.

316 Sequential approaches, by their nature, tackle long video generation segment by segment. Naive
317 sliding window techniques autoregressively generate a new chunk of latents conditioned on a lim-
318 ited history of prior latents (often re-noised to manage error). FIFO-Diffusion [16] offers a more
319 sophisticated sequential mechanism, processing a queue of latents with diagonally increasing noise
320 levels to output one clean latent per step, thereby aiming for better temporal consistency through
321 extended context. While these methods incorporate mechanisms to manage error, such as FIFO’s
322 broader context or the re-noising of historical data in naive sliding windows, they fundamentally
323 struggle with the *inevitable accumulation of errors* over very long sequences. This compounding
324 error degrades long-range coherence and overall video quality.

325 Our proposed **Parallel Multi-Window Denoising (PMWD)** is architecturally designed to overcome
326 this critical limitation. Instead of sequential generation, PMWD processes the entire target long video
327 (composed of multiple overlapping windows) *simultaneously* at each denoising step, with each latent
328 guided by its specific frame-level prompt. This parallel, holistic approach fundamentally disrupts
329 the chain of error propagation seen in sequential methods. The averaging of latents in overlapping
330 regions is a key aspect of PMWD. This, combined with parallel processing, not only fuses information
331 effectively but also transforms the strictly causal dependency of sequential models into a *bidirectional*
332 *contextual influence*, where information from temporally subsequent windows can refine earlier ones.
333 This capability is particularly advantageous for rendering naturalistic scene transitions and ensuring
334 global narrative consistency.



Figure 3: Complex Plot Generation. This figure illustrates the impact of different prompting strategies on visual storytelling performance in a complex narrative task based on The Ugly Duckling. The first row shows results generated using DF with a single global prompt, while the second row presents results from our proposed method that combines DF with multiple tailored prompts (multi-prompting). Our method demonstrates significantly improved coherence, reduced error accumulation, and less narrative confusion across the sequence. The images generated with multi-prompting maintain better stylistic and semantic consistency, showcasing its superiority over the global prompt approach in handling complex plot developments.

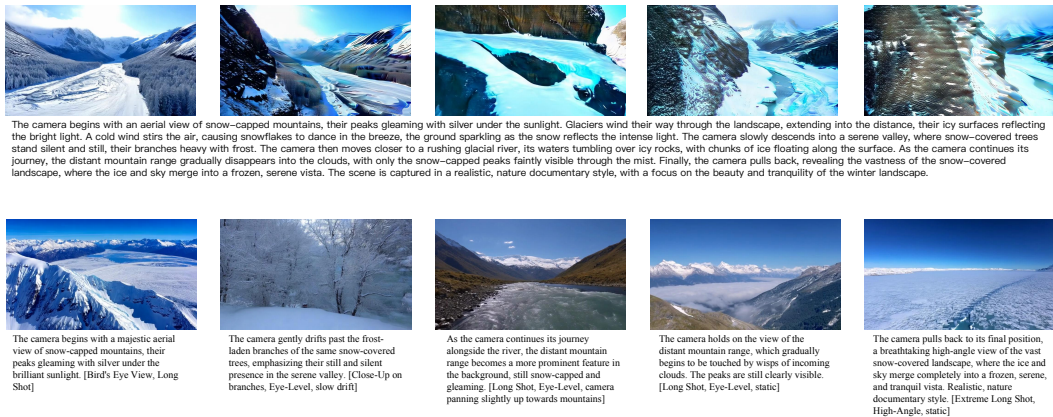


Figure 4: Complex Landscapes Generation. This figure compares two prompting strategies for generating complex scenes. The top row uses DF + global prompt, while the bottom row shows results from our method: DF + multi-prompt. Our approach significantly reduces content drift and error accumulation across frames. By using multiple prompts tailored to each scene segment, it achieves higher accuracy and coherence, capturing the complexity and progression of the winter landscape more effectively than the global prompt method.

335 6 Conclusion

336 Generating long, narratively complex videos with high fidelity remains challenging, primarily due
 337 to issues with coarse semantic guidance and the error accumulation inherent in common sequential
 338 generation techniques. We propose a comprehensive solution combining fine-grained frame-level
 339 annotations, novel training strategies, and a Parallel Multi-Window Denoising (PMWD) inference
 340 method. Our experiments on demanding VBench 2.0 benchmarks demonstrate that this integrated
 341 system significantly improves prompt adherence for complex, multi-scene narratives in ultra-long
 342 videos, achieving high-quality results with minimal error accumulation.

343 References

- 344 [1] Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat,
345 Junhwa Hur, Guanghui Liu, Amit Raj, et al. Lumiere: A space-time diffusion model for video
346 generation. In *SIGGRAPH*, pages 1–11, 2024.
- 347 [2] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Do-
348 minik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion:
349 Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 350 [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja
351 Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent
352 diffusion models. In *CVPR*, pages 22563–22575, 2023.
- 353 [4] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr,
354 Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators.
355 *OpenAI Blog*, 1:8, 2024.
- 356 [5] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet:
357 A large-scale video benchmark for human activity understanding. In *CVPR*, pages 961–970,
358 2015.
- 359 [6] Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent
360 Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in*
361 *Neural Information Processing Systems*, 37:24081–24125, 2024.
- 362 [7] Guibin Chen, Dixuan Lin, Jiangping Yang, Chunze Lin, Juncheng Zhu, Mingyuan Fan, Hao
363 Zhang, Sheng Chen, Zheng Chen, Chengchen Ma, et al. Skyreels-v2: Infinite-length film
364 generative model. *arXiv preprint arXiv:2504.13074*, 2025.
- 365 [8] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao,
366 Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m:
367 Captioning 70m videos with multiple cross-modality teachers. In *CVPR*, pages 13320–13331,
368 2024.
- 369 [9] Mengyu Chu, You Xie, Jonas Mayer, Laura Leal-Taixé, and Nils Thuerey. Learning temporal
370 coherence via self-supervision for gan-based video generation. *TOG*, 39(4):75–1, 2020.
- 371 [10] Rohit Girdhar, Mannat Singh, Andrew Brown, Quentin Duval, Samaneh Azadi, Sai Saketh
372 Rambhatla, Akbar Shah, Xi Yin, Devi Parikh, and Ishan Misra. Factorizing text-to-video
373 generation by explicit image conditioning. In *ECCV*, pages 205–224. Springer, 2024.
- 374 [11] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh
375 Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image
376 diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.
- 377 [12] Yuwei Guo, Ceyuan Yang, Ziyang Yang, Zhibei Ma, Zhijie Lin, Zhenheng Yang, Dahua Lin, and
378 Lu Jiang. Long context tuning for video generation. *arXiv preprint arXiv:2503.10589*, 2025.
- 379 [13] Roberto Henschel, Levon Khachatryan, Daniil Hayrapetyan, Hayk Poghosyan, Vahram Tade-
380 vosyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Streamingt2v: Consistent,
381 dynamic, and extendable long video generation from text. *arXiv preprint arXiv:2403.14773*,
382 2024.
- 383 [14] Emiel Hoogeboom, Alexey A Gritsenko, Jasmijn Bastings, Ben Poole, Rianne van den Berg,
384 and Tim Salimans. Autoregressive diffusion models. In *International Conference on Learning*
385 *Representations*, 2023.
- 386 [15] Haoyang Huang, Guoqing Ma, Nan Duan, Xing Chen, Changyi Wan, Ranchen Ming, Tianyu
387 Wang, Bo Wang, Zhiying Lu, Aojie Li, et al. Step-video-ti2v technical report: A state-of-the-art
388 text-driven image-to-video generation model. *arXiv preprint arXiv:2503.11251*, 2025.
- 389 [16] Jihwan Kim, Junoh Kang, Jinyoung Choi, and Bohyung Han. Fifo-diffusion: Generating infinite
390 videos from text without training, 2024. URL <https://arxiv.org/abs/2405.11473>.

- 391 [17] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin
392 Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video
393 generative models. *arXiv preprint arXiv:2412.03603*, 2024.
- 394 [18] Chengxuan Li, Di Huang, Zeyu Lu, Yang Xiao, Qingqi Pei, and Lei Bai. A survey on long
395 video generation: Challenges, methods, and prospects, 2024. URL [https://arxiv.org/
396 abs/2403.16407](https://arxiv.org/abs/2403.16407).
- 397 [19] Shuang Li, Yihuai Gao, Dorsa Sadigh, and Shuran Song. Unified video action model, 2025.
398 URL <https://arxiv.org/abs/2503.00200>.
- 399 [20] Bin Lin, Yunyang Ge, Xinhua Cheng, Zongjian Li, Bin Zhu, Shaodong Wang, Xianyi He, Yang
400 Ye, Shenghai Yuan, Liuhan Chen, et al. Open-sora plan: Open-source large video generation
401 model. *arXiv preprint arXiv:2412.00131*, 2024.
- 402 [21] Fuchen Long, Zhaofan Qiu, Ting Yao, and Tao Mei. Videostudio: Generating consistent-content
403 and multi-scene videos. In *ECCV*, pages 468–485. Springer, 2024.
- 404 [22] Guoqing Ma, Haoyang Huang, Kun Yan, Liangyu Chen, Nan Duan, Shengming Yin, Changyi
405 Wan, Ranchen Ming, Xiaoni Song, Xing Chen, Yu Zhou, Deshan Sun, Deyu Zhou, Jian Zhou,
406 et al. Step-video-t2v technical report: The practice, challenges, and future of video foundation
407 model, 2025. URL <https://arxiv.org/abs/2502.10248>.
- 408 [23] Guoqing Ma, Haoyang Huang, Kun Yan, Liangyu Chen, Nan Duan, Shengming Yin, Changyi
409 Wan, Ranchen Ming, Xiaoni Song, Xing Chen, et al. Step-video-t2v technical report: The
410 practice, challenges, and future of video foundation model. *arXiv preprint arXiv:2502.10248*,
411 2025.
- 412 [24] Gyeongrok Oh, Jaehwan Jeong, Sieun Kim, Wonmin Byeon, Jinkyu Kim, Sungwoong Kim,
413 and Sangpil Kim. Mevg: Multi-event video generation with text-to-video models. In *ECCV*,
414 pages 401–418. Springer, 2024.
- 415 [25] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, pages
416 4195–4205, 2023.
- 417 [26] Tianhao Qi, Jianlong Yuan, Wanquan Feng, Shancheng Fang, Jiawei Liu, Siyu Zhou, Qian He,
418 Hongtao Xie, and Yongdong Zhang. Mask2dit: Dual mask-based diffusion transformer for
419 multi-scene long video generation. *arXiv preprint arXiv:2503.19881*, 2025.
- 420 [27] Haonan Qiu, Menghan Xia, Yong Zhang, Yingqing He, Xintao Wang, Ying Shan, and Ziwei
421 Liu. Freenoise: Tuning-free longer video diffusion via noise rescheduling. *ICLR*, 2023.
- 422 [28] Team Seaweed, Ceyuan Yang, Zhijie Lin, Yang Zhao, Shanchuan Lin, Zhibei Ma, Haoyuan Guo,
423 Hao Chen, Lu Qi, Sen Wang, et al. Seaweed-7b: Cost-effective training of video generation
424 foundation model. *arXiv preprint arXiv:2504.08685*, 2025.
- 425 [29] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video
426 generator with the price, image quality and perks of stylegan2. In *CVPR*, pages 3626–3636,
427 2022.
- 428 [30] Gabriel Tseng, Ruben Cartuyvels, Ivan Zvonkov, Mirali Purohit, David Rolnick, and Hannah
429 Kerner. Lightweight, pre-trained transformers for remote sensing timeseries. *arXiv preprint
430 arXiv:2304.14065*, 2023.
- 431 [31] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing
432 motion and content for video generation. In *CVPR*, pages 1526–1535, 2018.
- 433 [32] Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao,
434 Jianxiao Yang, Jianyuan Zeng, et al. Wan: Open and advanced large-scale video generative
435 models. *arXiv preprint arXiv:2503.20314*, 2025.
- 436 [33] Bo Wang, Haoyang Huang, Zhiyin Lu, Fengyuan Liu, Guoqing Ma, Jianlong Yuan, Yuan Zhang,
437 and Nan Duan. Storyanchors: Generating consistent multi-scene story frames for long-form
438 narratives, 2025. URL <https://arxiv.org/abs/2505.08350>.

- 439 [34] Fu-Yun Wang, Wenshuo Chen, Guanglu Song, Han-Jia Ye, Yu Liu, and Hongsheng Li.
440 Gen-l-video: Multi-text to long video generation via temporal co-denoising. *arXiv preprint*
441 *arXiv:2305.18264*, 2023.
- 442 [35] Qiheng Wang, Yukai Shi, Jiarong Ou, Rui Chen, Ke Lin, Jiahao Wang, Boyuan Jiang, Haotian
443 Yang, Mingwu Zheng, Xin Tao, et al. Koala-36m: A large-scale video dataset improving con-
444 sistency between fine-grained conditions and video content. *arXiv preprint arXiv:2410.08260*,
445 2024.
- 446 [36] Weihang Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Shiyu Huang, Bin
447 Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Lvbench: An extreme long video understanding
448 benchmark, 2024.
- 449 [37] Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-Levy. Videoagent: Long-form
450 video understanding with large language model as agent, 2024. URL [https://arxiv.org/
451 abs/2403.10517](https://arxiv.org/abs/2403.10517).
- 452 [38] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex:
453 A large-scale, high-quality multilingual dataset for video-and-language research. In *ICCV*,
454 pages 4581–4591, 2019.
- 455 [39] Yaohui Wang, Piotr Bilinski, Francois Bremond, and Antitza Dantcheva. Imaginator: Condi-
456 tional spatio-temporal gan for video generation. In *WACV*, pages 1160–1169, 2020.
- 457 [40] Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu,
458 and Baining Guo. Advancing high-resolution video-language representation with large-scale
459 video transcriptions. In *CVPR*, pages 5036–5045, 2022.
- 460 [41] Xin Yan, Yuxuan Cai, Qiuyue Wang, Yuan Zhou, Wenhao Huang, and Huan Yang. Long video
461 diffusion generation with segmented cross-attention and content-rich video data curation. *arXiv*
462 *preprint arXiv:2412.01316*, 2024.
- 463 [42] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming
464 Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion
465 models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.
- 466 [43] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and
467 Yejin Choi. Merlot: Multimodal neural script knowledge models. *Neurips*, 34:23634–23651,
468 2021.
- 469 [44] Canyu Zhao, Mingyu Liu, Wen Wang, Jianlong Yuan, Hao Chen, Bo Zhang, and Chunhua
470 Shen. Moviedreamer: Hierarchical generation for coherent long visual sequence. *arXiv preprint*
471 *arXiv:2407.16655*, 2024.
- 472 [45] Dian Zheng, Ziqi Huang, Hongbo Liu, Kai Zou, Yinan He, Fan Zhang, Yuanhan Zhang, Jingwen
473 He, Wei-Shi Zheng, Yu Qiao, and Ziwei Liu. Vbench-2.0: Advancing video generation bench-
474 mark suite for intrinsic faithfulness, 2025. URL <https://arxiv.org/abs/2503.21755>.
- 475 [46] Mingzhe Zheng, Yongqi Xu, Haojian Huang, Xuran Ma, Yexin Liu, Wenjie Shu, Yatian Pang,
476 Feilong Tang, Qifeng Chen, Harry Yang, et al. Videogen-of-thought: A collaborative framework
477 for multi-shot video generation. *arXiv preprint arXiv:2412.02259*, 2024.
- 478 [47] Luowei Zhou, Chenliang Xu, and Jason Corso. Towards automatic learning of procedures from
479 web instructional videos. In *AAAI*, volume 32, 2018.
- 480 [48] Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion:
481 Consistent self-attention for long-range image and video generation. *Neurips*, 37:110315–
482 110340, 2024.

483 **NeurIPS Paper Checklist**

484 **1. Claims**

485 Question: Do the main claims made in the abstract and introduction accurately reflect the
486 paper's contributions and scope?

487 Answer: [\[Yes\]](#)

488 Justification: The claims made in the abstract and introduction are aligned with the theoretical
489 and experimental results presented in the paper. The introduction provides a concise
490 summary of what the reader can expect from the paper.

491 Guidelines:

- 492 • The answer NA means that the abstract and introduction do not include the claims
493 made in the paper.
- 494 • The abstract and/or introduction should clearly state the claims made, including the
495 contributions made in the paper and important assumptions and limitations. A No or
496 NA answer to this question will not be perceived well by the reviewers.
- 497 • The claims made should match theoretical and experimental results, and reflect how
498 much the results can be expected to generalize to other settings.
- 499 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
500 are not attained by the paper.

501 **2. Limitations**

502 Question: Does the paper discuss the limitations of the work performed by the authors?

503 Answer: [\[Yes\]](#)

504 Justification: Primarily covered in Section 6.

505 Guidelines:

- 506 • The answer NA means that the paper has no limitation while the answer No means that
507 the paper has limitations, but those are not discussed in the paper.
- 508 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 509 • The paper should point out any strong assumptions and how robust the results are to
510 violations of these assumptions (e.g., independence assumptions, noiseless settings,
511 model well-specification, asymptotic approximations only holding locally). The authors
512 should reflect on how these assumptions might be violated in practice and what the
513 implications would be.
- 514 • The authors should reflect on the scope of the claims made, e.g., if the approach was
515 only tested on a few datasets or with a few runs. In general, empirical results often
516 depend on implicit assumptions, which should be articulated.
- 517 • The authors should reflect on the factors that influence the performance of the approach.
518 For example, a facial recognition algorithm may perform poorly when image resolution
519 is low or images are taken in low lighting. Or a speech-to-text system might not be
520 used reliably to provide closed captions for online lectures because it fails to handle
521 technical jargon.
- 522 • The authors should discuss the computational efficiency of the proposed algorithms
523 and how they scale with dataset size.
- 524 • If applicable, the authors should discuss possible limitations of their approach to
525 address problems of privacy and fairness.
- 526 • While the authors might fear that complete honesty about limitations might be used by
527 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
528 limitations that aren't acknowledged in the paper. The authors should use their best
529 judgment and recognize that individual actions in favor of transparency play an impor-
530 tant role in developing norms that preserve the integrity of the community. Reviewers
531 will be specifically instructed to not penalize honesty concerning limitations.

532 **3. Theory assumptions and proofs**

533 Question: For each theoretical result, does the paper provide the full set of assumptions and
534 a complete (and correct) proof?

535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
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

Answer: [NA]

Justification: This is not a purely theoretical paper.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We explained our framework in Section 4 and experiment settings in Section 5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

587 Question: Does the paper provide open access to the data and code, with sufficient instruc-
588 tions to faithfully reproduce the main experimental results, as described in supplemental
589 material?

590 Answer: [Yes]

591 Justification: We will release code with instructions to reproduce the results.

592 Guidelines:

- 593 • The answer NA means that paper does not include experiments requiring code.
- 594 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
595 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 596 • While we encourage the release of code and data, we understand that this might not be
597 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
598 including code, unless this is central to the contribution (e.g., for a new open-source
599 benchmark).
- 600 • The instructions should contain the exact command and environment needed to run to
601 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
602 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 603 • The authors should provide instructions on data access and preparation, including how
604 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 605 • The authors should provide scripts to reproduce all experimental results for the new
606 proposed method and baselines. If only a subset of experiments are reproducible, they
607 should state which ones are omitted from the script and why.
- 608 • At submission time, to preserve anonymity, the authors should release anonymized
609 versions (if applicable).
- 610 • Providing as much information as possible in supplemental material (appended to the
611 paper) is recommended, but including URLs to data and code is permitted.

612 6. Experimental setting/details

613 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
614 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
615 results?

616 Answer: [Yes]

617 Justification: The paper specifies training and test details in Section 5

618 Guidelines:

- 619 • The answer NA means that the paper does not include experiments.
- 620 • The experimental setting should be presented in the core of the paper to a level of detail
621 that is necessary to appreciate the results and make sense of them.
- 622 • The full details can be provided either with the code, in appendix, or as supplemental
623 material.

624 7. Experiment statistical significance

625 Question: Does the paper report error bars suitably and correctly defined or other appropriate
626 information about the statistical significance of the experiments?

627 Answer: [Yes]

628 Justification: The paper reports means and variances of the experiments. We made sure that
629 for all the experiments conducted throughout the paper, we averaged across multi runs to
630 make sure that the results are reliable.

631 Guidelines:

- 632 • The answer NA means that the paper does not include experiments.
- 633 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
634 dence intervals, or statistical significance tests, at least for the experiments that support
635 the main claims of the paper.
- 636 • The factors of variability that the error bars are capturing should be clearly stated (for
637 example, train/test split, initialization, random drawing of some parameter, or overall
638 run with given experimental conditions).

- 639 • The method for calculating the error bars should be explained (closed form formula,
640 call to a library function, bootstrap, etc.)
- 641 • The assumptions made should be given (e.g., Normally distributed errors).
- 642 • It should be clear whether the error bar is the standard deviation or the standard error
643 of the mean.
- 644 • It is OK to report 1-sigma error bars, but one should state it. The authors should
645 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
646 of Normality of errors is not verified.
- 647 • For asymmetric distributions, the authors should be careful not to show in tables or
648 figures symmetric error bars that would yield results that are out of range (e.g. negative
649 error rates).
- 650 • If error bars are reported in tables or plots, The authors should explain in the text how
651 they were calculated and reference the corresponding figures or tables in the text.

652 8. Experiments compute resources

653 Question: For each experiment, does the paper provide sufficient information on the com-
654 puter resources (type of compute workers, memory, time of execution) needed to reproduce
655 the experiments?

656 Answer: [Yes]

657 Justification: We explained our settings in Section 5.

658 Guidelines:

- 659 • The answer NA means that the paper does not include experiments.
- 660 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
661 or cloud provider, including relevant memory and storage.
- 662 • The paper should provide the amount of compute required for each of the individual
663 experimental runs as well as estimate the total compute.
- 664 • The paper should disclose whether the full research project required more compute
665 than the experiments reported in the paper (e.g., preliminary or failed experiments that
666 didn't make it into the paper).

667 9. Code of ethics

668 Question: Does the research conducted in the paper conform, in every respect, with the
669 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

670 Answer: [Yes]

671 Justification: The paper adheres to the ethical guidelines set forth by NeurIPS. We ensured
672 that the research is conducted responsibly, with considerations for potential biases, fairness,
673 and the broader impact on society.

674 Guidelines:

- 675 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 676 • If the authors answer No, they should explain the special circumstances that require a
677 deviation from the Code of Ethics.
- 678 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
679 eration due to laws or regulations in their jurisdiction).

680 10. Broader impacts

681 Question: Does the paper discuss both potential positive societal impacts and negative
682 societal impacts of the work performed?

683 Answer: [Yes]

684 Justification: The paper discusses both potential positive societal impacts and negative
685 societal impacts of the work performed in Section 6.

686 Guidelines:

- 687 • The answer NA means that there is no societal impact of the work performed.
- 688 • If the authors answer NA or No, they should explain why their work has no societal
689 impact or why the paper does not address societal impact.

- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- 708
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

709 11. Safeguards

710 Question: Does the paper describe safeguards that have been put in place for responsible
711 release of data or models that have a high risk for misuse (e.g., pretrained language models,
712 image generators, or scraped datasets)?

713 Answer: [NA]

714 Justification: The work doesn't pose any risks. We either generated our own data or used a
715 justified reliable benchmark.

716 Guidelines:

- 717
- 718
- 719
- 720
- 721
- 722
- 723
- 724
- 725
- 726
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

727 12. Licenses for existing assets

728 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
729 the paper, properly credited and are the license and terms of use explicitly mentioned and
730 properly respected?

731 Answer: [Yes]

732 Justification: We credited the author for the code package and benchmark dataset that have
733 been used in the paper.

734 Guidelines:

- 735
- 736
- 737
- 738
- 739
- 740
- 741
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 742
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 743
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 744
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 745
- 746
- 747
- 748
- 749

750 13. **New assets**

751 Question: Are new assets introduced in the paper well documented and is the documentation
752 provided alongside the assets?

753 Answer: [Yes]

754 Justification: We will release code with detailed documentation and an appropriate license.

755 Guidelines:

- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 756
- 757
- 758
- 759
- 760
- 761
- 762
- 763

764 14. **Crowdsourcing and research with human subjects**

765 Question: For crowdsourcing experiments and research with human subjects, does the paper
766 include the full text of instructions given to participants and screenshots, if applicable, as
767 well as details about compensation (if any)?

768 Answer: [NA]

769 Justification: The paper does not involve crowdsourcing nor research with human subjects.

770 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 771
- 772
- 773
- 774
- 775
- 776
- 777
- 778

779 15. **Institutional review board (IRB) approvals or equivalent for research with human 780 subjects**

781 Question: Does the paper describe potential risks incurred by study participants, whether
782 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
783 approvals (or an equivalent approval/review based on the requirements of your country or
784 institution) were obtained?

785 Answer: [NA]

786 Justification: We do not involve crowdsourcing nor research with human subjects.

787 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- 788
- 789
- 790
- 791
- 792

793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: During inference, we employ LLMs to transform video-level prompts into frame-level prompts. During training, we similarly use the LLM to assist in constructing frame-level annotations, including generating shared or independent captions based on visual changes. The detailed LLM usage in both stages is described in the Section 3.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.