
Spiral Evolution of Visual World Model: Reclaiming Autoregression from the Diffusion Era

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent advances in video generation have been dominated by diffusion-based
2 models, which produce high-quality, prompt-faithful sequences through holistic
3 denoising. While this paradigm has achieved striking visual fidelity, it falls short
4 for real-time, interactive applications that require frame-level responsiveness and
5 causal coherence—cornerstones of practical world modeling. In this position paper,
6 we advocate for a strategic return to autoregressive generation as the foundational
7 architecture for building interactive world simulators. We argue that beyond of-
8 fering faster inference, autoregressive models bring critical structural advantages:
9 they naturally support predictive compression, enable causal disentanglement, and
10 offer a more responsive mechanism for integrating control signals in dynamic
11 settings. Unlike language-conditioned diffusion models, autoregression flexibly
12 accommodates frame-wise control inputs such as camera motion and joint actions,
13 making it ideally suited for agent-centric simulation. We further highlight emerging
14 techniques and promising directions—including selective denoising, adaptive reso-
15 lution, and postdictive coding—that address historical limitations of autoregression
16 and unlock new levels of interactivity. We contend that embracing autoregression
17 will be essential for developing practical, controllable, and truly intelligent world
18 models.

19 1 Introduction

20 In recent years, diffusion-based models [1, 2] have taken center stage in video generation, praised
21 for their ability to produce high-fidelity, prompt-aligned sequences through holistic denoising. This
22 paradigm has fueled breakthroughs in creative and cinematic applications, where visual quality and
23 global coherence are key. However, the very design that enables such visual excellence—global opti-
24 mization across the entire video—also makes diffusion models inherently slow and computationally
25 intensive. These limitations hinder their applicability in real-time, interactive settings.

26 Before the diffusion era, autoregressive models—often built on ConvRNN architectures [3, 4]—were
27 the standard, as illustrated in Fig. 1. These models generated videos frame by frame, embracing a
28 causal structure that naturally supported temporal reasoning and sequential feedback. Despite these
29 strengths, they fell out of favor due to limited visual quality and rigid control mechanisms. Ironically,
30 as the field shifts from offline synthesis to embodied intelligence and interactive simulation—domains
31 where responsiveness, causal coherence, and multimodal control are paramount—the foundational
32 strengths of autoregression are becoming increasingly relevant. The task ahead is not to revert to
33 obsolete architectures but to re-envision autoregressive generation with modern advances, positioning
34 it as the backbone of practical world models.

35 In this position paper, we argue that **video generation—particularly for world modeling—stands**
36 **to benefit significantly from a renewed focus on autoregressive methods.** While diffusion-based

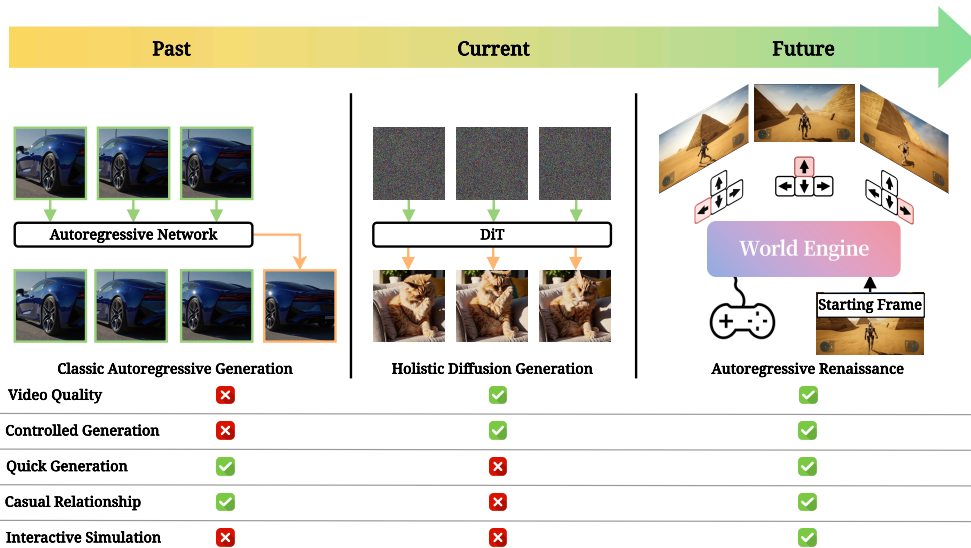


Figure 1: Illustration of the paradigm shift in video generation: we advocate revisiting the autoregressive paradigm as a foundation for building more powerful and interactive world models.

37 models excel at holistic scene synthesis and compositional generalization, they fall short in tasks
 38 demanding real-time feedback and fine-grained control. Autoregressive generation, by contrast,
 39 naturally meets these needs: it integrates frame-level signals, supports diverse conditioning, and
 40 enables continuous interaction. By framing generation as a step-by-step prediction task, it promotes
 41 compact, causally grounded representations that are computationally efficient and well-suited to
 42 embodied agents. Revisiting autoregression is not a regression but a necessary progression—one that
 43 enables the development of practical, controllable, and intelligent world models. We contend that
 44 the future of video generation—and embodied AI more broadly—belongs to models that predict,
 45 compress, and revise the unfolding movie of the world, one actionable frame at a time.

46 2 World models

47 Before we dive into the discussion of video generation, it is worth pausing to examine why an agent
 48 needs a world model in the first place. A clear appreciation of this motivation will illuminate the
 49 position we have taken late.

50 At its core, a world model [5] is a compressed encoding of how the external world evolves in space and
 51 time. By internally rehearsing possible futures—compressing past experience while predicting what
 52 comes next—an agent can uncover the underlying physics and causal structure of its environment.
 53 This rehearsal acts as a dry-run for real exploration, a preview of an era where agents primarily learn
 54 from their own trajectories. Since an agent’s sensorium provides an unlimited stream of training data,
 55 a scalable offline generative approach becomes essential: the model improves through imagination,
 56 not by risking failure in the real world.

57 Moreover, the human brain offers a compelling biological precedent for the necessity of world models
 58 [6, 7]. Our minds constantly engage in a similar "filling-in" mechanism, essentially constructing an
 59 internal model that mirrors the evolution of the external world to grasp its causal dynamics. This
 60 model functions as a continuous prediction engine, ceaselessly forecasting incoming sensory input
 61 and updating itself to minimize discrepancies between expected and actual outcomes. Concepts like
 62 the Bayesian brain and "controlled hallucination" further highlight this predictive and adaptive nature,
 63 where our brains aren’t just passively receiving information but actively generating and refining a
 64 coherent internal narrative of reality through processes akin to dreaming and predictive coding.

65 **3 Paradigm Shift of Video Generation**

66 Over the past decade, video generation and prediction have undergone a clear paradigm shift—from
67 early sequential modeling to today’s full-sequence diffusion frameworks. Along the way, several
68 initial assumptions proved flawed or misguided. Reflecting on these evolving perspectives offers
69 both intellectual insight and essential context, laying the groundwork for the position advanced in the
70 following sections.

71 **3.1 The ConvRNN Era: Predicting Frame by Frame**

72 Research into video generation began alongside the early successes of deep learning [8, 9, 10,
73 11]. Initially, researchers proceeded cautiously, working with simple datasets to test the feasibility
74 of predicting future frames. In 2014, Srivastava et al. [12] first introduced a benchmark called
75 Moving MNIST, derived from the MNIST dataset, and applied recurrent neural networks (RNNs),
76 specifically LSTMs [13], for future frame forecasting. Building on this, Shi et al. [3] later proposed
77 ConvLSTM, an architecture that combined convolutional operations with LSTM units. This allowed
78 the convolutional layers to model spatial relationships while the recurrent layers focused on capturing
79 temporal dynamics. Subsequent research extensively focused on improving this early framework.
80 Some works proposed architectural enhancements to RNNs [4, 14], while others explored lossless
81 compression techniques to compact video representations without sacrificing quality [15]. There
82 were also efforts to incorporate variational inference to better capture the stochastic nature of future
83 events, and methods that separately modeled foreground and background to enhance generation
84 quality [16, 17].

85 Although Moving MNIST may seem trivial today, achieving strong performance on it was still
86 challenging as late as 2020. At the time, models trained with Mean Squared Error (MSE) were known
87 to produce blurry outputs, as MSE tends to average over multiple uncertain futures, merging distinct
88 outcomes into one. Variational autoencoders (VAEs) [18], which compress visual data, were believed
89 to worsen this blurriness. As a result, many video prediction frameworks avoided compression, opting
90 for resolution-preserving designs. Meanwhile, the rise of GANs [19] in image generation heavily
91 influenced video research. Many efforts integrated adversarial loss to reduce MSE-induced blur.
92 While this led to sharper frames, it also introduced artifacts and disrupted temporal consistency. This
93 highlighted a tradeoff: MSE models yielded blurrier but temporally stable videos, while GAN-based
94 ones produced crisp yet often incoherent sequences.

95 **3.2 The Diffusion Era: Denoising Step by Step**

96 Before diffusion models rose to prominence, generative modeling had already begun shifting beyond
97 traditional GANs, particularly in image generation. Two-stage pipelines emerged: a VAE or VQ-VAE
98 [20] first learned to compress images into latent space, then an autoregressive transformer modeled
99 that space with text conditioning. Early successes like VQGAN [21], DALL-E [22], and CogView
100 [23] proved that models could reliably generate images aligned with arbitrary prompts—a major
101 leap forward. Around the same time, diffusion models [24, 25] introduced a different generative
102 approach: learning to reverse a noising process to turn noise into data via iterative denoising. Initially
103 operating in pixel space with U-Net architectures [26], these models achieved impressive quality
104 but at high computational cost. A turning point came with Stable Diffusion [27], which adopted a
105 two-stage latent-based design, greatly improving efficiency and making high-quality image generation
106 accessible on consumer hardware.

107 The natural next step was extending diffusion to video. Early efforts like AnimateDiff [28] and
108 SVD [2] reused pretrained 2D VAEs and augmented U-Nets with temporal modules—typically
109 convolutions or attention layers—to handle motion. These models applied holistic denoising over
110 entire video latents, not frame-by-frame prediction. While promising, they often suffered from
111 temporal artifacts like flickering, revealing the limitations of adapting image-centric architectures
112 to dynamic video data. A major advance came with the Diffusion Transformer (DiT) [29] in Sora
113 [1], which replaced U-Nets with transformers [30] operating on spatio-temporal patches, capturing
114 long-range dependencies and improving motion consistency. Just as crucial was Sora’s use of a
115 3D VAE, which compressed video into spacetime latents, enabling temporal compression and more
116 coherent dynamics. Unlike 2D VAEs, training 3D VAEs requires distinct methods—e.g., using
117 variable-resolution clips to enhance generalization. This methodological shift explains early reliance

118 on frame-wise 2D VAEs and underscores the growing consensus: temporal compression is vital for
119 scalable, coherent video generation.

120 **4 Driving Forces Behind Paradigm Shift**

121 Having outlined the past decade’s paradigm shift of video generation in the previous section, we now
122 turn to an analysis of the key forces that drove this transformation. Understanding these underlying
123 factors will set the stage for the following section, where we will argue for a return to autoregressive
124 generation.

125 **4.1 Why ConvRNNs paradigm fall short in video generation?**

126 The decline of the ConvRNN framework stems from several limitations. First, its task formulation
127 is inherently narrow: video prediction is typically framed as predicting the next K frames from N
128 past ones. This restricts the model to short-term extrapolation based solely on prior motion. For
129 example, in car-mounted videos, if the car is moving forward, the model will simply continue that
130 motion, with no capacity for control. Even if given a command—like turning right—it won’t know
131 what should appear in the new view. Will there be a pedestrian or a tree? Resolving this ambiguity
132 requires additional modalities to set expectations. Without them, the model faces high uncertainty
133 and struggles to generate coherent long-term content. Worse still, ConvRNNs are poorly equipped to
134 incorporate such multimodal signals, further limiting their usefulness.

135 Another key issue is their reliance on autoregressive generation. Since frames are produced one at
136 a time, small artifacts introduced early in the sequence accumulate and intensify, degrading output
137 quality—a problem known as compounding or drifting error.

138 **4.2 The rise of text-to-video generation**

139 The rise of diffusion models in video generation followed naturally from their breakthroughs in
140 image synthesis — especially in text-to-image generation [2, 27], where they achieved high visual
141 fidelity and strong compositionality. Building on this success, early video diffusion work focused on
142 text-to-video (T2V) generation, directly inheriting both architecture and conditioning strategies from
143 image diffusion models. But why has language — especially text — become the default conditioning
144 signal for large-scale video pretraining?

145 Language is uniquely expressive and complete. It can precisely describe any concept, event, or entity,
146 and its compositional nature allows for building complex ideas from simple components. Among all
147 modalities, language alone enables both abstract generalization and fine-grained control [31, 32, 33].
148 In theory, one could even describe every pixel in every frame using natural language, suggesting
149 a theoretical one-to-one mapping between text and video. This expressive completeness makes
150 language a powerful interface for guiding generative models.

151 In video generation, pairing diffusion with text conditioning triggered a paradigm shift. Text naturally
152 captures spatiotemporal events, aligning well with the holistic denoising process of diffusion. Unlike
153 autoregressive models that predict frames sequentially, diffusion generates entire video sequences
154 in latent space, promoting global consistency. Attention mechanisms support dense, bidirectional
155 interactions across space and time, letting the prompt influence every region of every frame —
156 ensuring semantic alignment with the input. As T2V dominates the field, diffusion has emerged as
157 the leading video generation framework.

Full-sequence diffusion models is a natural fit to text-to-video generation.

The shift from ConvRNNs to diffusion models in video generation stems from the rise of text-to-video synthesis, where full-sequence diffusion models excel at generating globally consistent, prompt-aligned videos through holistic denoising and strong attention control.

158

159 **5 The Autoregressive Renaissance**

160

161 While the integration of T2V tasks with diffusion models has delivered impressive results, video
162 generation remains imperfect. A major limitation is generation speed — reflecting the tradeoff
163 between computational cost and video quality. Though less critical for creative applications, this
164 becomes a bottleneck for real-time tasks like world modeling. Various acceleration methods (e.g.,
165 DDIM [34], DPM-Solver [35]) have been explored, but real-time feedback remains a challenge.
166 Meanwhile, conditional control is only beginning to reveal its full potential. Language’s unmatched
167 versatility has, in some cases, overshadowed the value of other modalities. Consider camera motion
168 [36, 37]: it can be conveyed via pose data or as text like “the camera pans left.” Does this diminish
169 the role of structured inputs? From a world modeling lens, clearly not — precise multimodal control
170 is essential for grounded, interactive, temporally coherent simulation.

171 We argue that the next shift in video generation will be a return from holistic diffusion to autoregressive
172 generation. This isn’t just about speed — though autoregressive models are faster — it’s about
173 enabling fine-grained, frame-level control. More importantly, autoregressive generation promotes
174 predictive compression: by learning to generate future frames from past ones, models internalize
175 compact, structured world representations. This predictive structure allows world models to move
176 beyond surface correlations and capture real causal dynamics. While “next-step” need not mean
177 “next-frame,” autoregression reintroduces interpretability, controllability, and a causally grounded
178 approach to understanding the world.

179 5.1 Autoregressive Imagination for Interactive World Modeling

180 The visual world model takes the form of a video generator for one reason: to let agents visually
181 imagine alternative futures — a key prerequisite for planning and interaction. But this imagination
182 must be fast. Asking a model to render a 10-second video over several minutes isn’t just inefficient
183 — it breaks the principle of agency. Real-world environments are dynamic and unpredictable. No
184 intelligent agent, robotic or biological, can afford to wait before acting. Interaction must be real-time
185 and continuous, just like how humans instinctively shift gaze or reach in response to unfolding events.

186 This immediacy demands more than abstract, long-horizon speculation — it requires short-horizon,
187 frame-by-frame anticipation grounded in the present. Language excels at conveying goals or imagined
188 outcomes, but it’s ill-suited for fine-grained, low-latency control. This reveals a key mismatch:
189 text operates at a coarse temporal scale, while decisions unfold at much finer resolutions. Here,
190 autoregressive generation emerges as the natural solution — offering the agility and responsiveness
191 needed for intelligent agents.

192 Crucially, not all conditioning signals are equal — they vary in modality and temporal granularity. A
193 text prompt might describe a five-second scene globally, while control signals like camera pose, joint
194 angles, or keypresses update every frame. Static attributes like background layout or identity must
195 remain consistent throughout. Current frameworks often stack these signals into fixed sequences
196 to align with global prompts, but this sacrifices reactivity and flexibility. Worse, signals across
197 timescales may conflict — e.g., a global prompt suggests panning right, while frame-wise pose
198 indicates motion left. These mismatches confuse the model and destabilize internal representations.
199 Ideally, local, real-time signals should override global ones — yet multi-timescale control remains
200 poorly understood.

201 Autoregressive generation, by contrast, naturally supports temporal hierarchy. It generates videos
202 one frame (or patch) at a time, allowing real-time feedback and causal integration of both high-level
203 intent and low-level updates. Like a game engine, it’s responsive, grounded, and adaptable. This
204 frame-wise structure also supports continual learning — enabling the model to update its internal
205 state on the fly, adapt to new environments, and improve without full retraining. If the goal is an
206 interactive world model — one that supports planning, adaptation, and real-time decision-making —
207 then the future lies not in static, offline diffusion, but in fast, flexible, autoregressive imagination.

208 5.2 Autoregression for Causality Learning

209 In its most straightforward form, autoregressive video modeling defines the generative process as a
210 factorized distribution over frames or spatiotemporal tokens:

$$p(x_{1:T}) = \prod_{t=1}^T p(x_t | x_{<t}),$$

211 where x_t denotes the video content (frame or token) at time t . This formulation forces the model to
 212 predict the future from the past, mirroring the causal flow of time. Unlike holistic denoising methods,
 213 which generate frames jointly and bidirectionally, autoregressive models must isolate which aspects
 214 of the past are predictive of future outcomes. This naturally encourages disentanglement of causal
 215 factors from statistical correlations, since only truly informative features support accurate prediction.

216 This sequential structure aligns with the dynamics of the real world, where observations emerge
 217 through state transitions governed by physical laws and actions [38, 39]. Autoregressive paradigms
 218 are well-suited to capture such structure, especially when paired with control-aware architectures.
 219 For example, modeling the system as a Markov process with hidden states s_t and control inputs a_t ,
 220 the model learns:

$$p(x_t | x_{<t}) \approx p(x_t | s_t), \quad \text{where} \quad s_t = f(s_{t-1}, a_{t-1}),$$

221 revealing that next-step prediction is not merely about extrapolation but about learning compact
 222 state representations that encode the rules of temporal evolution. Recently, a growing body of work
 223 has shown that video diffusion models often lack physical common sense — further reinforcing the
 224 necessity of learning causal structures.

225 5.3 Autoregressive Prediction as Compression

226 To illustrate how autoregressive future prediction functions as a form of compression, let us consider
 227 a large language model (LLM) as a concrete example. [40] (Note that the same line of reasoning can
 228 also be seamlessly extended to video generation.) Suppose Alice intends to transmit a massive dataset
 229 \mathcal{D} of length n . At a given point in time, the first t words, denoted as x_1, x_2, \dots, x_t , have already
 230 been transmitted. Without loss of generality, we assume that the dataset’s dictionary has a size of m .
 231 In the worst-case scenario, each word is one-hot encoded in binary, requiring $\log m$ bits per word.
 232 Thus, the total cost of transmitting the entire dataset \mathcal{D} using this naive encoding method is given by:

$$C_0 = |f_0| + (n - t) \log m, \quad (1)$$

233 where $|f_0|$ represents the constant overhead required for transmission.

234 Certainly, more advanced compression techniques can be employed to reduce the transmission cost.
 235 However, recent research has revealed an intriguing insight: training an LLM specifically to predict
 236 the next token is itself a highly efficient method of data compression.

237 Since we already have the transmitted words x_1, x_2, \dots, x_t , we can leverage the output of the LLM,
 238 which provides a probabilistic distribution over the next token:

$$P(x_{t+1} | x_{1:t}). \quad (2)$$

239 Using arithmetic coding, we can exploit this predictive distribution to reduce the number of bits
 240 required for transmission. The theoretical upper bound on the number of bits required for encoding a
 241 token is given by:

$$-\log P(x_{t+1} = x_{t+1}^* | x_{1:t}) + 1. \quad (3)$$

242 Thus, the total cost of transmitting the dataset under this predictive coding scheme becomes:

$$C_1 = |f_1| + (n - t) + \sum_{i=t}^n -\log P(x_{i+1} = x_{i+1}^* | x_{1:i}). \quad (4)$$

243 Here, $|f_1| + (n - t)$ represents a constant term related to the transmission overhead and the model’s
 244 internal encoding scheme. The main term that determines the efficiency of compression is:

$$\sum_{i=t}^n -\log P(x_{i+1} = x_{i+1}^* | x_{1:i}), \quad (5)$$

245 This depends on the model’s predictive capability. While no theoretical bound can be derived, it can
 246 be estimated empirically. Hoffmann’s experiments showed that a well-trained large language model
 247 (LLM) achieved a 14-fold data compression rate. Moreover, larger models yield higher compression,
 248 suggesting they better capture statistical regularities. In contrast, the best text compression algorithm
 249 from the Hutter Prize achieves only an 8.7-fold ratio, highlighting the efficiency of predictive coding
 250 in neural networks.

251 In sum, autoregressive modeling thus serves as both a predictive and compressive mechanism,
 252 encouraging the world model to internalize causal structures rather than overfit to surface patterns.

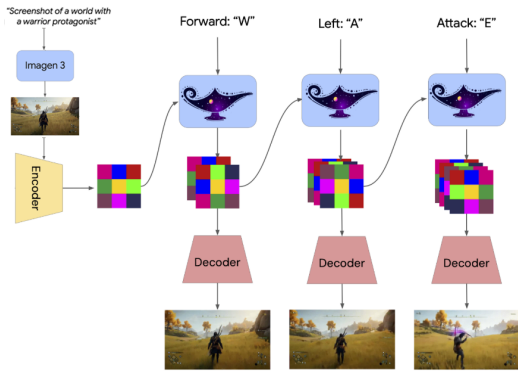


Figure 2: Autoregressive roll-out of Genie2

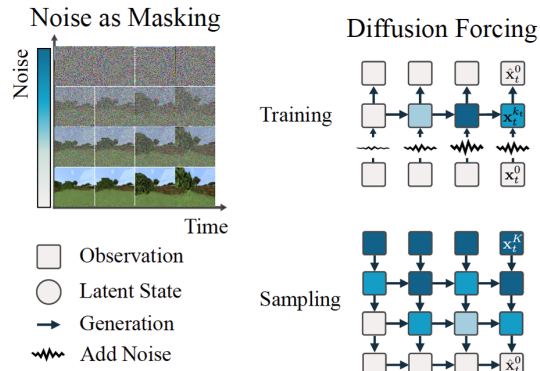


Figure 3: Diffusion Forcing denoising strategy

253 **5.4 Escaping Language’s Babysitting**

254 Finally, we propose a speculative yet compelling question: in the relationship between language and
 255 visual world models, which truly comes first? While language can resolve much of the uncertainty in
 256 video generation, it may be better seen as a byproduct of our structured understanding of 4D reality.
 257 That is, cognition of space and time likely grounds language—not vice versa. [41]

258 Though language can be learned autoregressively without visual input, grounding it in visual experi-
 259 ence may greatly improve sample efficiency. [42] This raises the possibility that the Text-to-Video
 260 paradigm may be somewhat inverted—we use language to guide video generation, rather than letting
 261 vision structure language. While language and vision likely co-evolved, overreliance on linguistic
 262 prompts may prevent video models from learning stable, intrinsic representations of the physical
 263 world.

264 Some argue that language serves as scaffolding—useful for structuring early learning. [43] But like
 265 scaffolding removed once a structure is complete, language may become less central. A mature world
 266 model may rely less on instruction and more on direct interaction—learning through perception,
 267 action, and embodied experience. Future paradigms should move beyond curated language data
 268 toward models grounded in raw experience.

Autoregressive generative models: the bedrock of the interactive world model.

Autoregressive video generation enables real-time, frame-level control and predictive com-
 pression, fostering causally grounded, interpretable world models.

270 **6 Next-Gen Autoregressive Video Models**

271 Having established the case for autoregressive generation as a foundation for world modeling, the
 272 next step is to explore how this shift can be practically realized. In practice, the reinforcement
 273 learning community, driven by a growing belief in the importance of world models, has already begun
 274 revisiting autoregressive video generation.

275 This section will therefore provide a concise overview and categorization of these nascent efforts,
 276 alongside proposing additional critical directions for future research. A central challenge in these
 277 early explorations stems from the quadratic computational cost of transformers — the basis for most
 278 state-of-the-art video generation models — as their cost increases with input length. As a result,
 279 much of the current research converges on a central question: how can we shorten latent sequences —
 280 or denoise only a small subset at each step — without compromising fidelity or control?

281 **6.1 Next-frame Prediction Using diffusion**

282 A simple yet effective way to implement autoregressive video generation is by adapting diffusion
 283 models to predict frames sequentially, generating videos one frame at a time using past latents and

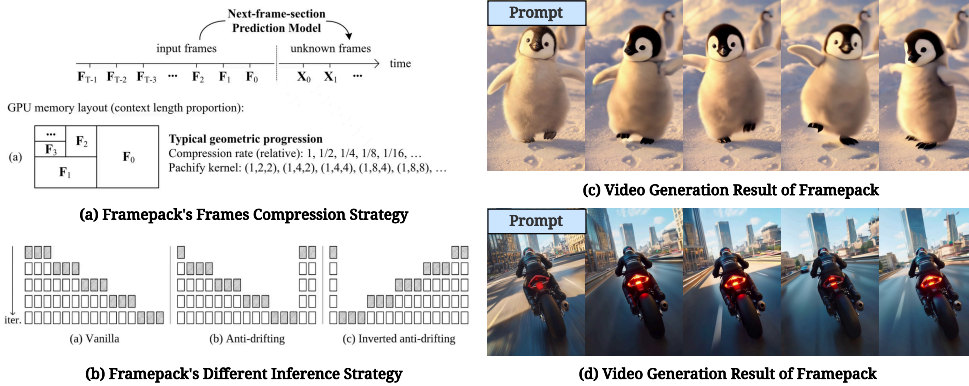


Figure 4: Illustration of Framepack: (a) One of the Framepack’s past frames compression schemes. (b) Framepack’s inference strategy, comprising causal autoregressive generation and two anti-drifting methods. (c) and (d) Example videos generated by Framepack.

284 actions. Genie 2 [44] follows this approach with a latent diffusion model combining an autoencoder
 285 and causal transformer to capture temporal dynamics, enabling framewise generation. Navigation
 286 World Model (NWM) [45] extends this to agent-centric navigation, using a Conditional Diffusion
 287 Transformer (CDiT) to simulate future egocentric views based on past inputs, modeling complex
 288 3D transitions. GameFactory [46] applies this to games, injecting keyboard and mouse actions into
 289 pretrained diffusion models and supporting long-horizon generation via continuous conditioning
 290 and variable noise schedules. Together, these models show how diffusion can be repurposed for
 291 autoregressive, temporally coherent, interactive video synthesis.

292 6.2 Selective Denoising

293 Selective denoising—where video segments are denoised incrementally rather than in a single
 294 global step—has emerged as a scalable, flexible alternative to traditional diffusion-based generation.
 295 Diffusion Forcing (DF) [47] exemplifies this by treating generation as partial unmasking: each token
 296 is noised independently, and the model learns to denoise arbitrary subsets using a shared next-token
 297 prediction setup. This enables flexible sequence lengths, adaptive schedules, and compositional
 298 generalization. Oasis [48] applies DF to complex settings like Minecraft, generating coherent
 299 frames conditioned on evolving states. Similarly, MAGI-1 [49] uses chunk-wise autoregression
 300 with diffusion transformers for temporally consistent, stepwise video generation. These methods
 301 showcase adaptive denoising’s benefits—greater temporal flexibility, lower compute, and enhanced
 302 interactivity—making it well-suited for real-time applications.

303 6.3 Adaptive Resolution

304 Another valid strategy to improve autoregressive video generation is to enable adaptive resolution.
 305 FramePack [50] achieves this by compressing past frames into a fixed-size memory, preserving only
 306 the most salient information while discarding redundant details. This allows autoregressive models to
 307 maintain constant compute regardless of sequence length, effectively acting as a temporal resolution
 308 filter. Framepack’s compression and inference strategy and generated samples can be seen in Fig. 4.
 309 If extended to video, models like VAR [51] could adopt similar strategies—using resolution-aware
 310 tokens or dynamic memory to scale autoregressive generation more efficiently. Inspired by human
 311 saccadic vision, where only the gaze center is rendered in high resolution, future autoregressive
 312 models might employ foveated rendering, focusing compute on regions of interest and blurring
 313 others—achieving adaptive resolution in both time and space.

314 6.4 Postdictive Coding

315 Most video generation research emphasizes predictive coding—inferring future states from prior ob-
 316 servations [52]—while largely neglecting postdictive coding [53], the process by which new sensory

317 input reshapes interpretations of earlier events. This retroactive revision is vital for autoregressive,
318 online world models, where predictions must be continuously updated as new data arrives. When
319 mismatches occur between earlier forecasts and later inputs, models need to revise latent memory
320 to preserve consistency. Cognitive science supports this: effects like the flash-lag and cutaneous
321 rabbit illusions show how the brain integrates sensory information over time to construct coherent,
322 post-hoc perceptions. Yet despite its biological relevance, postdictive coding is mostly absent from
323 current video generation frameworks, which typically use rigid feedforward designs—where the
324 past influences the future, but not vice versa. This limits real-time adaptability and self-correction.
325 Incorporating postdictive mechanisms could refine past representations with future evidence, en-
326 hancing causal reasoning and temporal coherence. A hybrid system combining predictive foresight
327 with postdictive revision promises more adaptive, cognitively grounded video models—capable of
328 learning not only by anticipating what comes next, but by reinterpreting what came before.

329 7 Alternative Views

330 **View #1: If autoregressive generation was once abandoned, why return to it now—and what** 331 **makes this time different?**

332 The limitations of early autoregressive video generation were not due to flaws in the paradigm itself,
333 but external constraints. The dominant task—predicting future frames from short histories—ignored
334 multimodal control, not from oversight, but due to missing interfaces, limited paired video data,
335 and the lack of strong generative backbones. Today, these barriers are dissolving. With growing
336 interest in fine-grained control and rich conditioning, autoregressive modeling is being revived—not
337 as regression, but as a path to real-time, interactive world simulation.

338 **View #2: Autoregressive generation still suffers from compounding errors—how do we address** 339 **that this time?**

340 Compounding error has long challenged autoregressive video generation, stemming from the accumu-
341 lation of small imperfections over time. Earlier models often produced blurry frames or artifacts that,
342 when recursively reused, degraded quality. Yet this issue isn't inherent to autoregression—it depends
343 on the fidelity and temporal consistency of predictions. When outputs are clean and coherent, error
344 accumulation becomes far less problematic.

345 Text conditioning further narrows the space of possible futures, injecting semantic structure and
346 reducing uncertainty. Importantly, revisiting autoregression today doesn't mean abandoning text
347 guidance—it signals a move toward uniting fine-grained temporal modeling with high-level intent.
348 Thanks to improved architectures, richer training data, and multimodal inputs, the factors that once
349 amplified compounding errors have been largely addressed.

350 **View #3: If my sole purpose for video generation is to create creative content, does the return to** 351 **autoregressive models still not concern me?**

352 Yes—autoregressive generation offers practical advantages beyond world modeling. Most notably, it
353 enables faster inference, giving creators near-instant feedback during creation. This speed accelerates
354 iteration and makes interactive, game-like content possible, pushing beyond static video playback
355 into dynamic, engaging media.

356 8 Conclusion

357 We have taken the position that video generation, particularly in the context of world modeling,
358 stands to benefit significantly from a return to autoregressive paradigms. While diffusion-based
359 models have excelled at holistic scene synthesis and compositional generalization, they fall short in
360 scenarios that demand real-time responsiveness and precise control. Autoregressive generation, in
361 contrast, is inherently better suited for integrating frame-wise signals, accommodating multimodal
362 conditions, and adapting to continuous feedback—all of which are essential for interactive and
363 causally coherent simulation. By framing generation as a sequential prediction task, autoregressive
364 models promote compact, causally grounded representations that are not only faster to compute
365 but also more aligned with the needs of embodied agents. As we look toward building practical,
366 interpretable, and controllable video models for future world-simulating systems, we believe the
367 autoregressive approach offers the most promising foundation.

368 **References**

- 369 [1] OpenAI. Sora. <https://www.openai.com/sora>, 2024. Sora is an AI model that can create
370 realistic and imaginative scenes from text instructions.
- 371 [2] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Do-
372 minik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion:
373 Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 374 [3] Xingjian Shi, Zhoung Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun
375 Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting.
376 In *Advances in neural information processing systems*, pages 802–810, 2015.
- 377 [4] Yunbo Wang, Mingsheng Long, Jianmin Wang, Zhifeng Gao, and S Yu Philip. Predrnn:
378 Recurrent neural networks for predictive learning using spatiotemporal lstms. In *Advances in*
379 *Neural Information Processing Systems*, pages 879–888, 2017.
- 380 [5] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 381 [6] Jakob Hohwy. *The Predictive Mind*. Oxford University Press, 2013.
- 382 [7] Anil Seth. *Being You: A New Science of Consciousness*. Faber Faber, 2021.
- 383 [8] Geoffrey E. Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep
384 belief nets. *Neural Computation*, 18:1527–1554, 2006.
- 385 [9] Yoshua Bengio and Yann LeCun. Scaling learning algorithms towards AI. In *Large Scale*
386 *Kernel Machines*. MIT Press, 2007.
- 387 [10] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
388 convolutional neural networks. In *Advances in neural information processing systems*, pages
389 1097–1105, 2012.
- 390 [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
391 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
392 pages 770–778, 2016.
- 393 [12] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video
394 representations using lstms. In *International conference on machine learning*, pages 843–852,
395 2015.
- 396 [13] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*,
397 9(8):1735–1780, 1997.
- 398 [14] Yunbo Wang, Zhifeng Gao, Mingsheng Long, Jianmin Wang, and Philip S Yu. Predrnn++:
399 Towards a resolution of the deep-in-time dilemma in spatiotemporal predictive learning. *arXiv*
400 *preprint arXiv:1804.06300*, 2018.
- 401 [15] Wei Yu, Yichao Lu, Steve Easterbrook, and Sanja Fidler. Efficient and information-preserving
402 future frame prediction and beyond. 2020.
- 403 [16] Alex X Lee, Richard Zhang, Frederik Ebert, Pieter Abbeel, Chelsea Finn, and Sergey Levine.
404 Stochastic adversarial video prediction. *arXiv preprint arXiv:1804.01523*, 2018.
- 405 [17] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior. *arXiv preprint*
406 *arXiv:1802.07687*, 2018.
- 407 [18] Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013.
- 408 [19] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
409 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural*
410 *information processing systems*, pages 2672–2680, 2014.
- 411 [20] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in*
412 *neural information processing systems*, 30, 2017.

- 413 [21] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution
414 image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
415 recognition*, pages 12873–12883, 2021.
- 416 [22] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
417 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on
418 machine learning*, pages 8821–8831. Pmlr, 2021.
- 419 [23] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin,
420 Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via
421 transformers. *Advances in neural information processing systems*, 34:19822–19835, 2021.
- 422 [24] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper-
423 vised learning using nonequilibrium thermodynamics. In *International conference on machine
424 learning*, pages 2256–2265. pmlr, 2015.
- 425 [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances
426 in neural information processing systems*, 33:6840–6851, 2020.
- 427 [26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks
428 for biomedical image segmentation. In *Medical image computing and computer-assisted
429 intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9,
430 2015, proceedings, part III 18*, pages 234–241. Springer, 2015.
- 431 [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
432 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF
433 conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 434 [28] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Ani-
435 matediff: Animate your personalized text-to-image diffusion models without specific tuning.
436 *arXiv preprint arXiv:2307.04725*, 2023.
- 437 [29] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings
438 of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- 439 [30] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
440 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information
441 processing systems*, 30, 2017.
- 442 [31] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image
443 prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.
- 444 [32] Yilun Du, Sherry Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans,
445 and Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in
446 Neural Information Processing Systems*, 36, 2024.
- 447 [33] Wei Yu, Wenxin Chen, Songheng Yin, Steve Easterbrook, and Animesh Garg. Modular action
448 concept grounding in semantic video prediction. In *Proceedings of the IEEE/CVF Conference
449 on Computer Vision and Pattern Recognition*, pages 3605–3614, 2022.
- 450 [34] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv
451 preprint arXiv:2010.02502*, 2020.
- 452 [35] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver:
453 A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in
454 Neural Information Processing Systems*, 35:5775–5787, 2022.
- 455 [36] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li, Zhipeng Huang, Xiangjun Gao,
456 Tien-Tsin Wong, Ying Shan, and Yonghong Tian. Viewcrafter: Taming video diffusion models
457 for high-fidelity novel view synthesis. *arXiv preprint arXiv:2409.02048*, 2024.
- 458 [37] Wei Yu, Songheng Yin, Steve Easterbrook, and Animesh Garg. Egocentric exploration
459 in virtual worlds with multi-modal conditioning. In *First Workshop on Controllable Video
460 Generation@ ICML24*.

- 461 [38] Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu
462 Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evaluating physical
463 commonsense for video generation. *arXiv preprint arXiv:2406.03520*, 2024.
- 464 [39] Bingyi Kang, Yang Yue, Rui Lu, Zhijie Lin, Yang Zhao, Kaixin Wang, Gao Huang, and Jiashi
465 Feng. How far is video generation from world model: A physical law perspective. *arXiv*
466 *preprint arXiv:2411.02385*, 2024.
- 467 [40] Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christo-
468 pher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, et al.
469 Language modeling is compression. *arXiv preprint arXiv:2309.10668*, 2023.
- 470 [41] Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai,
471 Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, et al. Experience
472 grounds language. *arXiv preprint arXiv:2004.10151*, 2020.
- 473 [42] Daniel Casasanto. Language, cognition and space. In *Proceedings of the 30th Annual Meeting*
474 *of the Cognitive Science Society*, pages 1090–1095, 2008.
- 475 [43] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang,
476 Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions.
477 *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- 478 [44] Yeqing Lin, Minji Lee, Zhao Zhang, and Mohammed AlQuraishi. Out of many, one: Designing
479 and scaffolding proteins at the scale of the structural universe with genie 2. *arXiv preprint*
480 *arXiv:2405.15489*, 2024.
- 481 [45] Amir Bar, Gaoyue Zhou, Danny Tran, Trevor Darrell, and Yann LeCun. Navigation world
482 models. *arXiv preprint arXiv:2412.03572*, 2024.
- 483 [46] Jiwen Yu, Yiran Qin, Xintao Wang, Pengfei Wan, Di Zhang, and Xihui Liu. Gamefactory:
484 Creating new games with generative interactive videos. *arXiv preprint arXiv:2501.08325*, 2025.
- 485 [47] Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent
486 Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in*
487 *Neural Information Processing Systems*, 37:24081–24125, 2024.
- 488 [48] Ziyi Yang, Zaibin Zhang, Zirui Zheng, Yuxian Jiang, Ziyue Gan, Zhiyu Wang, Zijian Ling,
489 Jinsong Chen, Martz Ma, Bowen Dong, et al. Oasis: Open agents social interaction simulations
490 on one million agents. *arXiv preprint arXiv:2411.11581*, 2024.
- 491 [49] Hansi Teng, Hongyu Jia, Lei Sun, Lingzhi Li, Maolin Li, Mingqiu Tang, Shuai Han, Tianning
492 Zhang, WQ Zhang, Weifeng Luo, et al. Magi-1: Autoregressive video generation at scale. *arXiv*
493 *preprint arXiv:2505.13211*, 2025.
- 494 [50] Lvmin Zhang and Maneesh Agrawala. Packing input frame context in next-frame prediction
495 models for video generation. *arXiv preprint arXiv:2504.12626*, 2025.
- 496 [51] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive
497 modeling: Scalable image generation via next-scale prediction. *Advances in neural information*
498 *processing systems*, 37:84839–84865, 2024.
- 499 [52] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video
500 prediction and unsupervised learning. *arXiv preprint arXiv:1605.08104*, 2016.
- 501 [53] Shinsuke Shimojo. Postdiction: its implications on visual awareness, hindsight, and sense of
502 agency. *Frontiers in Psychology*, 5:196, 2014.