

INFINITE-RESOLUTION INTEGRAL NOISE WARPING FOR DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Adapting pretrained image-based diffusion models to generate temporally consistent videos has become an impactful generative modeling research direction. Training-free noise-space manipulation has proven to be an effective technique, where the challenge is to preserve the Gaussian white noise distribution while adding in temporal consistency. Recently, Chang et al. (2024) formulated this problem using an integral noise representation with distribution-preserving guarantees, and proposed an upsampling-based algorithm to compute it. However, while their mathematical formulation is advantageous, the algorithm incurs a high computational cost. Through analyzing the limiting-case behavior of their algorithm as the upsampling resolution goes to infinity, we develop an alternative algorithm that, by gathering increments of multiple Brownian bridges, achieves their infinite-resolution accuracy while simultaneously reducing the computational cost by orders of magnitude. We prove and experimentally validate our theoretical claims, and demonstrate our method’s effectiveness in real-world applications. We further show that our method can readily extend to the 3-dimensional space.

1 INTRODUCTION

The success of diffusion models for image generation and manipulation (Rombach et al., 2022; Nichol et al., 2021; Ho et al., 2020; Zhang et al., 2023a) has spurred significant interest in lifting these capacities to the video domain (Singer et al., 2022; Durrett, 2019; Gupta et al., 2023; Blattmann et al., 2023; Ho et al., 2022; Guo et al., 2024). While building video models trained directly on spatiotemporal data is a natural idea, practical concerns such as limited availability of large-scale video data and high computational cost have motivated investigations into training-free alternatives. One such training-free approach is to use image models to directly generate video frames, utilizing various techniques such as cross-frame attention, feature injection and hierarchical sampling to promote cross-frame temporal consistency (Zhang et al., 2023b; Khachatryan et al., 2023; Cong et al., 2023).

One particularly effective consistency-promoting technique is the temporally consistent initialization of noise across frames. However, most existing approaches either result in a loss of Gaussianity in the noise image (and subsequently introduces a domain gap at inference time), or restrict themselves to simple manipulations of the noise image (*e.g.* filtering and blending). Recently, Chang et al. (2024) proposed a method that preserves both Gaussian white noise distribution and cross-frame temporal correlations via *integral noise warping*: each warped noise pixel is obtained by integrating a continuous noise field, where the integration is implemented by summing over a polygonal deformed pixel region in the upsampled prior noise image. However, the time and memory costs of this algorithm grows quadratically with the upsampling resolution, prohibiting the adoption of the method in certain applications (Kwak et al., 2024) due to this computational expense.

Regarding this challenge, our key insight is that when adopting an Eulerian perspective as opposed to the original Lagrangian perspective, the limiting-case algorithm of Chang et al. (2024) for computing a warped noise pixel reduces to summing over increments from multiple Brownian bridges (Durrett, 2019, Section 8.4). In place of the costly upsampling procedure, sampling the increments of a Brownian bridge can be done efficiently in an autoregressive manner (2). We thus propose *infinite-resolution integral noise warping* (Algorithm 1), which can directly and efficiently resolve noise transport in the continuous space, when supplied with an oracle that returns the overlapping area between a pixel square and a deformed pixel region (Section 2.3).

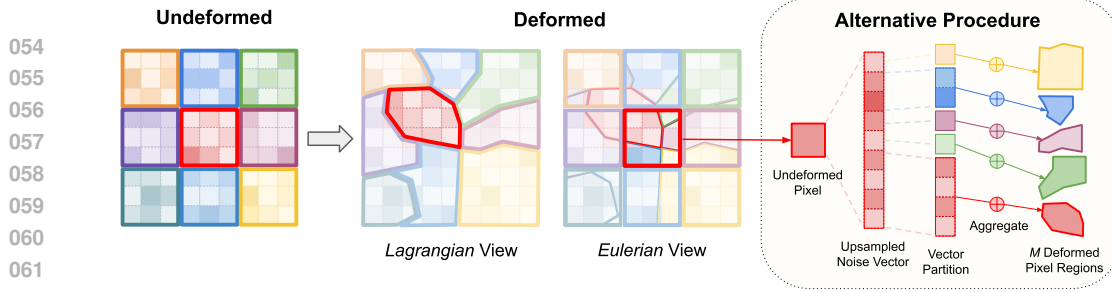


Figure 1: When the image grid deforms, the Lagrangian view tracks a deformed pixel region, while the Eulerian view tracks the undeformed pixel square as it gets partitioned into multiple regions. On the right, we leverage the exchangeability of upsampled subpixels to convert the Lagrangian gathering procedure into scattering noise subpixels to overlapped deformed pixel regions.

We propose two concrete methods for computing this oracle, leading to a *grid-based* and a *particle-based* variant. Following Chang et al. (2024), the *grid-based* variant (Algorithm 2) computes the area by explicitly constructing per-pixel deformed polygons, and is exactly equivalent to the existing approach (Chang et al., 2024) with an infinite upsampling resolution, while running $8.0\times$ to $19.7\times$ faster and using $9.22\times$ less memory. Inspired by hybrid Eulerian-Lagrangian fluid simulation (Brackbill et al., 1988), our novel *particle-based* variant (Algorithm 3) computes area in a fuzzy manner, which not only offers a further $5.21\times$ speed-up *over our grid-based variant*, but is also agnostic to non-injective maps. In real-world scenarios, the particle-based variant shows no compromise in generation quality compared to the grid-based one (see video results), while offering superior robustness, efficiency, simplicity, and extensibility to higher dimensions.

In summary, we propose a novel method for computing temporally correlated noise to facilitate consistent video generation with image-based diffusion models. Our algorithm computes the integral noise from Chang et al. (2024) at infinite resolution, warping a 1024×1024 noise image in ~ 0.045 s (grid variant) and ~ 0.0086 s (particle variant), achieving orders of magnitude speed-up compared to Chang et al. (2024) while retaining the distribution-preserving and temporally-coherent properties.

2 METHODOLOGY

In this section, we introduce our method as follows:

- We present an equivalent Eulerian interpretation (Figure 1) for the method by Chang et al. (2024), which was developed from a Lagrangian viewpoint.
- We show that the limiting algorithm of the Eulerian formulation as upsampling level goes to infinity is equivalent to sampling increments of Brownian bridges.
- We present our main algorithm (Algorithm 1) which, given a partition record that returns the overlapping area between a pixel square and a deformed pixel region, samples increments of Brownian bridges and scatters the increments to form the warped noise image.
- We propose two practical algorithms for computing the overlap areas. The *grid-based* Algorithm 2 extends Chang et al. (2024) to infinite resolution without the overhead of upsampling. The *particle-based* Algorithm 3 departs from grid-based discretization and uses particles instead, resulting in a simpler algorithm that is robust to degenerate maps.

2.1 NOISE WARPING: AN ALTERNATIVE EULERIAN PERSPECTIVE

Given a $D \times D$ prior noise image $I_W \in \mathbb{R}^{D \times D^1}$ and a deformation map $\psi : [0, 1]^2 \rightarrow [0, 1]^2$, the noise-warping algorithm (Chang et al., 2024) computes the warped noise image $\tilde{I}_W \in \mathbb{R}^{D \times D}$ with upsampling level $N \in \mathbb{Z}_{\geq 1}$ as follows:

¹In practice, the noise image has multiple channels. Since channels are always treated independent of one another, to simplify the notation, we will assume the image has a single channel.

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

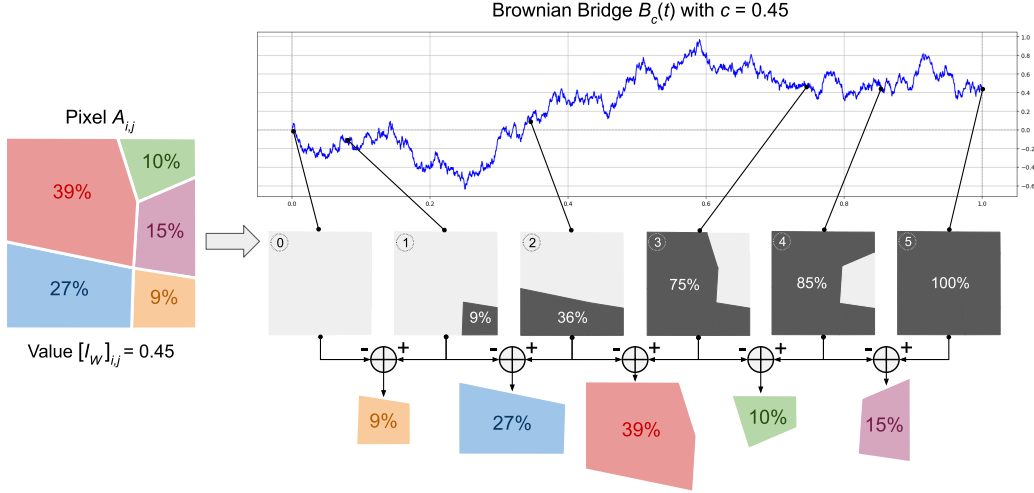


Figure 2: Connection between Eulerian noise-warping and increments of a Brownian bridge for a fixed prior noise pixel $[I_W]_{i,j}$. The overlapping area of each colored warped region becomes the time increment for the Brownian bridge. Hence, sampling the Brownian bridge at these times and taking consecutive differences yields integral noise that is scattered to form each warped noise pixel.

1. For $i, j = 1, \dots, D$, upsample noise pixel $[I_W]_{i,j}$ to an $N \times N$ subimage $[\widehat{I}_W]_{i,j} \in \mathbb{R}^{N \times N}$:

$$[\widehat{I}_W]_{i,j} = \frac{[I_W]_{i,j}}{N^2} + \frac{1}{N} \left(Z - \frac{S}{N^2} \right), \text{ with } Z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ and } S = \sum_{k=1}^{N^2} Z_k. \quad (1)$$

The subimage for each pixel assembles into an $ND \times ND$ upsampled noise image \widehat{I}_W .

2. For $i, j = 1, \dots, D$, the pixel square $A_{i,j} := [\frac{i-1}{D}, \frac{i}{D}] \times [\frac{j-1}{D}, \frac{j}{D}]$ is warped to a deformed pixel region $\widetilde{A}_{i,j} := \psi(A_{i,j})$, and the warped noise pixel $[\widetilde{I}_W]_{i,j}$ is set to be the sum of all subpixels in \widehat{I}_W covered by $\widetilde{A}_{i,j}$ divided by $\sqrt{|\widetilde{A}_{i,j}|}$, where $|A|$ denotes the Lebesgue measure of a Borel set $A \subset \mathbb{R}^2$.

We describe an alternative but equivalent procedure by making the following two observations, which are illustrated in Figure 1.

Gathering Noise \rightarrow Scattering Noise. While the original procedure computes the warped noise image by *gathering* the upsampled noise subpixels in each deformed pixel region $\widetilde{A}_{i,j}$ in a *Lagrangian* fashion, we can instead use an alternative procedure by *scattering* the upsampled noise subpixels in each pixel square $A_{i,j}$ to overlapping deformed pixel regions. This new *Eulerian* procedure does not change the output, but it yields new insights in conjunction with our second observation.

Scattering Noise \rightarrow Counting Overlapping Subpixels. Observe that the $N \times N$ subpixels in $[\widehat{I}_W]_{i,j}$, for every i, j , are correlated only through their sum S when conditioning on $[I_W]_{i,j}$ (1), so they are exchangeable. Hence, when scattering these upsampled noise subpixels to deformed pixel regions, the order of scattering does not matter, and we only need to count *the number of subpixels* covered by each deformed pixel region.

Alternative Eulerian Procedure. Putting both observations together, we now describe an alternative procedure to Chang et al. (2024) with unaltered output:

1. For each noise image pixel $[I_W]_{i,j}$, draw an upsampled subimage, now represented as a 1D vector $X \in \mathbb{R}^{N^2}$ using (1). Then, compute a prefix sum $H_{i,j}$ via $[H_{i,j}]_k := \sum_{q=1}^k X_q$ for $k = 1, \dots, N^2$.
2. Warp each pixel square and compute deformed pixel regions $\widetilde{A}_{i,j}$ as before.
3. For each $A_{i,j}$, let M denote the number of deformed pixel regions that overlap with $A_{i,j}$. With index $k = 1, \dots, M$, we use l_k, m_k to denote the coordinates of the k^{th} overlap, whose pixel region is \widetilde{A}_{l_k, m_k} and pixel value $[\widetilde{I}_W]_{l_k, m_k}$. Form $L \in \mathbb{Z}_{\geq 0}^M$ where L_k represents the

number of upsampled subpixels covered by \tilde{A}_{ℓ_k, m_k} . Then, compute a prefix sum $[C_{i,j}]_k := \sum_{q=1}^k L_q$. For $k = 1, \dots, M$, accrue $[H_{i,j}]_{[C_{i,j}]_k} - [H_{i,j}]_{[C_{i,j}]_{k-1}}$ to $[\tilde{I}_W]_{\ell_k, m_k}$.

4. Divide each warped noise pixel $[\tilde{I}_W]_{i,j}$ by $\sqrt{|\tilde{A}_{i,j}|}$.

Discussion. Compared to the original procedure by Chang et al. (2024), this alternative but equivalent algorithm highlights how the upsampled subpixels of $[I_W]_{i,j}$ are scattered to form the warped noise pixels. In particular, each warped noise pixel receives the *sum of a continuous segment* in $H_{i,j}$. Since $H_{i,j}$ is a summation of weakly correlated and exchangeable subpixels, once conditioned on $[I_W]_{i,j}$, can we avoid explicitly instantiating every single subpixel, but instead model the *sum* of these weakly correlated subpixels?

The key insight of this paper is that when the upsampling resolution $N \rightarrow \infty$, the scaling limit of the prefix sum $H_{i,j}$ (with proper interpolation and time scaling to a continuous function) is precisely the Brownian bridge (Durrett, 2019, Section 8.4) conditioned on $[I_W]_{i,j}$. Once this connection is established, it is easy to progressively sample increments of the Brownian bridge, resulting in a clean and efficient noise-warping algorithm that bypasses the need for upsampling in Chang et al. (2024).

2.2 INFINITE-RESOLUTION NOISE SCATTERING

In this section, we first derive a scaling limit result to Brownian bridges. We then illustrate that the limiting version of the Eulerian procedure from the previous section matches precisely this scaling limit result. Lastly, we demonstrate an autoregressive way to sample increments of a Brownian bridge that is linear in runtime in terms of the number of increments.

Theorem 1 (Scaling limit to Brownian bridge). Let $\{Z_n\}$ be a sequence of i.i.d. random variables with finite variance that are normalized such that $\mathbb{E}[Z_n] = 0$ and $\text{Var}(Z_n) = 1$. For $c \in \mathbb{R}$, define

$$S_n := \sum_{i=1}^n Z_i, \quad X_{i,n} := \frac{c}{n} + \frac{1}{\sqrt{n}} \left(Z_i - \frac{S_n}{n} \right).$$

Consider the sequence of random continuous functions $\{H_n(t)\} \subset C[0, 1]$ defined as

$$H_n(t) := \sum_{i=1}^{\lfloor nt \rfloor} X_{i,n} + (nt - \lfloor nt \rfloor) X_{\lfloor nt \rfloor + 1, n}.$$

Then the sequence $\{H_n\}$ converges in distribution under the sup-norm metric on $C[0, 1]$ to $B_c(t) := W(t) - tW(1) + tc$, the Brownian bridge ending at c , where $W(t)$ is standard Brownian motion. Moreover, in distribution, we have $B_c(t) \stackrel{d}{=} (W(t) \mid W(1) = c)$, where $(W(t) \mid W(1) = c)$ is the disintegrated measure (Pachl, 1978) of $W(t)$ on $W(1) = c$.

We prove Theorem 1 in Appendix A. To connect the Eulerian procedure with the setup in Theorem 1, let us fix a pixel $[I_W]_{i,j}$, and let $B := B_{[I_W]_{i,j}}$, $H := H_{i,j}$, $C := C_{i,j}$ to simplify the notation. By setting $n = N^2$ and $c = [I_W]_{i,j}$, the sequence $\{X_{k,n}\}$ from the theorem has exactly the same law as the upsampled subpixels in $[\tilde{I}_W]_{i,j}$. Moreover, $H_{nt} = H_n(t)$ when $nt \in \mathbb{Z}_{\geq 1}$. By taking $N \rightarrow \infty$, implying $n \rightarrow \infty$, for any $t_1, \dots, t_M \in [0, 1]$, we have the convergence in distribution of $(H_{\lfloor nt_1 \rfloor}, \dots, H_{\lfloor nt_M \rfloor}) \stackrel{d}{\rightarrow} (B(t_1), \dots, B(t_M))$. Recall in the Eulerian procedure, we only need to access the prefix sum H at indices $\{C_k\}_{k=1}^M$, where C_k counts the number of upsampled subpixels covered by the first k overlaps. This suggests that if we choose

$$t_k = \lim_{N \rightarrow \infty} \frac{C_k}{N^2} = \sum_{k'=1}^k \left| A_{i,j} \cap \tilde{A}_{\ell_{k'}, m_{k'}} \right|,$$

and use $B(t_k)$ in place of H_k , then we just need to sample from B at times t_1, \dots, t_M — precisely the limiting algorithm of the Eulerian procedure. We illustrate this connection in Figure 2.

Autoregressive Sampling of Brownian Bridges. Since a Brownian bridge is a Markov process (Oksendal, 2013, Exercise 5.11), we can sample the vector $(B_c(t_1), \dots, B_c(t_M))$ in an autoregressive fashion, each time sampling $B_c(t_{k+1})$ conditioned on $B_c(t_k)$:

$$(B_c(t_{k+1}) \mid B_c(t_k) = q) \stackrel{d}{=} \mathcal{N} \left(\frac{1-t_{k+1}}{1-t_k} q + \frac{t_{k+1}-t_k}{1-t_k} c, \frac{(t_{k+1}-t_k)(1-t_{k+1})}{1-t_k} \right). \quad (2)$$

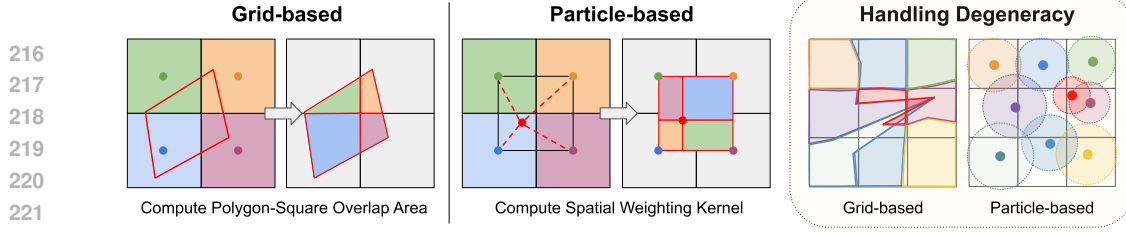


Figure 3: In the grid-based variant, we compute the overlapping area by explicitly constructing the polygonal region of the deformed pixel region. In the particle-based variant, we approximate the overlapping area with a linear kernel-based interpolation. On the right, we show how the two variants differ when the map is non-injective: the grid-based variant yields intersecting regions, causing stability issues, whereas the particle-based variant remains stable.

Once the Brownian bridge at times t_k is sampled, we just need to accrue the increments $B_c(t_k) - B_c(t_{k-1})$ to $[\tilde{I}_W]_{\ell_k, m_k}$, the k^{th} overlapped warped noise pixel. This allows us to present Algorithm 1. Compared to the discrete procedures described earlier, we no longer need upsampling. In addition, we exploited the autoregressive nature of Brownian bridges to bring down the time complexity to linear in the number of overlapping warped pixel regions.

Algorithm 1 Infinite-Resolution Integral Noise Warp

Input: prior noise image $I_W \in \mathbb{R}^{D \times D}$, deformation map $\psi : [0, 1] \rightarrow [0, 1]$

Output: warped noise image $\tilde{I}_W \in \mathbb{R}^{D \times D}$

Build a partition record \mathcal{P} from ψ (Section 2.3)

Initialize $\mathcal{A}_{i,j} \leftarrow 0$ for all $i, j = 1, \dots, D$

$\triangleright \mathcal{A}_{i,j}$ will eventually be the area of $\tilde{\mathcal{A}}_{i,j}$

parallel for each $u, v = 1, \dots, D$ **do**

$t, q, M \leftarrow 0, 0, |\mathcal{P}_{u,v}|$

for $k = 1, \dots, M$ **do**

$(a, i, j) \leftarrow [\mathcal{P}_{u,v}]_k$ $\triangleright a$ is the overlapping area between $\mathcal{A}_{i,j}$ and $\tilde{\mathcal{A}}_{u,v}$

Sample $q' \sim (B_c(t+a) | B_c(t) = q)$ by (2) with $c = [I_W]_{u,v}$

$[\tilde{I}_W]_{i,j} \leftarrow [\tilde{I}_W]_{i,j} + (q' - q)$

$\mathcal{A}_{i,j} \leftarrow \mathcal{A}_{i,j} + a$

$q, t \leftarrow q', t + a$

Normalize $[\tilde{I}_W]_{i,j} \leftarrow \mathcal{A}_{i,j}^{-\frac{1}{2}} [\tilde{I}_W]_{i,j}$ for all $i, j = 1, \dots, D$

return \tilde{I}_W

Preservation of Gaussian White Noise. A central desideratum of noise warping is that the resulting warped noise image \tilde{I}_W needs to have pixels that are i.i.d. standard Gaussians when the prior noise image I_W is Gaussian white noise. This ensures that the warped noise is in-distribution for a pre-trained diffusion model. Our algorithm automatically guarantees this preservation of Gaussianity, as long as the warping function ψ is injective. To see this, the injectivity of ψ implies that the warped pixel regions are non-overlapping in the square $[0, 1]^2$. For each $\mathcal{A}_{i,j}$, since $[I_W]_{i,j} \stackrel{d}{=} \mathcal{N}(0, 1) \stackrel{d}{=} W(1)$, by the conditional interpretation of Brownian bridges (1), when marginalizing out $[\tilde{I}_W]_{i,j}$, the Brownian bridge $B_{[I_W]_{i,j}}$ reduces to standard Brownian motion. Since the increments of the Brownian motion are independent Gaussians, the contribution to a deformed pixel region is simply a zero-mean Gaussian with variance equal to the overlapping area. Therefore, each deformed pixel region will receive the sum of a number of independent Gaussians whose variances sum to the area of the region. The scaling by the inverse square root of the area in Algorithm 1 thus makes each warped noise pixel an i.i.d. standard Gaussian.

2.3 BUILDING PARTITION RECORDS

In this section we present two algorithms, one grid-based and one particle-based, for building partition records that return the area between each pixel square and its overlapping deformed pixel regions. Both versions are outlined in Algorithm 2 and Algorithm 3, where we use i, j to index pixel regions in the deformed space, and u, v for those in the undeformed space.

Algorithm 2 Grid-based Partition

Input: Deformation map ψ
Output: Partition record \mathcal{P}

```

1: parallel for each  $i, j$  do
2:    $A^* \leftarrow \text{DiscretizeSquare}(A_{i,j})$ 
3:    $S \leftarrow \psi(A^*)$ 
4:    $u^-, u^+, v^-, v^+ \leftarrow \text{AABB}(S)$ 
5:   for  $u \in [u^-, u^+]$  do
6:     for  $v \in [v^-, v^+]$  do
7:        $a \leftarrow \text{PolygonArea}(\text{Clip}(S, u, v))$ 
8:        $\mathcal{P}_{u,v} \leftarrow \mathcal{P}_{u,v} + [(a, i, j)]$ 
9: return  $\mathcal{P}$ 

```

Algorithm 3 Particle-based Partition

Input: Deformation map ψ
Output: Partition record \mathcal{P}

```

1: parallel for each  $i, j$  do
2:    $(x, y) \leftarrow \psi(\frac{i+0.5}{D}, \frac{j+0.5}{D})$ 
3:    $\alpha_{0,0}, \alpha_{0,1}, \alpha_{1,0}, \alpha_{1,1} \leftarrow \text{BilinearWeights}(X)$ 
4:   for  $s, t \in [0, 1]$  do
5:      $u, v \leftarrow \lfloor x \rfloor + s, \lfloor y \rfloor + t$ 
6:      $\mathcal{P}_{u,v} \leftarrow \mathcal{P}_{u,v} + [(\alpha_{s,t}, i, j)]$ 
7: parallel for each  $u, v$  do
8:   Normalize total area of  $\mathcal{P}_{u,v}$  to  $D^{-2}$ 
9: return  $\mathcal{P}$ 

```

As illustrated in Figure 3, our grid-based method (left) follows the same approach as Chang et al. (2024), treating each deformed pixel as an octagon and each undeformed pixel as a grid cell; our particle-based method (middle) borrows from the grid-to-particle (G2P) technique in fluid particle-in-cell method methods (Brackbill et al., 1988), where we treat each deformed pixel as a particle and each undeformed pixel as grid cell. Each particle requests area from nearby grid cells based on distance; upon receiving requests, each grid cell normalizes the requests to ensure partition-of-unity, and distributes its area to contacting particles.

Discussion. Conceptually, our grid-based and particle-based methods correspond to two different interpretations of ψ when it is available only as discrete samples (e.g. optical flow image). The grid-based method implicitly reconstructs the continuous ψ field by linear interpolation, whereas the particle-based method makes no such interpolation and assumes ψ is only known point-wise. This implies that when ψ is smooth, linear interpolation works well and the grid-based method will yield a higher-quality warp as seen in Figure B.4. But when ψ is non-smooth, which is usually the case in real-world data, linear interpolation becomes problematic and can lead to degenerate polygons as illustrated on the right of Figure 3. These degenerate polygons can violate the assumption that deformed pixel regions do not overlap, which is required for independence in the warped noise. We note that, in practice, both Chang et al. (2024) and our grid-based method are equipped with fail-safes that ensures this independence by “patching up” degenerate regions with new Gaussian noise, but these mechanisms do not handle the overlaps in a principled manner, while the particle-based variant fundamentally circumvents such overlaps with its topology-free nature.

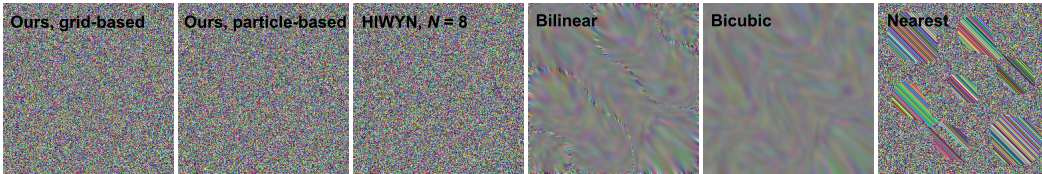
In addition, we also highlight the simplicity of the particle-based approach, as it boils down the partition record computation to evaluating bilinear kernels, which is highly efficient and parallelizable. Such simplicity offers a unique opportunity to extend the noise warping to higher dimensions, as it effectively only requires the change from the bilinear kernel to its higher-dimensional counterpart. We showcase this possibility using a 3-dimensional noise warp example as shown in Figure B.5.

3 RESULTS

We verify our theoretical claims by showing that both variants of our method preserve Gaussian white noise distribution, and that Chang et al. (2024) (HIWYN) converges to our grid-based variant as N increases. We analyze the behaviors of our grid-based and particle-based variants under diffeomorphic and non-diffeomorphic deformations. We then apply our method in video generation and benchmark against existing methods (Ge et al., 2023; Chen et al., 2023; Chang et al., 2024). Finally, we extend our method to warping volumetric noise and demonstrate a use case in 3D graphics.

Gaussian White Noise Preservation. In Figure 4, we iteratively warp a noise image by the same deformation map for 50 timesteps. We gauge the output noise’s resemblance to Gaussian white noise by measuring normality using one-sample Kolmogorov-Smirnov (K-S) test and detecting spatial correlation using Moran’s I . Our results show that both HIWYN and our method generate warped noise images indistinguishable from Gaussian white noise which significantly improve upon baselines.

Convergence of Chang et al. (2024). We validate that our method is the limiting case of HIWYN. Starting with an 8×8 prior noise image and a flow map (Figure 5, top left), we run our method along with HIWYN for $N \in \{2, 4, 8, \dots, 256\}$ for 100,000 independent runs to estimate the distribution of



Distribution Preservation Metrics

Method	Ours, grid-based	Ours, particle-based	HIWYN, $N=8$	Bilinear	Bicubic	Nearest Neighbor
Moran's I	5.103e-4 / 0.849	-1.995e-3 / 0.475	3.215e-3 / 0.243	0.612 / 0	0.983 / 0	2.974e-2 / 6.103e-27
K-S Test	3.410e-3 / 0.430	3.023e-3 / 0.586	3.274e-3 / 0.482	0.366 / 0	0.422 / 0	9.806e-3 / 6.681e-06

Figure 4: Preservation of Gaussian white noise achieved by different warping methods. We report scores and p-values for both Moran's I (spatial correlation) and K-S test (normality). We show that results from our method (both variants) and HIWYN are indistinguishable from white Gaussian noise, while generic warping methods lead to corrupted noise.

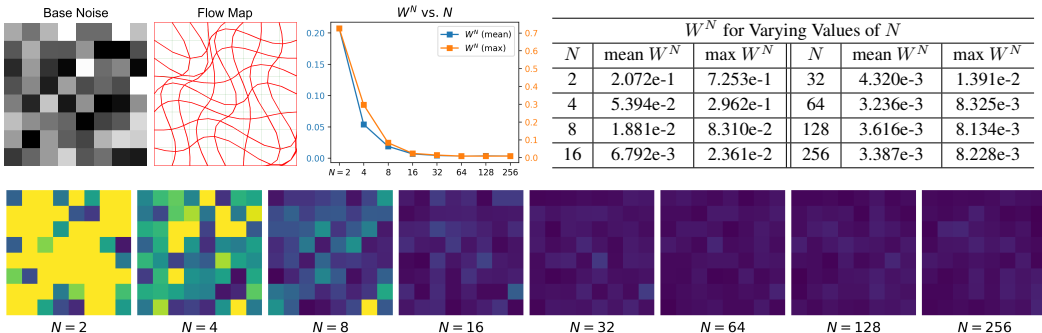


Figure 5: Convergence of HIWYN to our method as N increases. Top left: experimental setup with prior noise and deformation map. Top middle: 2-Wasserstein distance W^N between the output of HIWYN and ours. Top right: statistics table. Bottom: W^N difference image between the output of HIWYN and ours as N increases. Notice W^N becomes statistically insignificant for $N \geq 64$.

the warped noise image. For each upsampling resolution N , we compute the 2-Wasserstein distance W^N between the output of HIWYN and that of our method. The results in Figure 5 demonstrate the convergence of HIWYN to our method as N increase, and reveal that $N = 8$ (recommended by Chang et al. (2024)) is not yet in the converged phase to yield a negligible W^N .

Performance Comparison. For our methods and HIWYN with upsampling levels $N \in \{2, 4, 8\}$, we perform 100 independent runs on a 1024×1024 image. We report the kernel time with CPU and GPU backends (Figure 7) as well as the memory usage. The runtime and memory usage of our methods are largely comparable with those of HIWYN with $N = 2$. Compared to HIWYN with $N = 8$, both our methods offer order-of-magnitude improvements in runtime and memory usage. Specifically, our grid-based method achieves infinite upsampling resolution while being $19.7 \times$ faster on CPU and $8.0 \times$ faster on GPU, using $9.22 \times$ less memory, and our particle-based method, albeit not strictly equivalent to HIWYN at $N = \infty$, achieves a $41.7 \times$ speedup on GPU. In the following sections, we show that the particle-based version consistently achieves comparable quality to the grid-based version in real-world scenarios.

Comparison between Grid-Based and Particle-Based Variants. In Figure B.4, we compare both variants when the deformation map is diffeomorphic under different levels of distortion. Visually, the difference between the two variants is negligible at frame 25 and becomes noticeable in frame 100. We measure this difference by comparing the deformed regions for each pixel in terms of IoU and weighted Chamfer distance. We additionally compare the particle-based result with that of an identity-map baseline (right column in Figure B.4), which shows that the gap between the two variants remains small even under large distortion. In Figure 8, we stress test both variants under non-diffeomorphic maps obtained using optical flow (Teed & Deng, 2020) on a real-world video (Brox & Malik, 2011). In images 3 and 4, we see that the real-world flow map induces inverted meshes for the grid-based variant and clustered particles for the particle-based variant.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

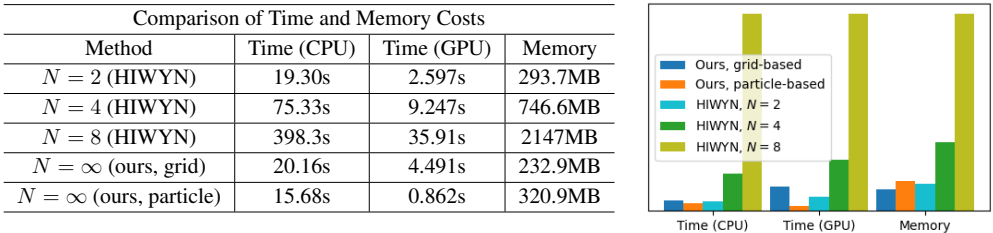


Figure 7: Runtime and memory usage of our method vs. HIWYN with $N = 2, 4, 8$. We compare total allocated memory and kernel time on a CPU/GPU. The computation is done on a laptop with Intel i7-12700H and GeForce RTX 3070 Ti.

While clustered particles are always assigned disjoint regions due to the continuous nature of our algorithm, mesh inversions cause area contention issue due to overlaps. In images 5 and 6, we mark the grid cells with area contention in orange, which occurs in the grid-based version but not in the particle-based version.

Video Super-resolution with I²SB. We integrate our method with I²SB (Liu et al., 2023) and adapt its pre-trained image $4\times$ super-resolution model (bicubic) to perform *video* super-resolution. We show our results in Figures 9 and B.2, and we refer to our supplementary video for best visualization of these results. Since I²SB is an image-to-image bridge model, it well preserves the low-frequency structures of the input images regardless of noise scheme. But as seen in our video, without noise warping, the results either show strong flickering in the high-frequency details (random noise) or sticking artifacts (fixed noise). Noise warping allows high-frequency details to transport with the optical flow, making the result significantly more consistent. We also validate that both our variants yield visual quality on par with HIWYN across all tested scenarios while being much more efficient.

Conditional Video Generation with SDEdit. We apply our method to conditional video generation by adapting SDEdit (Meng et al., 2021), a conditional image generation method, to producing consistent video frames. We apply Perturbed-Attention Guidance (Ahn et al., 2024) to the unconditional models with scale 3.0. Our two inputs are a conditioning video (generated by applying a median filter to real-world videos similar to Chen et al. (2023)) and an optical flow field (Teed & Deng, 2020). Without noise manipulation, if we run SDEdit frame-by-frame (Figure B.7, bottom row), the details (e.g. in the tower and trees) would result in strong flickering. By warping the noise using the optical flow, the temporal consistency is much improved. As shown in Figure B.7, our methods (both variants) and HIWYN yield comparable visual qualities. Full experiments that shows comparison with Control-A-Video (Chen et al., 2023) and PYoCo (Ge et al., 2023) and additional baselines are provided in Figures B.8 and B.9 with generation quality metrics reported in Figure B.1. Further results that additionally integrate cross-frame attention (Ceylan et al., 2023) (anchor every 3 frames) are shown in Figure 10 and B.3. We refer to our supplementary video for better visualization.

3D Noise Warp. We extend our particle-based algorithm to 3D by replacing the bilinear kernel with the bicubic kernel in Algorithm 3 and apply it to GaussianCube (Zhang et al., 2024), which denoises a dense 3D noise grid to reconstruct 3D Gaussians. We adapt it to perform conditional generation a la SDEdit. Starting with a 3D pickup truck generated unconditionally, we condition the model to generate vehicles with smaller and larger cabins by deforming the truck with a horizontal shear velocity field. We compare the results from using random noise to those using noise warped with our particle-based method. Using the warped noise improves the consistency, reducing the flickering of the cars’ geometries and textures. We show the results in Figure B.5 and refer to our supplementary video for better visualization.

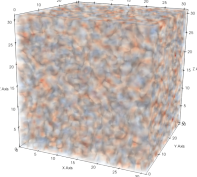
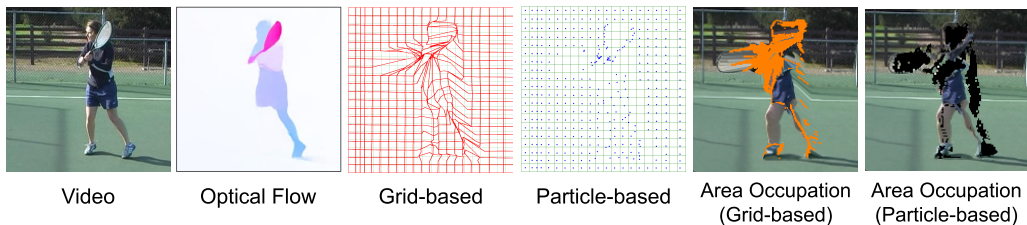


Figure 6: 3D noise warped by our particle variant.

4 RELATED WORKS

Noise in Diffusion Models. Diffusion models generate images from input noise, and noise can thus be considered the counterpart to the latent codes utilized in GAN models. As such, the outputs of diffusion models have dependencies and correlations to the initial input noise, making noise a

432
433
434
435
436
437
438
439440
441
442
Figure 8: Comparison of grid-based vs. particle-based variants under non-diffeomorphic optical flow. Pixels with detected overlaps are colored in orange. Further results are given in Figure B.6.443
444
445
446
447
448
449
450
451
452
453
454455
456
457
458
Figure 9: Video 4 \times super-resolution by integrating our method (particle) with I²SB. Top row shows the low resolution input video; bottom row shows the output video. Additional results are shown in Figure B.2. We refer to our supplementary video for better visualization of these results.459
460
461
462
463
464
465
466
467
468
469
470

useful handle to control temporal consistency (Khachatryan et al., 2023). In addition to Chang et al. (2024) which this work was inspired by and improves upon, there are various other temporal noise manipulation techniques that do not preserve Gaussian noise distribution— some methods (Ma et al. (2024); Ren et al. (2024)) blend high frequency Gaussian noise with low frequency motion, while others (Mokady et al. (2022); Wallace et al. (2022)) rely on approximating the inversion of noise from temporally coherent image sequences. Pandey et al. (2024) goes one step further and manipulates inverted noise in 3D space. These approaches are flexible but degrade the output of the diffusion model due to the domain gap between inference time noise and training time noise, and as such, have occasionally been accompanied by mitigation strategies such as anisotropic diffusion (Yu et al. (2024)). Noise manipulation is also not limited to the generation and stylization of videos, but has various applications in image editing (Hou et al. (2024); Pandey et al. (2024)) and 3D mesh texturing (Richardson et al. (2023)) as well.

471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Noise in Computer Graphics. While our noise warping work draws main inspiration from simulation techniques, spatial noise manipulation has been extensively studied in the graphics community through applications in animation and rendering. Works like (Kass & Pesare, 2011; Burley et al., 2024) present 2D noise manipulation techniques that add a stylized organic hand-drawn look to computer-generated animation via dynamic noise textures (Perlin, 1985). In order to make sure the stylization is temporally consistent and visually pleasing, noise textures are deformed in a way that makes them consistent with the underlying animation, but little emphasis is given to the preservation/rigor of the noise distribution. On the other hand, properties of 2D spatial noise have been extensively and rigorously studied in rasterization and raytracing literature (Cook, 1986; Lagae & Dutré, 2008), originating from the idea of using dithering to reduce banding and quantization artifacts in image signal processing (Roberts, 1962). In particular, the lack of low frequency details and clumping in blue noise as opposed to white Gaussian noise has made it the choice of foundational antialiasing methods such as Poisson disc sampling (McCool & Fiume, 1992), and recent progress made in this line of antialiasing research has close ties with our methodology. For example, Wolfe et al. (2022) look at accelerating rendering tasks by extending spatial blue noise to the temporal domain, while Huang et al. (2024) show promising results in supplementing white noise with blue noise during diffusion model training.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

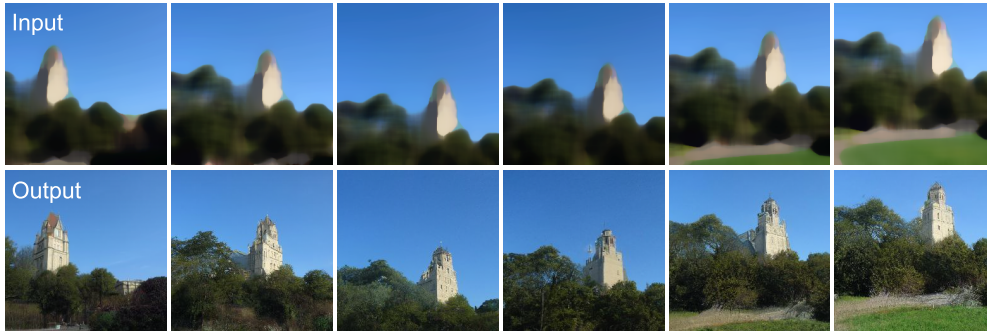


Figure 10: Conditional video generation results by integrating our method (particle) with SDEdit and cross-frame attention (Ceylan et al., 2023). Top row shows the input video prompt in a stroke painted style which is converted into a video of photorealistic style (bottom). Additional results are shown in Figure B.3. We refer to our supplementary video for better visualization of these results.

5 DISCUSSION AND CONCLUSION

We present infinite-resolution integral noise warping (Algorithm 1), a novel and efficient algorithm for warping a prior noise image into a sequence of noise frames while preserving the Gaussianity and temporal consistency. This is achieved by a theoretically motivated analysis of the infinite-resolution case of the integral noise warping algorithm by Chang et al. (2024), which enables an orders-of-magnitude improvement in efficiency with no trade-offs in accuracy.

Usability of Noise Warping We highlight that the noise warping problem that we address is a recurring subtask in generative modeling, and our method is hence a general-purpose tool that can be integrated in a variety of ways that extend well beyond the ones we showcase in the paper. First, noise warping, which excels at controlling high-frequency details, is orthogonal and thus combinable with *feature-level*, structure-preserving techniques (e.g. Ceylan et al. (2023); Cong et al. (2023)) to achieve consistency across the frequency spectrum. Our drastic cost-saving makes noise warping an affordable and harm-free add-on to all such existing and future techniques. In addition, the concurrent work by Daras et al. (2024) shows that noise warping can be combined with equivariance guidance to gain further consistency and integrate with latent diffusion models like SDXL (Podell et al., 2023). Beyond video generation, Kwak et al. (2024) showcases the usefulness of noise warping in 3D generation by combining with score distillation sampling (SDS). The advanced noise warping algorithm that we propose presents itself as a desirable candidate across these diverse tasks.

Significance of Our Speed-up We argue that the drastic speed-up our method offers has profound practical significance. While the standard denoising diffusion setting requires only a single noise warp operation per image, there exist many use cases that require noise warping to be computed more intensively, which renders our speed-up critical. For example, the combination with bridge models (e.g. I²SB) requires one noise warp per iteration. With its reported $> 0.6s$ time cost per warp, preparing the noise using HIWYN would cost $\sim 4\times$ the time to actually run the image generation model, increasing the total inference time from ~ 24 minutes to ~ 2 hours. In comparison, our method (particle) prepares the noise in 40.6s (wall time), effectively making the overhead negligible. Similarly, combining noise warping with SDS also requires one noise warp per iteration, which makes HIWYN computationally intractable (Kwak et al., 2024) and our improvements called for. Our speed-up hence makes integral noise warping deployable in a much broader class of problems.

We note some directions for future work. Since our particle-based variant does not leverage the deformation gradient of ψ , it does not account for area contraction and expansion. Voronoi re-partitioning may address this problem at the cost of extra computation. More broadly, since our method relies on the consistency of the deformation map and its alignment with the conditional video, it can be limited by the availability and accuracy of flow extraction techniques. Finally, the effectiveness of warped volumetric noise in 3D generation and editing tasks remains to be studied.

REFERENCES

- 540
541
542 Donghoon Ahn, Hyoungwon Cho, Jaewon Min, Wooseok Jang, Jungwoo Kim, SeonHwa Kim,
543 Hyun Hee Park, Kyong Hwan Jin, and Seungryong Kim. Self-rectifying diffusion sampling with
544 perturbed-attention guidance. *arXiv preprint arXiv:2403.17377*, 2024.
- 545 Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler,
546 and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion
547 models. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- 548 Jeremiah U Brackbill, Douglas B Kothe, and Hans M Ruppel. Flip: a low-dissipation, particle-in-
549 cell method for fluid flow. *Computer Physics Communications*, 48(1):25–38, 1988.
- 550
551 T. Brox and J. Malik. Large displacement optical flow: descriptor matching in variational motion
552 estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(3):500–513,
553 2011. URL [http://lmb.informatik.uni-freiburg.de/Publications/2011/](http://lmb.informatik.uni-freiburg.de/Publications/2011/Brolla)
554 [Brolla](http://lmb.informatik.uni-freiburg.de/Publications/2011/Brolla).
- 555 Brent Burley, Brian Green, and Daniel Teece. Dynamic screen space textures for coherent styliza-
556 tion. In *ACM SIGGRAPH 2024 Talks, SIGGRAPH '24*, New York, NY, USA, 2024. Association
557 for Computing Machinery. ISBN 9798400705151. doi: 10.1145/3641233.3664321. URL
558 <https://doi.org/10.1145/3641233.3664321>.
- 559
560 Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. Pix2video: Video editing using image
561 diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
562 23206–23217, 2023.
- 563 Pascal Chang, Jingwei Tang, Markus Gross, and Vinicius C Azevedo. How i warped your noise: a
564 temporally-correlated noise prior for diffusion models. In *The Twelfth International Conference*
565 *on Learning Representations*, 2024.
- 566 Weifeng Chen, Yatai Ji, Jie Wu, Hefeng Wu, Pan Xie, Jiashi Li, Xin Xia, Xuefeng Xiao, and Liang
567 Lin. Control-a-video: Controllable text-to-video generation with diffusion models. *arXiv preprint*
568 *arXiv:2305.13840*, 2023.
- 569
570 Yuren Cong, Mengmeng Xu, Christian Simon, Shoufa Chen, Jiawei Ren, Yanping Xie, Juan-Manuel
571 Perez-Rua, Bodo Rosenhahn, Tao Xiang, and Sen He. Flatten: optical flow-guided attention for
572 consistent text-to-video editing. *arXiv preprint arXiv:2310.05922*, 2023.
- 573 Robert L Cook. Stochastic sampling in computer graphics. *ACM Transactions on Graphics (TOG)*,
574 5(1):51–72, 1986.
- 575
576 Giannis Daras, Weili Nie, Karsten Kreis, Alex Dimakis, Morteza Mardani, Nikola Borislavov Ko-
577 vachki, and Arash Vahdat. Warped diffusion: Solving video inverse problems with image diffu-
578 sion models. *arXiv preprint arXiv:2410.16152*, 2024.
- 579 Rick Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019.
- 580
581 Songwei Ge, Seungjun Nah, Guilin Liu, Tyler Poon, Andrew Tao, Bryan Catanzaro, David Jacobs,
582 Jia-Bin Huang, Ming-Yu Liu, and Yogesh Balaji. Preserve your own correlation: A noise prior for
583 video diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer*
584 *Vision*, pp. 22930–22941, 2023.
- 585 Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh
586 Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image dif-
587 fusion models without specific tuning. *International Conference on Learning Representations*,
588 2024.
- 589 Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang,
590 and José Lezama. Photorealistic video generation with diffusion models. *arXiv preprint*
591 *arXiv:2312.06662*, 2023.
- 592
593 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
neural information processing systems, 33:6840–6851, 2020.

- 594 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
595 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
596 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- 597
598 Xingzhong Hou, Boxiao Liu, Yi Zhang, Jihao Liu, Yu Liu, and Haihang You. Easydrag: Efficient
599 point-based manipulation on diffusion models. In *Proceedings of the IEEE/CVF Conference on*
600 *Computer Vision and Pattern Recognition (CVPR)*, pp. 8404–8413, June 2024.
- 601 Xingchang Huang, Corentin Salaun, Cristina Vasconcelos, Christian Theobalt, Cengiz Oztireli, and
602 Gurprit Singh. Blue noise for diffusion models. In *ACM SIGGRAPH 2024 Conference Papers*,
603 pp. 1–11, 2024.
- 604 Michael Kass and Davide Pesare. Coherent noise for non-photorealistic rendering. *ACM Transac-*
605 *tions on Graphics (TOG)*, 30(4):1–6, 2011.
- 606
607 Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang
608 Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models
609 are zero-shot video generators. In *Proceedings of the IEEE/CVF International Conference on*
610 *Computer Vision*, pp. 15954–15964, 2023.
- 611 Min-Seop Kwak, Donghoon Ahn, Ines Hyeonsu Kim, Jin-wha Kim, and Seungryong Kim.
612 Geometry-aware score distillation via 3d consistent noising and gradient consistency modeling.
613 *arXiv preprint arXiv:2406.16695*, 2024.
- 614
615 Ares Lagae and Philip Dutré. A comparison of methods for generating poisson disk distributions.
616 In *Computer Graphics Forum*, volume 27, pp. 114–129. Wiley Online Library, 2008.
- 617 Guan-Horng Liu, Arash Vahdat, De-An Huang, Evangelos A Theodorou, Weili Nie, and Anima
618 Anandkumar. I²sb: Image-to-image schrödinger bridge. *arXiv preprint arXiv:2302.05872*, 2023.
- 619
620 Xin Ma, Yaohui Wang, Gengyu Jia, Xinyuan Chen, Yuan-Fang Li, Cunjian Chen, and Yu Qiao. Cin-
621 emo: Consistent and controllable image animation with motion diffusion models. *arXiv preprint*
622 *arXiv:2407.15642*, 2024.
- 623 Michael McCool and Eugene Fiume. Hierarchical poisson disk sampling distributions. In *Graphics*
624 *interface*, volume 92, pp. 94–105, 1992.
- 625
626 Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon.
627 Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint*
628 *arXiv:2108.01073*, 2021.
- 629
630 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for
631 editing real images using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022.
- 632
633 Peter Mörters and Yuval Peres. *Brownian motion*, volume 30. Cambridge University Press, 2010.
- 634
635 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
636 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
637 text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 638
639 Bernt Oksendal. *Stochastic differential equations: an introduction with applications*. Springer
640 Science & Business Media, 2013.
- 641
642 Jan K Pachl. Disintegration and compact measures. *Mathematica Scandinavica*, pp. 157–168, 1978.
- 643
644 Karran Pandey, Paul Guerrero, Matheus Gadelha, Yannick Hold-Geoffroy, Karan Singh, and Niloy J
645 Mitra. Diffusion handles enabling 3d edits for diffusion models by lifting activations to 3d. In *Pro-*
646 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7695–
647 7704, 2024.
- 648
649 Ken Perlin. An image synthesizer. *ACM Siggraph Computer Graphics*, 19(3):287–296, 1985.
- 650
651 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
652 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
653 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.

- 648 Weiming Ren, Harry Yang, Ge Zhang, Cong Wei, Xinrun Du, Stephen Huang, and Wenhu
649 Chen. Consisti2v: Enhancing visual consistency for image-to-video generation. *arXiv preprint*
650 *arXiv:2402.04324*, 2024.
- 651
- 652 Elad Richardson, Gal Metzger, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-
653 guided texturing of 3d shapes. In *ACM SIGGRAPH 2023*, 2023.
- 654
- 655 L. Roberts. Picture coding using pseudo-random noise. *IRE Transactions on Information Theory*, 8
656 (2):145–154, 1962. doi: 10.1109/TIT.1962.1057702.
- 657
- 658 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
659 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
660 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 661
- 662 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry
663 Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video
664 data. *arXiv preprint arXiv:2209.14792*, 2022.
- 665
- 666 Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer*
667 *Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings,*
668 *Part II 16*, pp. 402–419. Springer, 2020.
- 669
- 670 Bram Wallace, Akash Gokul, and Nikhil Naik. Edict: Exact diffusion inversion via coupled trans-
671 formations, 2022. URL <https://arxiv.org/abs/2211.12446>.
- 672
- 673 Alan Wolfe, Nathan Morrical, Tomas Akenine-Möller, and Ravi Ramamoorthi. Spatiotemporal Blue
674 Noise Masks. In Abhijeet Ghosh and Li-Yi Wei (eds.), *Eurographics Symposium on Rendering*.
675 The Eurographics Association, 2022. ISBN 978-3-03868-187-8. doi: 10.2312/sr.20221161.
- 676
- 677 Xi Yu, Xiang Gu, Haozhi Liu, and Jian Sun. Constructing non-isotropic gaussian diffusion model
678 using isotropic gaussian diffusion model for image editing. *Advances in Neural Information*
679 *Processing Systems*, 36, 2024.
- 680
- 681 Bowen Zhang, Yiji Cheng, Jiaolong Yang, Chunyu Wang, Feng Zhao, Yansong Tang, Dong Chen,
682 and Baining Guo. Gaussiancube: Structuring gaussian splatting using optimal transport for 3d
683 generative modeling. *arXiv preprint arXiv:2403.19655*, 2024.
- 684
- 685 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
686 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
687 pp. 3836–3847, 2023a.
- 688
- 689 Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. Con-
690 trolvideo: Training-free controllable text-to-video generation. *arXiv preprint arXiv:2305.13077*,
691 2023b.

690 A PROOF OF THEOREM 1

692 *Proof.* By unrolling the definitions, for $t \in [0, 1]$, we have

$$693$$

$$694 H_n(t) = S_n^*(t) - tS_n^*(1) + tc, \quad S_n^*(t) := \frac{1}{\sqrt{n}} \left(\sum_{i=1}^{\lfloor nt \rfloor} Z_i + (nt - \lfloor nt \rfloor)Z_{\lfloor nt \rfloor + 1} \right).$$

695

696

697

698 By Mörters & Peres (2010, Theorem 5.22), $\{S_n^*\}_{n \in \mathbb{Z}_{\geq 1}}$ converges in distribution to $W(t)$ under the
699 sup-norm metric of $C[0, 1]$. To lift this convergence to the sequence $\{H_n\}_{n \in \mathbb{Z}_{\geq 1}}$, observe that the
700 function $g : C[0, 1] \rightarrow C[0, 1]$ defined by

$$701 g(x(t)) := x(t) - tx(1) + tc$$

is continuous under the sup-norm metric. To verify this, suppose $\lim_{n \rightarrow \infty} f_n = f$ for $\{f_n\}_{n \in \mathbb{Z}_{\geq 1}}, f \in C[0, 1]$. Then

$$\begin{aligned} \|g(f_n) - g(f)\|_\infty &= \sup_{t \in [0, 1]} |(f_n(t) - tf_n(1) + tc) - (f(t) - tf(1) + tc)| \\ &\leq \|f_n - f\|_\infty + \|f_n(1) - f(1)\| \leq 2\|f_n - f\|_\infty \rightarrow 0. \end{aligned}$$

Hence, by the continuous mapping theorem,

$$g(S_n^*) = H_n \xrightarrow{d} B(t) - tB(1) + tc.$$

To show

$$W(t) - tW(1) + tc = (W(t) \mid W(1) = c),$$

first of all, the conditioning $(W(t) \mid W(1) = c)$ is interpreted as the limit of $(W(t) \mid |W(1) - c| < \epsilon)$ as $\epsilon \rightarrow 0$. Denote $Y(t) := W(t) - tW(1)$, so that $W(t) = Y(t) + tW(1)$. Since $\text{Cov}(Y(t), tW(1)) = \text{Cov}(W(t) - tW(1), tW(1)) = t\text{Cov}(W(t), W(1)) - t^2\text{Var}(W(1), W(1)) = 0$ and that $Y(t), tW(1)$ are jointly Gaussian, they are independent. Therefore,

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} (W(t) \mid |W(1) - c| < \epsilon) &= \lim_{\epsilon \rightarrow 0} (Y(t) + tW(1) \mid |W(1) - c| < \epsilon) \\ &= Y(t) + \lim_{\epsilon \rightarrow 0} (tW(1) \mid |W(1) - c| < \epsilon) \\ &= W(t) - tW(1) + tc. \end{aligned}$$

□

B ADDITIONAL RESULTS

In this section we include additional visual and numerical results. In Figure B.8 and Figure B.9, we showcase full comparisons with the addition of Control-A-Video (Chen et al., 2023) and PYoCo (Ge et al., 2023), along with baselines with fixed noise and interpolated noise using bilinear and nearest interpolating schemes. The corresponding quantitative metrics for both church and cat scenes are reported in Figure B.1. In Figure B.6, we use additional examples to showcase the area contention issue caused by degenerate meshes that applies similarly to our grid-based variant and Chang et al. (2024), and highlight the robustness of our particle-based variant. In Figure B.5, we show additional results when combining our particle-based 3D noise warp with GaussianCube (Zhang et al., 2024).

Video Generation Quality (Church)									
Metric	Ours (G)	Ours (P)	HIWYN	PYoCo	CaV	Random	Fixed	Bilinear	Nearest
Consistency ↓	9.868e-2	1.065e-1	1.060e-1	1.175e-1	1.359e-1	1.538e-1	1.120e-1	8.114e-2	1.305e-1
Realism ↓	4.643e-2	5.180e-2	4.959e-2	4.119e-2	4.069e-2	3.731e-2	3.911e-2	2.301e-1	7.012e-2
Faithfulness ↓	3.872e-2	4.309e-2	4.377e-2	3.764e-2	4.169e-2	3.976e-2	3.264e-2	5.623e-2	9.321e-2
Video Generation Quality (Cat)									
Metric	Ours (G)	Ours (P)	HIWYN	PYoCo	CaV	Random	Fixed	Bilinear	Nearest
Consistency ↓	6.001e-2	5.898e-2	5.807e-2	6.383e-2	4.280e-2	1.219e-1	3.950e-2	3.503e-2	1.058e-1
Realism ↓	1.559e-1	1.496e-1	1.528e-1	1.506e-1	1.486e-1	1.221e-1	1.588e-1	3.687e-1	3.343e-1
Faithfulness ↓	2.039e-2	2.064e-2	2.022e-2	2.023e-2	1.817e-2	2.077e-2	1.972e-2	3.809e-2	2.201e-1

Figure B.1: We show the quality metrics for conditional video generating using SDEdit. The consistency is measured using Warp MSE, and the realism and faithfulness are measured as in Meng et al. (2021).

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

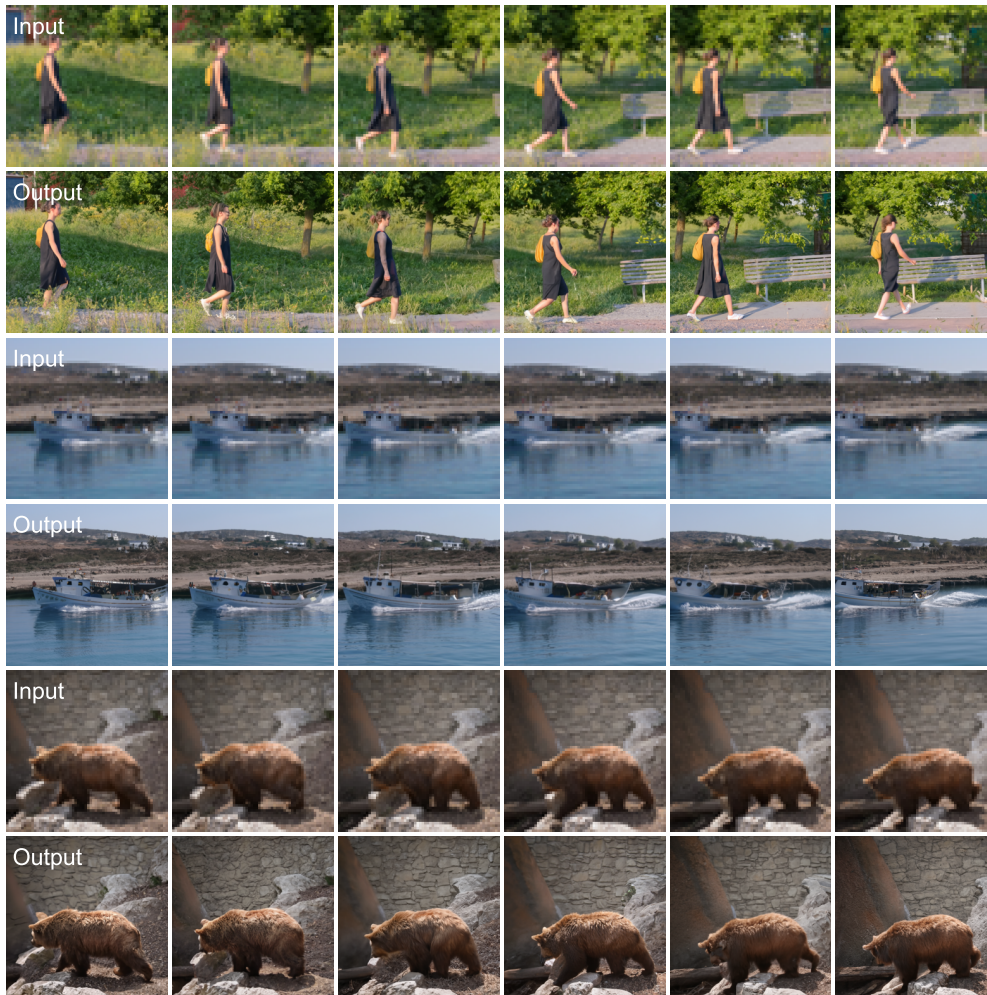


Figure B.2: Additional results generated by performing 4x video super-resolution with I²SB. For each scenario, the upper row represents the low resolution input video, and the lower row represents the high resolution output video. We refer to our supplementary video for better visualization of these results.

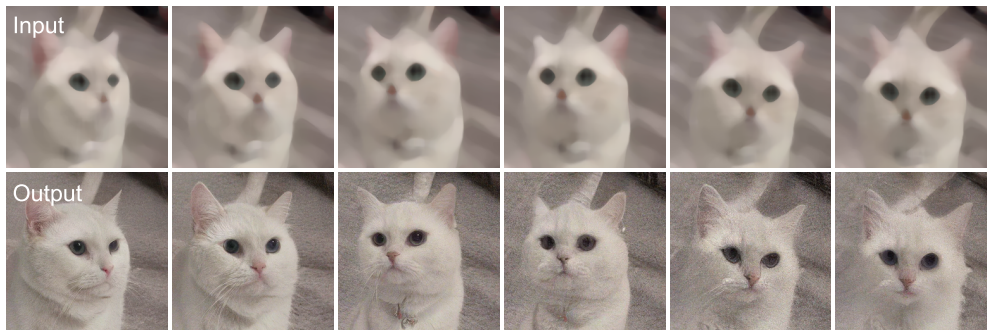


Figure B.3: Additional conditional video generation results by integrating our method (particle) with SDEdit and cross-frame attention (Ceylan et al., 2023). The top row shows the input video prompt in a stroke painted style which is converted into a video of photorealistic style (bottom). We refer to our supplementary video for better visualization of these results.

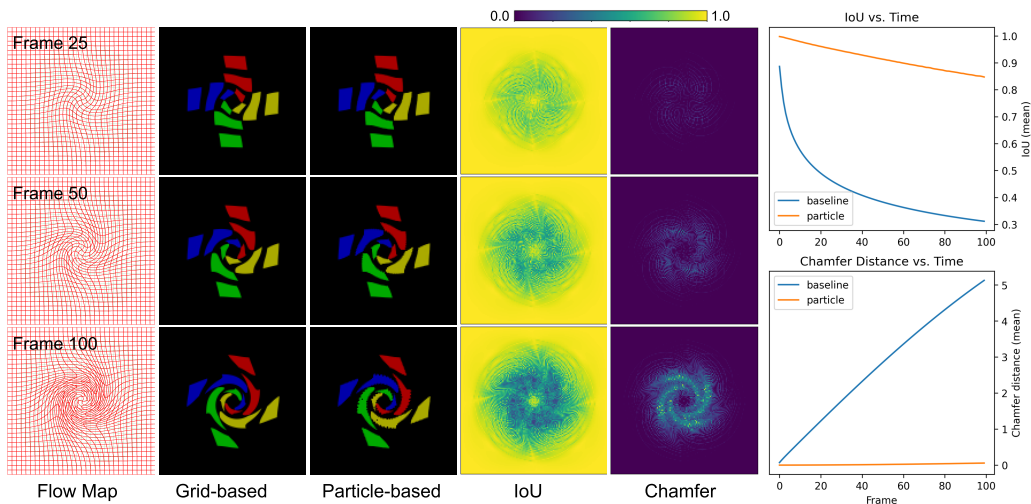


Figure B.4: Comparison between the grid-based and particle-based variants for building partition records when the deformation map is diffeomorphic. The first column shows the deformation map at different frames. The second and third columns visualize warped pixel regions for the two methods. The right two columns show IoU (larger is better) and Chamfer distance (smaller is better) between the outputs from both variants. We plot the distance between particle and grid variants alongside a baseline, which is the distance between identity map and the grid variant, which show that the particle-based version remains close to the grid-based version even under large distortion.

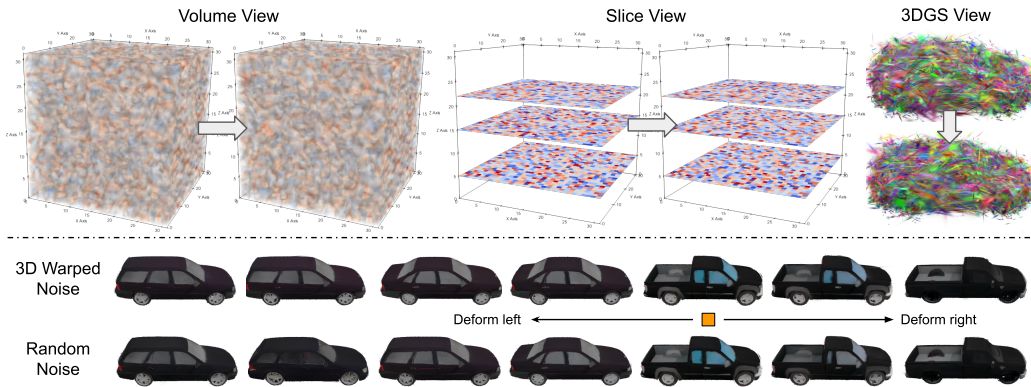


Figure B.5: Extension of our particle-based method to warping 3D noise. We show the volume render on the top left, the slice view on the top middle, and represent it as 3D Gaussian as used in GaussianCube representation on the top right. We then show that warping the noise in 3D space noticeably facilitates temporal consistency over random baseline when we perform 3D editing, which can be observed from the flickering of the color of the window in the bottom row. For best viewing of this experiment, please refer to our supplementary video.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

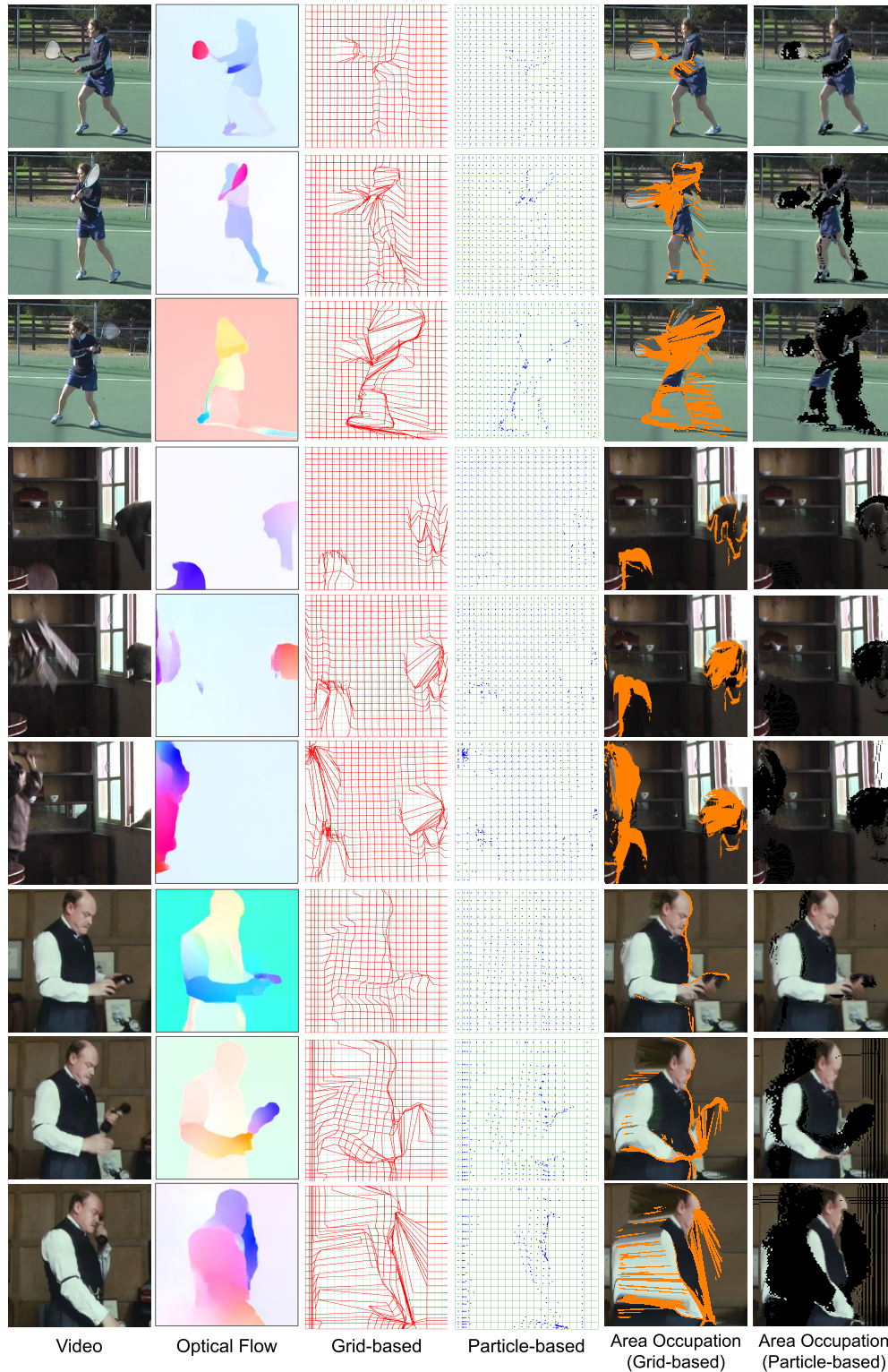


Figure B.6: Comparison of grid-based and particle-based variants under non-diffeomorphic deformation maps seen in real-world scenarios. The orange pixels are the invalid pixels where area overlap occurs. Flow maps are downsampled 10× for better visualization. Image sequence comes from Brox & Malik (2011) while optical flow is computed via Teed & Deng (2020).

918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

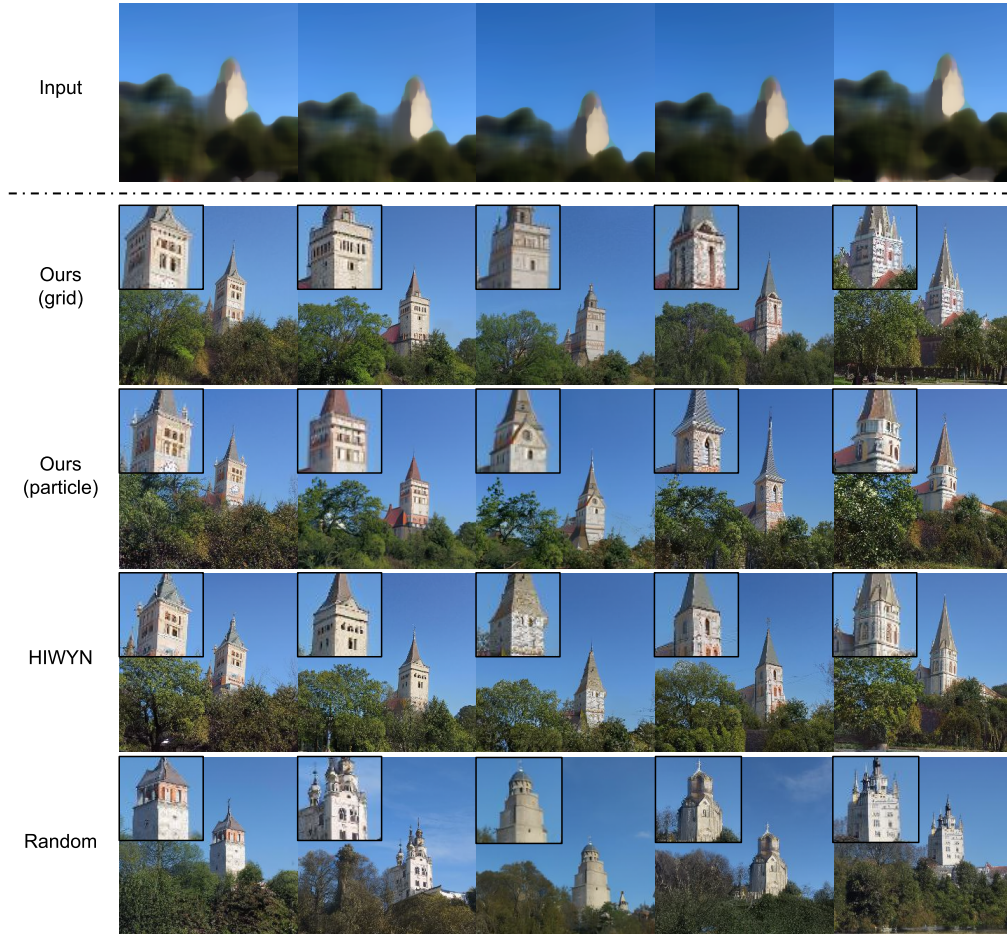


Figure B.7: We compare the consistency-preserving efficacy with different noise warp schemes. We use SDEdit using the conditional signal on the top image. We show here our method (both variants) and HIWYN, while further results are given in Figure B.8. We highlight the details of the tower, which is preserved to similar extents by the integral noise-based methods.

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

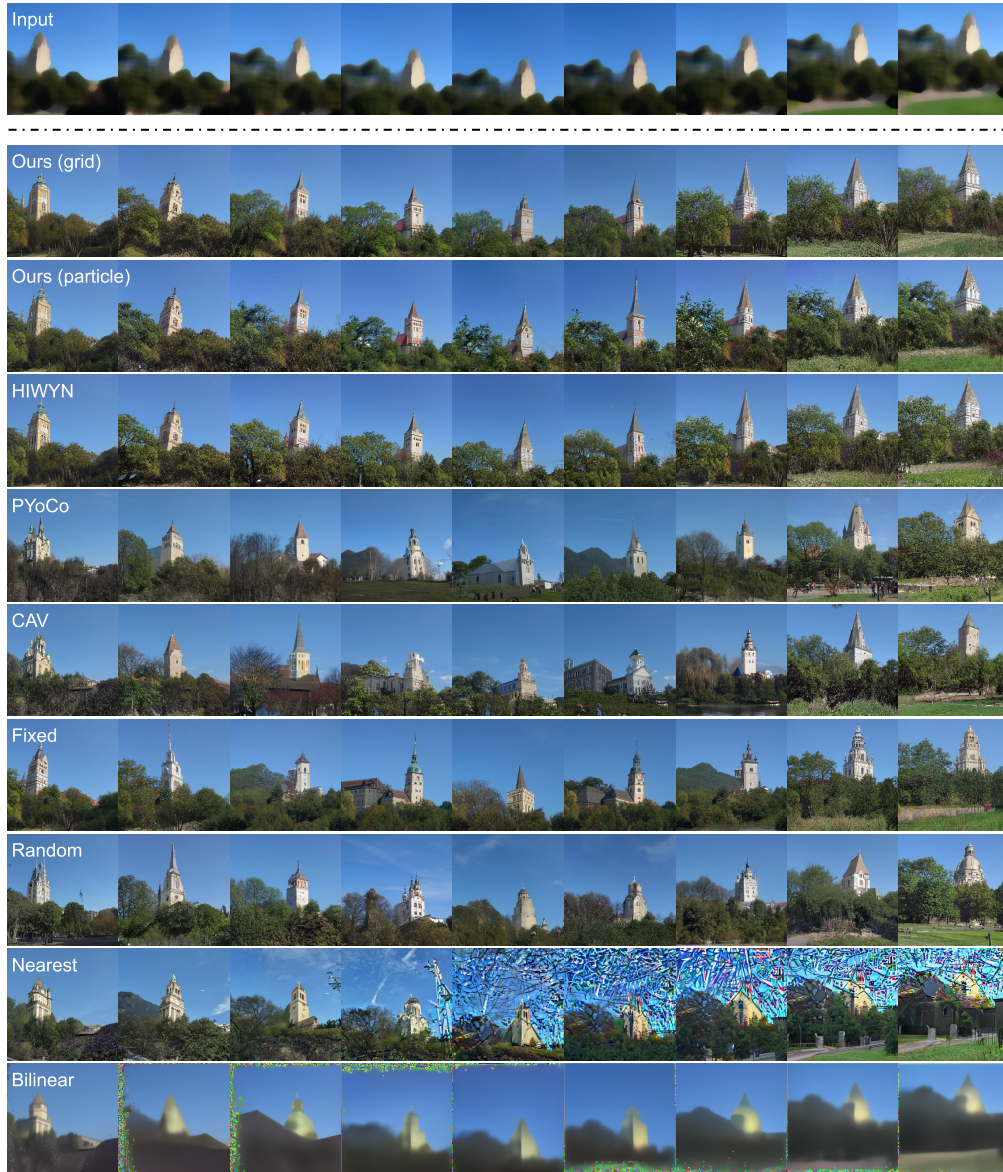


Figure B.8: The full results generated by all the compared methods on the church scene. The interpolation baselines yield noticeably corrupted results, but other yield similar quality *on the images level*. The difference lies in how the details are preserved across frames. Apart from the details of the main tower, the tree on the bottom left also exposes the interesting difference between noise initialization schemes. We refer to our supplementary video for better visualization of these results.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

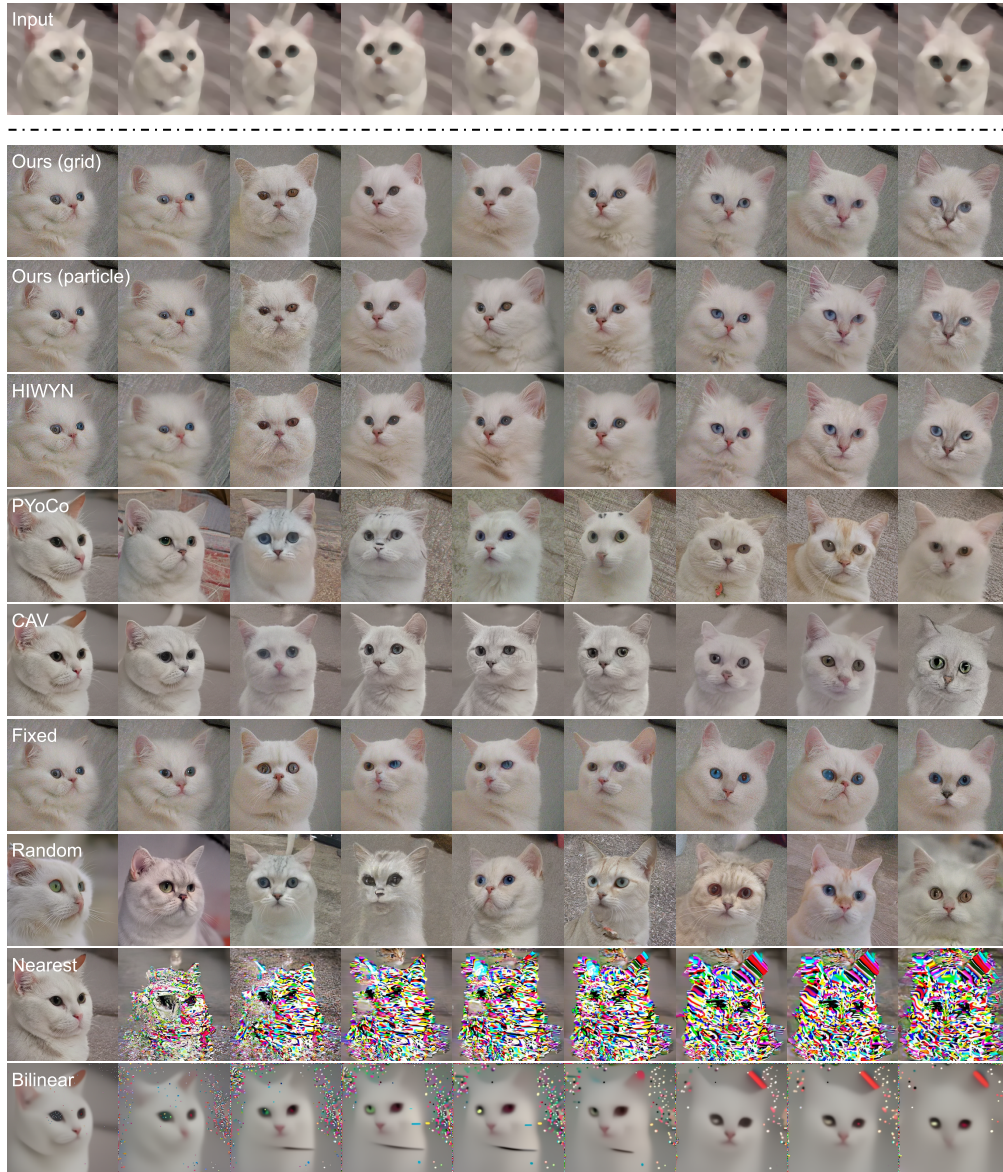


Figure B.9: The full results generated by all the compared methods on the cat scene. We observe that HIWYN and our methods (both variants) yield highly similar results. While both our variants are much faster and memory-efficient than HIWYN, this observation makes a particularly strong case for our particle-based variant considering its significantly improved simplicity and efficiency. We refer to our supplementary video for better visualization of these results.