

SphereDiff: Tuning-free Omnidirectional Panoramic Image and Video Generation via Spherical Latent Representation

Minho Park* Taewoong Kang* Jooyeol Yun Sungwon Hwang Jaegul Choo
 Korea Advanced Institute of Science and Technology (KAIST)
 {m.park, keh0t0, jchoo}@kaist.ac.kr

"Fireworks, City, Nightscape, Buildings, River, etc."

Figure 1. 360-degree panoramic video generated by *SphereDiff*. Click to play the animation clips. Best viewed with Acrobat Reader.

Abstract

The increasing demand for AR/VR applications has highlighted the need for high-quality 360-degree panoramic content. However, generating high-quality 360-degree panoramic images and videos remains a challenging task due to the severe distortions introduced by equirectangular projection (ERP). Existing approaches either fine-tune pre-trained diffusion models on limited ERP datasets or attempt tuning-free methods that still rely on ERP latent representations, leading to discontinuities near the poles. In this paper, we introduce *SphereDiff*, a novel approach for seamless 360-degree panoramic image and video generation using state-of-the-art diffusion models without additional tuning.

We define a spherical latent representation that ensures uniform distribution across all perspectives, mitigating the distortions inherent in ERP. We extend *MultiDiffusion* to spherical latent space and propose a spherical latent sampling method to enable direct use of pretrained diffusion models. Moreover, we introduce distortion-aware weighted averaging to further improve the generation quality in the projection process. Our method outperforms existing approaches in generating 360-degree panoramic content while maintaining high fidelity, making it a robust solution for immersive AR/VR applications. The code is available [here](#).

1. Introduction

The growing demand for AR/VR applications has significantly increased the need for high-quality immersive con-

* indicates equal contributions.

tent. AR/VR technologies offer highly engaging environments, providing a sense of presence that traditional displays (*e.g.*, phones and laptops) cannot. A key element in delivering such experiences is the $360^\circ \times 180^\circ$ panoramic scene, or 360-degree panorama, which provides an omnidirectional view of the virtual world. This allows users to explore their surroundings from any perspective, setting it apart from standard visual content. However, as 360-degree panoramas require specialized cameras, their availability is limited, and VR users have yet to experience a broad range of realistic content beyond simulations.

Recently, diffusion models have demonstrated remarkable performance in generating standard images and videos [9, 13, 21, 26]. Given their success, there is a growing interest in synthesizing 360-degree panoramic images or videos by leveraging the recent state-of-the-art diffusion models. 360-degree panorama is typically represented using an equirectangular projection (ERP), which maps spherical imagery onto a 2D rectangular plane, *e.g.*, the projection of a globe onto a world map. Due to the limitations of the 2D rectangular representation, the ERP introduces significant distortion, known as ERP distortion, in which high-latitude regions appear disproportionately large due to projection effects. For example, as shown in Fig. 1, the fireworks near the pole appear significantly larger than the others because our generated panoramic video is visualized in ERP. This ERP distortion causes a significant distribution shift in 360-degree panorama compared to standard perspective images or videos, making it challenging for standard diffusion models to generate omnidirectional panoramic content.

To mitigate the distribution shift, several previous studies have fine-tuned pretrained diffusion models using ERP datasets [5, 16, 24, 30]. However, due to the limited availability of text-ERP pairs, they often failed to generate seamless 360-degree panoramas, particularly near the poles, as shown in Fig. 2. On the other hand, some studies [2, 18] have proposed methods for generating arbitrarily sized panoramas without tuning pretrained models based on the MultiDiffusion framework [2]. Nevertheless, they are also limited to the ERP latent representation, which causes discontinuities near the poles.

In this paper, we present a novel framework, *SphereDiff*, which effectively generates 360-degree panoramic images and videos by leveraging recent state-of-the-art diffusion models without additional tuning. First, we define a spherical latent representation that pairs each latent with its corresponding position on the spherical surface. Then, we extend MultiDiffusion [2] to the spherical latent space. Since the spherical representation is equally distributed across all perspectives, we can handle every perspective uniformly, allowing for seamless image generation even near the poles. While the number of spherical latents is nearly equal across all perspectives, the projected spherical latents

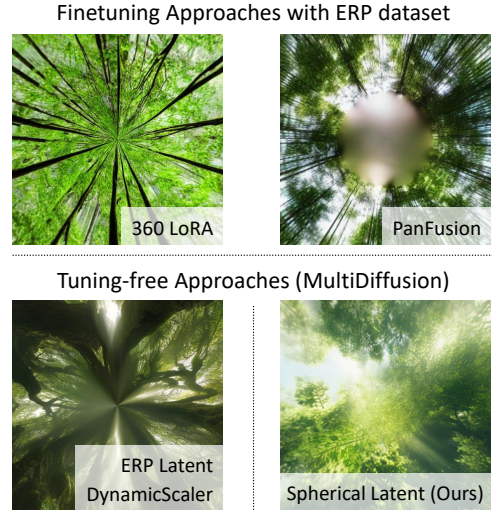


Figure 2. **Motivation.** Previous finetuning approaches [15, 30] often fail to generate continuous scenes near the pole due to the limited ERP dataset. The tuning-free approach [18] also fails to generate a seamless frame due to the ERP latent representation.

are distributed in continuous positions. Thus, we propose a novel spherical latent sampling method to discretize the spherical latents onto a 2D grid, enabling the use of readily available state-of-the-art diffusion models [9, 13, 14, 26]. Lastly, we introduce distortion-aware weighted averaging to further improve the minor distortions caused by spherical-to-perspective projection.

Our contributions can be summarized as follows:

- We propose *SphereDiff*, a tuning-free pipeline for high-quality and seamlessly continuous 360-degree panoramic image and video generation with minimal distortion.
- We introduce dynamic latent sampling that discretizes spherical latents onto a 2D grid, which enables the utilization of state-of-the-art diffusion models.
- We propose a distortion-aware weighted averaging technique that guarantees seamless content generation with superior visual quality.
- Extensive experiments demonstrate that *SphereDiff* highly outperforms existing methods in terms of visual quality and robustness to distortion.

2. Related Work

Latent Diffusion Models. Recent advancements in diffusion models have enabled the generation of high-quality images [4, 14, 21, 26] and videos [9, 13, 17, 23, 28, 31], achieving impressive visual results across various video generation tasks within the standard perspective of visual content. However, generating content beyond the standard perspective, such as regular or 360-degree panoramas, remains relatively underexplored. In this paper, we aim to generate 360-degree panoramas, which differ significantly from standard perspective scenes, by solely leveraging pre-

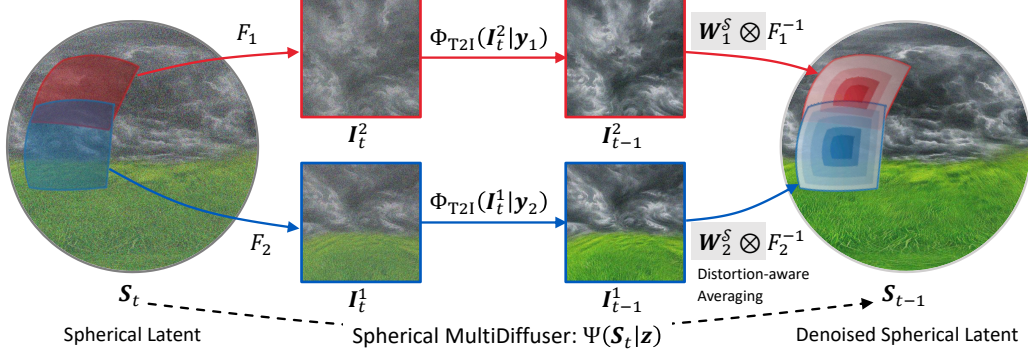


Figure 3. **Overall Pipeline.** We initialize uniform spherical latents and extract perspective latents for multiple views at each denoising step using dynamic latent sampling. These latents are then denoised and fused using the MultiDiffusion [2] with distortion-aware weighted averaging. This process enables seamless and distortion-free 360-degree panoramic image and video generation in a tuning-free manner.

trained diffusion models designed for standard perspectives. Furthermore, as our method is tuning-free, it is highly adaptable and can be seamlessly integrated with state-of-the-art diffusion-based image and video generation models.

360-degree Panoramic Scene Generation. Most panoramic generation methods rely on equirectangular projection (ERP), which maps spherical panorama coordinates onto a 2D rectangular plane, where latitude and longitude correspond to vertical and horizontal coordinates, respectively. However, ERP inherently introduces severe nonlinear distortions, particularly near the poles, degrading both visual quality and the processing of equirectangular panoramas. Although previous studies [5, 7, 16, 24, 25, 30, 32] attempt to address this issue by fine-tuning on panoramic ERP datasets, they often fail to generate seamless panoramas, especially near the poles, or struggle with text controllability due to the domain-specific nature of these datasets (e.g., indoor environments). CubeDiff [12] introduces an alternative approach using cube map representations for panoramic image generation. While this method effectively reduces distortions near the poles, it still struggles with discontinuities at cube-face boundaries. Another line of research [6, 16, 29, 32] focuses on scene-level generation, primarily producing static 3D scenes using Gaussian-based representations. While these methods achieve high-quality static outputs, they require optimization procedures that reduce efficiency. In contrast, we replace the ERP latent representation with a spherical latent representation, providing a natural solution to eliminate distortion across all perspectives.

Tuning-free Panorama Generation. For 360-degree panoramic video generation, far fewer datasets exist compared to ERP images. As a result, recent research trends favor utilizing perspective models without additional training. DynamicScaler [18] attempts to mitigate ERP’s inher-

ent distortions by employing panoramic-projected denoising, leveraging the MultiDiffusion framework [2] with adjusted windows. On the other hand, 4K4DGen [16] seeks to avoid distortion by utilizing ERP images with an image-to-video model, limiting its applicability to more complex dynamic content generation. However, both methods struggle to generate seamless scenes near the poles due to severe interpolation or sampling artifacts. In contrast, our method overcomes these limitations by directly utilizing uniformly distributed spherical latents, ensuring both efficiency and reduced distortion without requiring additional training.

3. Method

In this section, we introduce *SphereDiff*, a novel tuning-free framework for generating 360-degree panoramic images and videos. First, we present the spherical latent representation and spherical-to-perspective projection in Section 3.1. Next, we extend the MultiDiffusion framework [2] to the spherical latent space in Section 3.2. We then introduce spherical latent sampling, which discretizes the continuous coordinates of the perspective-projected spherical latent onto a 2D grid in Section 3.3. Finally, we propose a distortion-aware weighted averaging method to mitigate minor distortions from the spherical-to-perspective projection in Section 3.4. The overall pipeline is illustrated in Figure 3.

3.1. Spherical Latent Representation

Definition. We introduce a spherical representation of latent features for generating 360-degree panoramas. We define a latent feature $\mathbf{f} \in \mathbb{R}^C$ paired with the corresponding spherical coordinate \mathbf{d} on a spherical surface. The set of spherical coordinates can be represented as follows:

$$\mathbb{S}^2 = \{\mathbf{d} = (x, y, z) \mid x, y, z \in \mathbb{R}, \|\mathbf{d}\| = 1\}. \quad (1)$$

Then, we pair each latent feature with its associated position, i.e., $s = (\mathbf{d}, \mathbf{f})$, referred to as *spherical latent*. For N

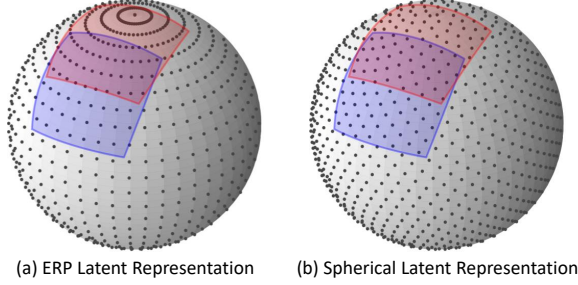


Figure 4. **Comparison of ERP and Spherical Latent Representations.** When changing perspective, the ERP latent representation shows significant density variations in latent density depending on position, especially near the poles, while our spherical representation maintains a nearly uniform density across all perspectives.

spherical latents, we define the spherical latents \mathcal{S} as:

$$\mathcal{S} = \{\mathbf{s}_i = (\mathbf{d}_i, \mathbf{f}_i) \mid \mathbf{d}_i \in \mathbb{S}^2, \mathbf{f}_i \in \mathbb{R}^C, \text{ for } i \in [1, N]\}. \quad (2)$$

which is now composed of multiple latents similar to standard 2D or 3D latent features. We refer to the domain of spherical latents as \mathcal{S} , i.e., $\mathbf{S} \in \mathcal{S}$.

Equirectangular Projection (ERP) latents also can be written in our spherical latent representation. However, as shown in Fig. 4 (a), due to the 2D grid constraint of ERP latent, its spherical coordinates are not uniformly distributed on the sphere’s surface. In contrast, we define the spherical latents using the Fibonacci Lattice [10], which offers the number of spherical latents is nearly equal across all perspectives, as shown in Fig. 4 (b).

Perspective Latent Representation. Since standard diffusion models operate in perspective space, we utilize a *spherical-to-perspective projection* which transforms the spherical coordinate to the perspective coordinate. To achieve this, we first define the domain of perspective coordinates as a discretized 2D plane as

$$\mathbb{P}^2 = \left\{ \mathbf{u} = \left(\frac{2j}{H}, \frac{2k}{W} \right) \mid j \in \left[-\frac{H}{2}, \frac{H}{2} \right], k \in \left[-\frac{W}{2}, \frac{W}{2} \right] \right\}, \quad (3)$$

where H, W indicates the height and width of the bounded 2D perspective plane, respectively. We use a view direction $\mathbf{v} \in \mathbb{S}^2$ and a predefined focal length f to define the spherical-to-perspective projection function $\mathbf{u} = \mathcal{T}_{\mathbb{S}^2 \rightarrow \mathbb{P}^2}(\mathbf{d}|\mathbf{v}, f)$. The detailed formula of the projection function is available in Section B.1.

3.2. MultiDiffusion for Spherical Latent.

The MultiDiffusion [2] framework is often utilized for generating arbitrary-shaped images by leveraging pretrained diffusion models [21] trained on standard perspective images. In this section, we introduce an extension of the MultiDiffusion framework to the spherical latent representation.

The goal of this framework is to construct the *Spherical MultiDiffuser* $\Psi : \mathcal{S} \times \mathcal{Z} \rightarrow \mathcal{S}$, which takes a noisy spherical latent \mathbf{S}_t and a set of text conditions \mathbf{z} as inputs and produces the denoised spherical latent \mathbf{S}_{t-1} , as illustrated in Figure 3. Based on the MultiDiffuser, a clean spherical latent \mathbf{S}_0 can be obtained from pure noise \mathbf{S}_T through an iterative denoising process using diffusion models as:

$$\mathbf{S}_T, \mathbf{S}_{T-1}, \dots, \mathbf{S}_0 \quad \text{s.t.} \quad \mathbf{S}_{t-1} = \Psi(\mathbf{S}_t|\mathbf{z}). \quad (4)$$

To construct MultiDiffuser Ψ , we first leverage a pretrained diffusion model trained on standard perspective latents $\Phi : \mathcal{I} \times \mathcal{Y} \rightarrow \mathcal{I}$, which takes a noisy latent \mathbf{I}_t and a text condition \mathbf{y} as inputs and produces the denoised latent \mathbf{I}_{t-1} . The pretrained diffusion model gradually denoises the pure Gaussian noise $\mathbf{I}_T \sim \mathcal{N}$ into a clean image \mathbf{I}_0 .

$$\mathbf{I}_T, \mathbf{I}_{T-1}, \dots, \mathbf{I}_0 \quad \text{s.t.} \quad \mathbf{I}_{t-1} = \Phi(\mathbf{I}_t|\mathbf{y}) \quad (5)$$

Next, we define a mapping function between the spherical and perspective latent spaces, $F_i : \mathcal{S} \rightarrow \mathcal{I}$, along with a corresponding condition mapping $\lambda_i : \mathcal{Z} \rightarrow \mathcal{Y}$, where $i \in \{1, \dots, n\}$. The mapping function F_i can be formulated in various ways, which will be discussed in Section 3.3.

$$\mathbf{I}_t^i = F_i(\mathbf{S}_t), \quad \mathbf{y}_i = \lambda_i(\mathbf{z}) \quad (6)$$

Finally, The denoising step in of MultiDiffuser can be formulated by a closed-form [2].

$$\Psi(\mathbf{S}_t|\mathbf{z}) = \sum_{i=1}^n \mathbf{W}_{\mathcal{S}}^i \otimes F_i^{-1}(\Phi(\mathbf{I}_t^i|\mathbf{y}_i)). \quad (7)$$

In the following sections, we define F_i (Section 3.3) and the per-pixel weight $\mathbf{W}_{\mathcal{S}}^i$ (Section 3.4) to ensure a seamless and consistent 360-degree panorama.

3.3. Sampling Spherical Latent

Based on the MultiDiffusion framework [2], we define the mapping function F that transforms a spherical latent representation into a perspective latent space \mathcal{I} . To define mapping function F , we first apply the transformation $\mathcal{T}_{\mathcal{S} \rightarrow \mathcal{I}}$ based on the view direction $\mathbf{v} \in \mathbb{S}^2$ and focal length f , which projects the coordinates of the spherical latents onto the perspective plane \mathcal{P} . Formally, the spherical-to-perspective latent transformation can be written as

$$\mathcal{T}_{\mathcal{S} \rightarrow \mathcal{P}}(\mathbf{S}|\mathbf{v}, f) = \mathbf{P} = \{\mathbf{p}_i = (\mathbf{u}_i, \mathbf{f}_i) | \mathbf{u}_i \in [-1, 1]^2\}, \quad (8)$$

where $\mathbf{u}_i = \mathcal{T}_{\mathbb{S}^2 \rightarrow \mathbb{P}^2}(\mathbf{d}_i|\mathbf{v}, f)$. Note that, \mathbf{P} does not contain N elements since the perspective are cropped by $[-1, 1]$.

Since visual diffusion models based on DiT [19] support continuous 1D representations, we explore a denoising process based on continuous coordinates with RoPE [9, 22]. However, this approach introduces unstructured outputs due

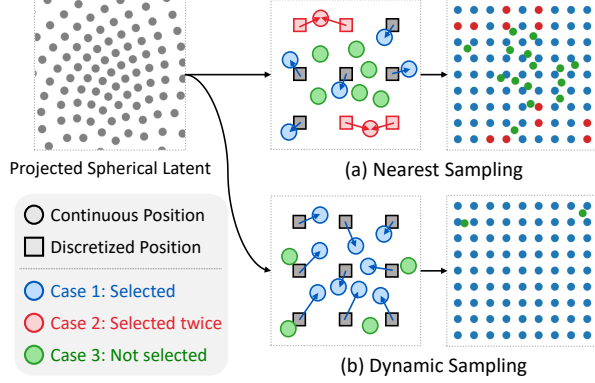


Figure 5. **Comparison of Nearest and Dynamic Sampling.** Nearest sampling often **resamples the selected latents** or **omits central ones**, while dynamic sampling selects latents from the center outward, discarding only the outermost ones.

to a distribution shift in positional embeddings, as detailed in Section B.2. To address this, we apply an additional discretization transformation \mathcal{D} that maps the continuous latent coordinates to \mathbb{P}^2 , i.e., $\mathbf{I} = \mathcal{D}(\mathcal{T}_{S \rightarrow \mathcal{I}}(S|\mathbf{v}, f))$. In the following sections, we introduce two simple yet effective methods for discretizing the perspective coordinates.

Nearest Point Sampling. A straightforward approach to discretizing continuous coordinates is nearest-neighbor sampling, where the nearest projected spherical latent is selected for each pixel position. Specifically, the latent closest to the center of a $H \times W$ grid is retrieved and used as input for denoising, as illustrated in Fig. 5 (a). Despite its simplicity, this method introduces two critical issues. First, the same latent may be selected multiple times, altering the latent distribution, which often degrades generation performance [3]. Second, some spherical latents may not be chosen even if they fall within the field of view of the current camera view direction. This phenomenon, referred to as the “undersampling problem”, has particularly detrimental consequences for generating a seamless panorama.

Undersampling Problem. The undersampling of spherical latents disrupts information flow across neighboring windows. As illustrated in Fig. 5, the green points lack information from the current field of view (FoV) since they are not denoised in this step. If the next window’s FoV captures a green point, not the blue point, it receives no information from the current window, causing discontinuities even when there is a large overlap. To address this, we propose a dynamic latent sampling algorithm that effectively collects projected spherical latents.

Dynamic Latent Sampling. As mentioned earlier, we aim to ensure that all points within the field of view (FoV) are selected, enabling seamless panorama generation. To

Algorithm 1 Dynamic Latent Sampling (in case of H, W are even and equal)

Input : Projected Latent $\mathbf{P} = \mathcal{T}_{S \rightarrow \mathcal{I}}(S|\mathbf{v}, f)$

Output: Arranged Perspective Latent \mathbf{I}

$\mathbf{I}' \leftarrow$ Sort the latents of \mathbf{P} by $\|\mathbf{u}_i\|$.

$M \leftarrow |\mathbf{P}|$ \triangleright Get the number of latents

$H, W \leftarrow \lfloor \sqrt{M} \rfloor, \lfloor \sqrt{M} \rfloor$ \triangleright Get smaller box

$\mathbf{I} \leftarrow \emptyset^{H \times W}$ \triangleright Initialize

for $i \in [1, H/2]$ **do**

$n \leftarrow (2i)^2 - (2i - 2)^2$ \triangleright count of i -th border

$\mathbf{l} \leftarrow$ first n latents from the sorted \mathbf{I}'

 Set i -th border of \mathbf{I} to \mathbf{l} \triangleright fill the center first

 Pop first n latents from \mathbf{I}'

end

return \mathbf{I} \triangleright Ignore $M - H \times W$ elements of \mathbf{I}

achieve this, we propose a novel dynamic latent sampling algorithm that first selects the center-positioned points and then ignores the remaining points at the outermost region of the current window. The dynamic latent sampling comprises three major components: 1) a queue, 2) FoV adjustment, and 3) center-first selection.

First, we prevent selecting the same spherical latent twice by leveraging a queue. Once a latent is selected, it is removed from the queue. Second, we dynamically adjust the FoV, ensuring that H and W are not fixed within our framework. Lastly, we prioritize selecting the center-positioned spherical latent first. Since our framework already supports a dynamic FoV, reducing the FoV is not an issue, as it can be compensated by increasing the density of view directions. However, if the selected latents within the FoV are unevenly distributed, even closely spaced view directions may fail to exchange information effectively, leading to problems. The entire algorithm is demonstrated in Algorithm 1.

3.4. Distortion-Aware Weighted Averaging

Although dynamic latent sampling improves seamless continuity, minor distortion remains in the spherical-to-perspective projection. While the spherical-to-perspective distortion is relatively smaller than the ERP distortion, it can still cause latent position misalignment for other viewpoints. To address this, we proposed distortion-aware weighted averaging within the MultiDiffusion framework [2]. Specifically, we adjust the per-pixel weight $\mathbf{W}_{\mathcal{I}}^i$ to account for the spherical-to-perspective distortion.

Since distortion increases with distance from the origin of the perspective image, we introduce a simple yet effective exponential weighting function in the image space \mathcal{I} .

$$\mathbf{W}_{\mathcal{I}}^i = [W_{jk}^i]_{j \in [1, H], k \in [1, W]} \in \mathbb{R}^{H \times W}, \quad (9)$$

$$W_{jk}^i = \exp(-\|\mathbf{u}_{jk}\|/\tau), \quad (10)$$

where $\|\mathbf{u}_{jk}\|$ is the distance from the origin in the perspec-

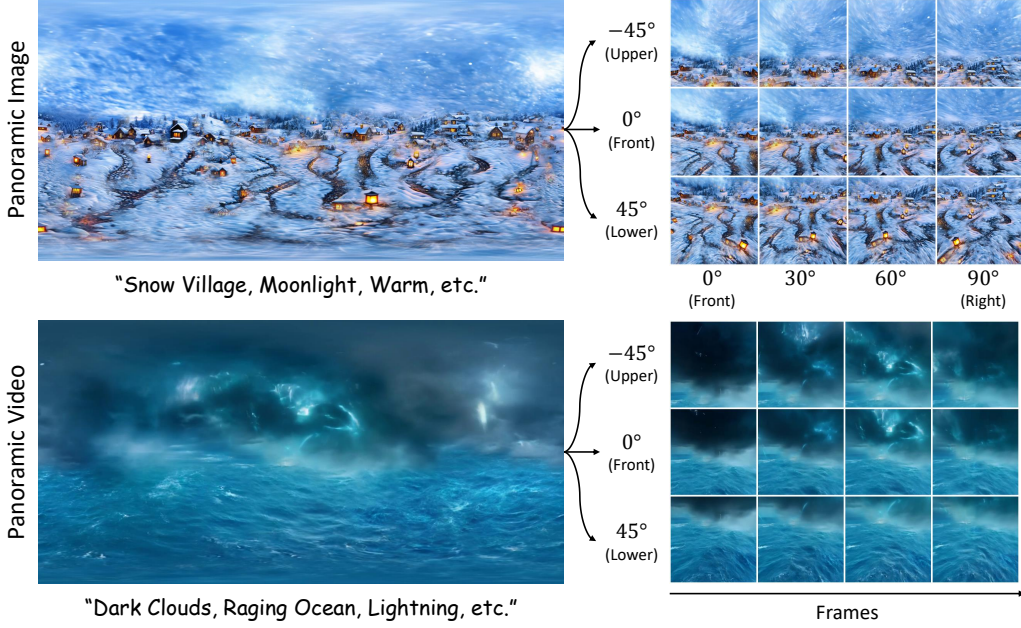


Figure 6. **Qualitative results of our method.** Visualization results for the entire scene using the ERP representation and 3 perspectives views across various elevation multiple perspective images or frames. Additional results are available in the supplementary materials.

tive image, and τ is a scaling factor controlling how quickly the weight decays toward the edges. Then, the weighting function can be represented as $\mathbf{W}_S^i = F^{-1}(\mathbf{W}_T^i)$. This per-pixel weighted average enables seamless panorama generation by assigning higher weights to the center of the FoV, where distortion is minimal.

4. Experiments

Our method is applicable to both image and video diffusion models, and we conduct experiments on both modalities. We utilize SANA [26], and LTX Video [9] for image and video generation, respectively. In all MultDiffusion setups [2], the base height and width are fixed at 512 pixels, and the temporal frame length is set to 121 frames for video experiments. Additionally, we use 2,600 points on the sphere and 89 view directions for denoising, where each perspective overlaps by 60%. To visualize the results, we decode the denoised latent representation \mathbf{S}_0 for each view direction using a VAE decoder and stitch them together to construct an ERP image. During this decoding process, we apply distortion-aware weighted averaging techniques to ensure seamless integration. Additional implementation details are provided in Section E. All source code and configurations will be made publicly available.

4.1. Experimental Setup

For evaluation, we use 20 fixed text prompts designed for generating outdoor scenes and assess all methods based

on perceptual image quality. Metrics are measured using images corresponding to 14 view directions evenly distributed across the sphere. Perspective transformations are performed with a fixed 90-degree FoV and a resolution of 512x512 pixels. We evaluate our method across four domains: panorama, image, text adherence, and video, using two metrics per category. For panorama, we assess distortion level and end continuity, while for image quality, we measure visual fidelity and aesthetic appeal. These aspects are quantitatively evaluated using LLM-based visual scores from the GPT-4o model [1], with detailed instructions in Section C.1. For text adherence and video quality, we adopt evaluation criteria from VBench [11] and CLIP-Score [20]. Additionally, we conduct a user study to assess panorama, image, and video quality. Further details on the evaluation protocol and user study are provided in Section C.

Baselines. For image panorama baselines, we include 360LoRA [15], a LoRA finetuned model from CivitAI, Text2Light [5], and Panfusion [30]. While DynamicScaler [18] targets panoramic video generation, it can be used to generate panoramic images. So we extend it to image generation and use it as a baseline for comparison. For video generation, we use 360DVD [24] and DynamicScaler as primary baselines. Furthermore, we apply 360LoRA for inference using AnimateDiff [8] and include it as an additional baseline. While DynamicScaler was originally implemented on DynamicCraft [27], we reimplemented it on the

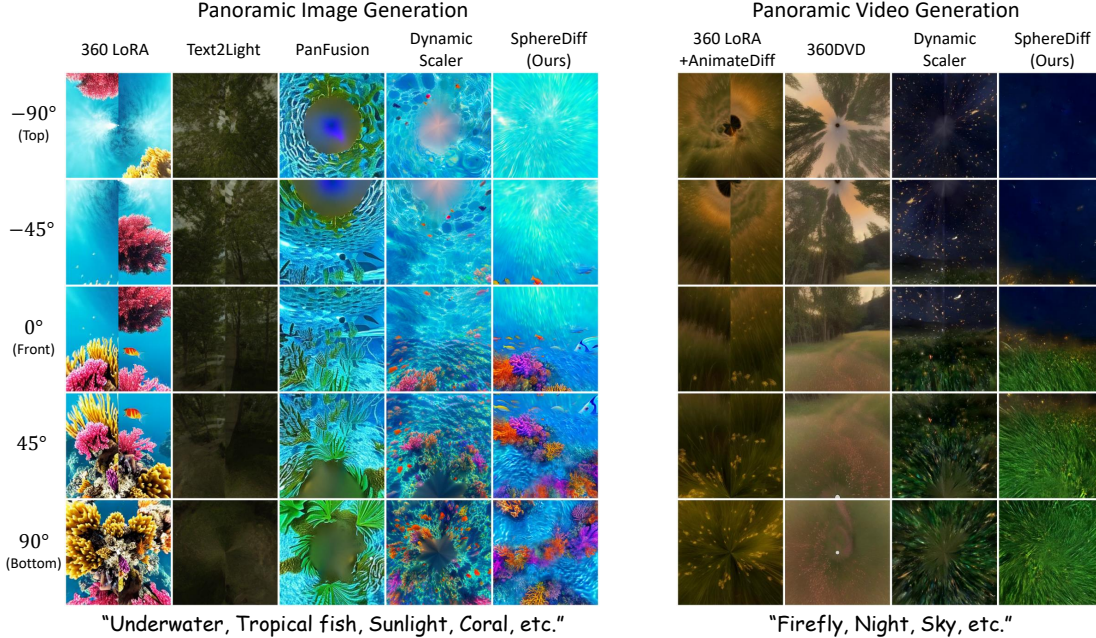


Figure 7. **Qualitative comparison of all image and video baselines.** Each sample presents perspective images from the top view to the bottom view, highlighting end-to-end continuity and distortion. Other methods exhibit noticeable artifacts, such as split seams, severe distortions near the poles, blurriness, or spots due to inadequate handling of these issues. In contrast, our method generates seamless, high-quality panoramic content without such artifacts.

Method	Panorama Criteria		Image Criteria		Text Adherence		Video Criteria	
	Distortion \uparrow	End Continuity \uparrow	Image Quality \uparrow	Aesthetic Appearance \uparrow	Scene \uparrow	CLIP-Score \uparrow	Motion Smoothness \uparrow	Temporal Flickering \uparrow
360 LoRA	2.027	3.423	2.965	3.492	<u>0.2875</u>	26.40	-	-
Text2Light [5]	2.381	3.454	2.415	2.777	0.0000	21.03	-	-
PanFusion [30]	1.965	3.696	2.819	3.450	0.2125	25.70	-	-
DynamicScaler [18]	<u>2.854</u>	<u>3.985</u>	4.496	<u>4.577</u>	0.2750	<u>26.63</u>	-	-
SphereDiff (Ours)	3.238	4.892	4.496	4.685	0.5875	28.65	-	-
360 LoRA (+ AnimateDiff [8])	1.939	3.482	3.179	3.571	0.2914	26.34	0.9908	0.9847
360DVD [24]	<u>2.086</u>	3.246	2.929	3.396	0.0719	23.13	0.9843	0.9777
DynamicScaler [18]	1.971	2.971	2.711	3.236	<u>0.4836</u>	<u>26.89</u>	<u>0.9943</u>	<u>0.9918</u>
SphereDiff (Ours)	2.579	4.496	<u>3.050</u>	3.593	0.5703	27.52	0.9956	0.9941

Table 1. **Quantitative comparison** with the best results in bold and the second best underlined. SphereDiff outperforms existing methods across all metrics except image quality for video generation, where it ranks second.

same models, SANA and LTX Video, for a fair comparison. For our method and DynamicScaler, we use identical text prompts, while other models accept only a single prompt, so we provide the main reference prompt for consistency.

4.2. Results

Qualitative Results. As shown in Tab. 3, all baseline methods operate in the equirectangular latent space. The detailed comparison is available in Section D. These methods do not fully address the inherent distortions of equirectangular projection. This limitation is evident in Fig. 7, where noticeable artifacts appear near the poles, which become even more pronounced when viewed in a perspective image. In contrast, our method effectively mitigates these distortions, producing significantly improved results. Another

major limitation of ERP-based methods is the requirement for end continuity, where the left and right edges of the image must seamlessly connect. This issue is not fully resolved in baselines, particularly when observed from specific view directions. In contrast, our method performs denoising on a uniform spherical representation, ensuring consistency across all viewing angles and eliminating such discontinuities as shown in Fig. 6.

Quantitative Results As presented in Tab. 1, our method outperforms all baselines across most of metrics for both image and video generation. Notably, our approach achieves significantly better scores in distortion and end continuity. This demonstrates its effectiveness in produc-

Method	Panorama Creteria		Image Creteria		Video Creteria	
	Distortion \uparrow	End Continuity \uparrow	Image Quality \uparrow	Text Alignment \uparrow	Motion Smoothness \uparrow	Temporal Flickering \uparrow
360 LoRA	21.43	23.81	21.43	20.24	-	-
Text2Light [5]	5.95	4.76	8.33	5.95	-	-
PanFusion [30]	14.29	10.71	10.71	16.67	-	-
DynamicScaler [18]	20.24	<u>25.00</u>	<u>25.00</u>	<u>27.38</u>	-	-
SphereDiff (Ours)	38.10	35.71	34.52	29.76	-	-
360 LoRA (+ AnimateDiff [8])	25.00	<u>27.38</u>	<u>27.38</u>	<u>27.38</u>	27.38	23.81
360DVD [24]	11.90	13.10	16.67	20.24	8.33	14.29
DynamicScaler [18]	<u>30.95</u>	26.19	28.57	22.62	<u>28.57</u>	<u>29.76</u>
SphereDiff (Ours)	32.14	33.33	<u>27.38</u>	29.76	35.71	32.14

Table 2. **User study results.** All images and videos are presented as perspective images, and the evaluation is conducted using a multiple-alternative forced-choice survey. Our method demonstrates the highest preference in most performance aspects, with particularly in distortion and end continuity.

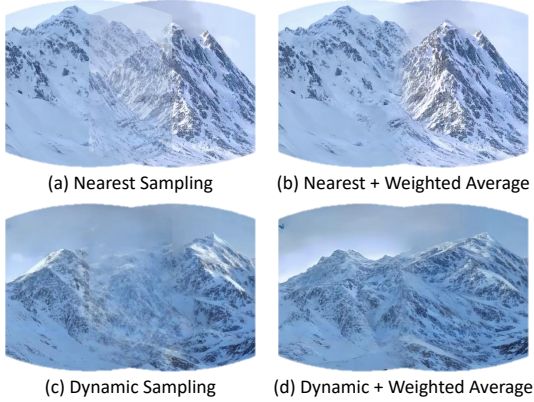


Figure 8. **Ablation on latent sampling and weighted averaging.** Nearest sampling lacks information exchange between views, leading to inconsistencies and visible overlap artifacts caused by *undersampling problem*. In contrast, dynamic sampling facilitates information sharing, resulting in more integrated outputs. With weighted averaging, both sampling methods improve seamlessness. However, nearest sampling still fails to maintain connectivity between adjacent regions, leading to discontinuities.

ing high-quality panoramic content. However, while the image quality of our video generation is lower than that of AnimateDiff, it remains comparable. The overall video generation score is lower than the image generation score, which may be attributed to the denoising model’s quality. Since our method is tuning-free, it can achieve higher performance by leveraging a more advanced denoising model, given sufficient memory availability.

4.3. User Study

To further evaluate the effectiveness of our panorama generation method, we conducted a user study using a multiple-alternative forced-choice survey. Participants were asked to compare 20 pairs of images and videos based on quality,

text alignment, distortion, end continuity, motion smoothness, and temporal flickering. The study involved 21 participants, each selecting their preferred option for each metric. As shown in Tab. 2, our method outperformed all baselines across almost all metrics in both image and video evaluations. While DynamicScaler achieve a higher image quality rating, it does so by sacrificing the 360-degree panoramic constraint. Specifically, it employs early stopping to prioritize image quality at the cost of panoramic continuity. In contrast, our approach preserves 360-degree consistency while achieving comparable image quality, balancing panoramic integrity and visual fidelity.

4.4. Ablation Study

We conducted ablation studies to evaluate the impact of each component in SphereDiff, spherical latent sampling methods, as well as the effect of distortion-aware weighted averaging. For clarity, we performed denoising on only two views and visualized the results in ERP representation. As shown in Fig. 8, nearest sampling fails to facilitate information exchange between different views, leading to unnatural transitions in overlapping regions. Our proposed dynamic latent sampling method improves information exchange between views, generating more seamless images. Furthermore, by incorporating our distortion-aware weighted averaging technique, we achieve significantly cleaner and more coherent outputs for both sampling methods. These results demonstrate that each component in our framework plays a critical role in seamlessly integrating multiple perspectives into a single spherical representation, ensuring high-quality panoramic generation.

5. Conclusion

We introduced SphereDiff, a tuning-free framework for omnidirectional panoramic generation that leverages spherical latent representations to minimize distortions and ensure

seamless continuity. Our method integrates dynamic latent sampling and distortion-aware weighted averaging, significantly improving panoramic content quality. Despite these advancements, limitations remain. Each viewpoint is currently processed independently, leading to a lack of global context. Future work will focus on addressing this challenge by incorporating global context-aware refinement methods. Nevertheless, we believe our framework is a solid starting point for 360-degree panorama generation.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 6
- [2] Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: fusing diffusion paths for controlled image generation. In *Proceedings of the 40th International Conference on Machine Learning*, pages 1737–1752, 2023. 2, 3, 4, 5, 6
- [3] Pascal Chang, Jingwei Tang, Markus Gross, and Vinicius C Azevedo. How i warped your noise: a temporally-correlated noise prior for diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024. 5
- [4] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 2
- [5] Zhaoxi Chen, Guangcong Wang, and Ziwei Liu. Text2light: Zero-shot text-driven hdr panorama generation. *ACM Transactions on Graphics (TOG)*, 41(6):1–16, 2022. 2, 3, 6, 7, 8, 12
- [6] Jaeyoung Chung, Suyoung Lee, Hyeongjin Nam, Jaerin Lee, and Kyoung Mu Lee. Luciddreamer: Domain-free generation of 3d gaussian splatting scenes. *arXiv preprint arXiv:2311.13384*, 2023. 3
- [7] Mohammad Reza Karimi Dastjerdi, Yannick Hold-Geoffroy, Jonathan Eisenmann, Siavash Khodadadeh, and Jean-François Lalonde. Guided co-modulated gan for 360 $\{\backslash\text{deg}\}$ field of view extrapolation. *arXiv preprint arXiv:2204.07286*, 2022. 3
- [8] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 6, 7, 8
- [9] Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson, Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, et al. Ltx-video: Realtime video latent diffusion. *arXiv preprint arXiv:2501.00103*, 2024. 2, 4, 6, 13
- [10] Doug P Hardin, TJ Michaels, and Edward B Saff. A comparison of popular point configurations on \mathbb{S}^2 . *arXiv preprint arXiv:1607.04590*, 2016. 4
- [11] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 6
- [12] Nikolai Kalischek, Michael Oechsle, Fabian Manhardt, Philipp Henzler, Konrad Schindler, and Federico Tombari. Cubediff: Repurposing diffusion-based image models for panorama generation. In *The Thirteenth International Conference on Learning Representations*, 2025. 3, 12
- [13] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 2
- [14] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 2
- [15] LatentLabs360. Latentlabs360. <https://civitai.com/models/10753/latentlabs360>, 2023. 2, 6
- [16] Renjie Li, Panwang Pan, Bangbang Yang, Dejia Xu, Shijie Zhou, Xuanyang Zhang, Zeming Li, Achuta Kadambi, Zhangyang Wang, Zhengzhong Tu, et al. 4k4dgen: Panoramic 4d generation at 4k resolution. *arXiv preprint arXiv:2406.13527*, 2024. 2, 3, 12
- [17] Dongyang Liu, Shicheng Li, Yutong Liu, Zhen Li, Kai Wang, Xinyue Li, Qi Qin, Yufei Liu, Yi Xin, Zhongyu Li, et al. Lumina-video: Efficient and flexible video generation with multi-scale next-dit. *arXiv preprint arXiv:2502.06782*, 2025. 2
- [18] Jinxiu Liu, Shaoheng Lin, Yinxiao Li, and Ming-Hsuan Yang. Dynamicscaler: Seamless and scalable video generation for panoramic scenes. *arXiv preprint arXiv:2412.11100*, 2024. 2, 3, 6, 7, 8, 12
- [19] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022. 4, 11
- [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021. 6
- [21] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 4
- [22] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. 4
- [23] Genmo Team. Mochi 1. <https://github.com/genmoai/models>, 2024. 2
- [24] Qian Wang, Weiqi Li, Chong Mou, Xinhua Cheng, and Jian Zhang. 360dvd: Controllable panorama video generation with 360-degree video diffusion model. In *Proceedings of*

- the *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6913–6923, 2024. [2](#), [3](#), [6](#), [7](#), [8](#), [12](#)
- [25] Tianhao Wu, Chuanxia Zheng, and Tat-Jen Cham. Panodiffusion: 360-degree panorama outpainting via diffusion. *arXiv preprint arXiv:2307.03177*, 2023. [3](#)
- [26] Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li, Ligeng Zhu, Yao Lu, et al. Sana: Efficient high-resolution image synthesis with linear diffusion transformers. *arXiv preprint arXiv:2410.10629*, 2024. [2](#), [6](#), [13](#)
- [27] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Wangbo Yu, Hanyuan Liu, Gongye Liu, Xintao Wang, Ying Shan, and Tien-Tsin Wong. Dynamicrafter: Animating open-domain images with video diffusion priors. In *European Conference on Computer Vision*, pages 399–417. Springer, 2024. [6](#)
- [28] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. [2](#)
- [29] Hong-Xing Yu, Haoyi Duan, Charles Herrmann, William T Freeman, and Jiajun Wu. Wonderworld: Interactive 3d scene generation from a single image. *arXiv preprint arXiv:2406.09394*, 2024. [3](#)
- [30] Cheng Zhang, Qianyi Wu, Camilo Cruz Gambardella, Xiaoshui Huang, Dinh Phung, Wanli Ouyang, and Jianfei Cai. Taming stable diffusion for text to 360 panorama image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6347–6357, 2024. [2](#), [3](#), [6](#), [7](#), [8](#), [12](#)
- [31] Shilong Zhang, Wenbo Li, Shoufa Chen, Chongjian Ge, Peize Sun, Yida Zhang, Yi Jiang, Zehuan Yuan, Binyue Peng, and Ping Luo. Flashvideo: Flowing fidelity to detail for efficient high-resolution video generation. *arXiv preprint arXiv:2502.05179*, 2025. [2](#)
- [32] Shijie Zhou, Zhiwen Fan, Dejia Xu, Haoran Chang, Pradyumna Chari, Tejas Bharadwaj, Suyu You, Zhangyang Wang, and Achuta Kadambi. Dreamscene360: Unconstrained text-to-3d scene generation with panoramic gaussian splatting. In *European Conference on Computer Vision*, pages 324–342. Springer, 2024. [3](#)

SphereDiff: Tuning-free Omnidirectional Panoramic Image and Video Generation via Spherical Latent Representation

Supplementary Material

A. Appendix

A Appendix	11
B Spherical Latent Sampling	11
B.1. Perspective Latent Representation.	11
B.2. Directly Use Continuous Position	11
C Experimental Details	11
C.1. LLM-based visual score	12
C.2. User Study Details	12
D Additional Qualitative/Quantitative Comparison	12
D.1. Detailed Comparison with Spherical and ERP	12
E Implementation Details	13
E.1. Text Prompt Examples	13

B. Spherical Latent Sampling

B.1. Perspective Latent Representation.

For perspective latent representation, we define a virtual camera centered at the origin. The points in world coordinates, denoted as $d = (x, y, z)^\top$, are projected onto image space using the projection matrix $P = K[R|t]$. Here, K is the intrinsic camera matrix derived from a predefined focal length f , R represents the viewing direction, and t is set to zero. The spherical-to-perspective projection function $\mathcal{T}_{\mathbb{S}^2 \rightarrow \mathbb{P}^2}$ can be formulated as:

$$\tilde{u} = K[R|t]\tilde{d}, \quad (11)$$

where $\tilde{d} = (x, y, z, 1)^\top$ is the homogeneous coordinate representation of the 3D point d , and $\tilde{u} = (u', v', w')^\top$ represents the projected homogeneous coordinates in image space. The final 2D perspective coordinates $u = (u, v)^\top$ are obtained via perspective division:

$$u = \left(\frac{u'}{w'}, \frac{v'}{w'} \right). \quad (12)$$

To ensure proper visibility, points located behind the view direction are masked out using their inner product values, retaining only the points in the frontal view.

B.2. Directly Use Continuous Position

Recent visual generation models, including those based on DiT [19], support continuous 1D representations through corresponding positional embeddings. A naive approach

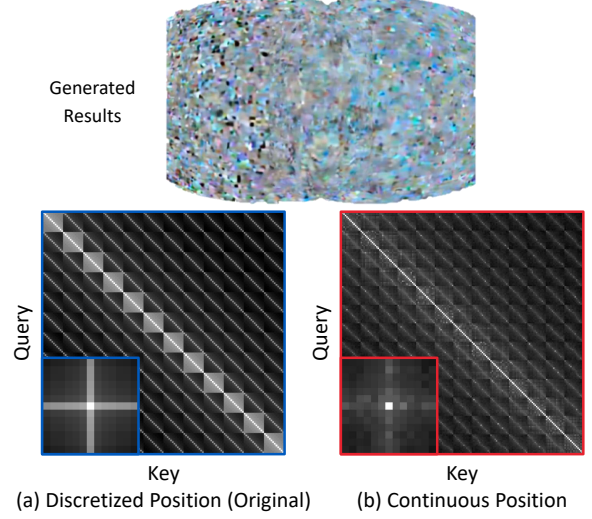


Figure 9. **Similarities between positional embeddings calculated with discrete and continuous positions.** The small squares represent the similarity distribution when the central pixel is used as a query. As shown in the figure, when discretization is applied, within the same row or column, resulting in high similarity. However, when positions vary slightly (as in the continuous case), the similarity drops significantly.

would be to leverage this property and treat the latent representations as continuous without discretization. However, this approach leads to unstructured outputs due to a distribution shift in positional embeddings. Specifically, in the original positional embedding space, latent similarities are high within the same row or column, ensuring spatial consistency. In contrast, when using continuous positional embeddings, the similarity between two adjacent points is not necessarily high, even if their spatial coordinates are close, as shown in Fig. 9. This discrepancy causes the model to fail in generating structured content. Although DiT can process continuous inputs, discretization remains essential for tuning-free panoramic visual generation to maintain structured and consistent latent relationships.

C. Experimental Details

All metrics were measured using perspective images. The selected view directions include four azimuth angles (0° , 90° , 180° , 270°) at elevation angles of 0° , 45° , and -45° , as well as one view each at elevations of 90° and -90° , resulting in a total of 14 view directions.

Method	Latent Space	Tuning-Free	Open-Sourced
Panoramic image generation			
Text2Light [5]	ERP	✗	✓
PanFusion [30]	ERP	✗	✓
Cubediff [12]	Cube Map	✗	✗
Panoramic video generation			
360DVD [24]	ERP	✗	✓
4K4DGen [16]	ERP	✗	✗
Panoramic image and video generation			
DynamicScaler [18]	ERP	✓	✗
SphereDiff (Ours)	Spherical	✓	✓

Table 3. Comparison of 360-degree panorama generation approach. Most existing panoramic generation models perform denoising in the equirectangular latent space, except for CubeDiff, which utilizes a cube map representation. In contrast, our method leverages spherical latents. Among the compared methods, only DynamicScaler and ours support tuning-free generation.

C.1. LLM-based visual score

GPT-4o evaluation prompt used for assessing image quality, aesthetic appeal, distortion level, and connectivity. Tab. 4 shows the prompt that we used for evaluate.

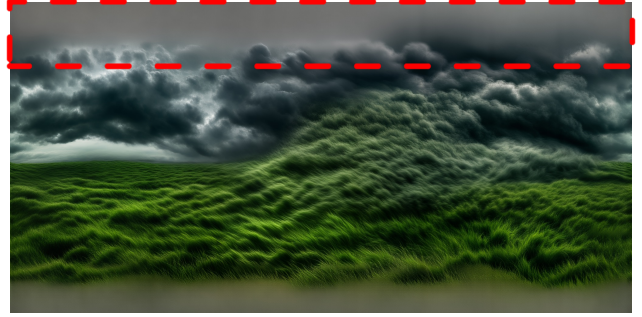
C.2. User Study Details

We divided the 20 pairs of images and videos into five sets, each containing four pairs, allowing users to select one set for evaluation based on their convenience. To accurately assess distortion and end-continuity, images and videos were presented from a viewpoint with an azimuth angle ($\theta = 0^\circ$) while varying the elevation angle (ϕ) across five positions: $-90^\circ, -45^\circ, 0^\circ, 45^\circ, 90^\circ$. The presentation order was randomized to minimize bias. For image evaluation, participants selected the most suitable model from five available options, while for video evaluation, they chose from four models, selecting the one they found most appropriate for each criterion. The evaluation criteria were instructed to be used similarly to those described in Tab. 4, which were designed for evaluation using GPT-4o.

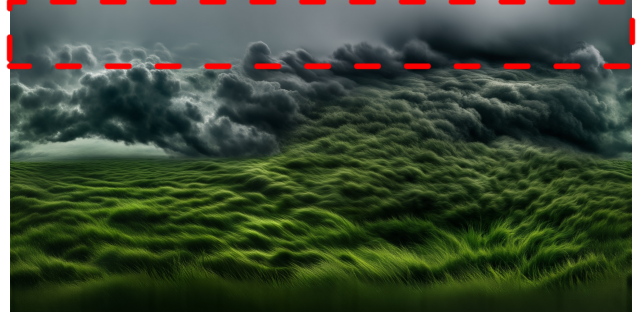
D. Additional Qualitative/Quantitative Comparison

D.1. Detailed Comparison with Spherical and ERP

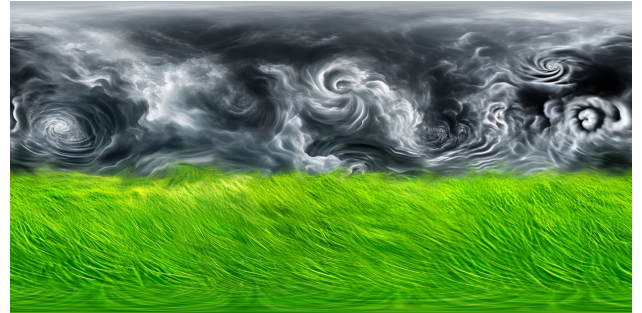
ERP projects latitude and longitude coordinates on the sphere onto vertical and horizontal coordinates on a rectangular grid, forming an $H \times W \times C$ representation, where H and W denote the resolution, and C represents the number of channels. While this format maintains a uniform distribution in the 2D rectangular grid, it results in a non-uniform distribution on the sphere. Specifically, points become densely concentrated near the poles, as shown in



(a) Dynamic Scaler w/ low view direction number



(b) Dynamic Scaler w/ high view direction number



(c) Ours

Figure 10. **Comparison with spherical and ERP.** The red box highlights blurry artifacts that appear near the polar regions. As the number of view directions decreases, more latents remain un-sampled and unprocessed during denoising, making the issue more severe.

Fig. 4. If this characteristic is ignored and the 2D grid is directly used for tuning, it leads to severe artifacts when viewed in perspective. As illustrated in Fig. 2, a significant number of points are concentrated in certain areas, creating an artifact where points appear to be pulled toward the center, disrupting spatial consistency.

To address these issues, recent works [16, 18] attempt to leverage the spherical properties of ERP for generating 360-degree panoramic videos. However, these methods also encounter challenges during the projection between perspective and ERP spaces. Interpolation distorts the latent distribution, while sampling-based methods result in missing

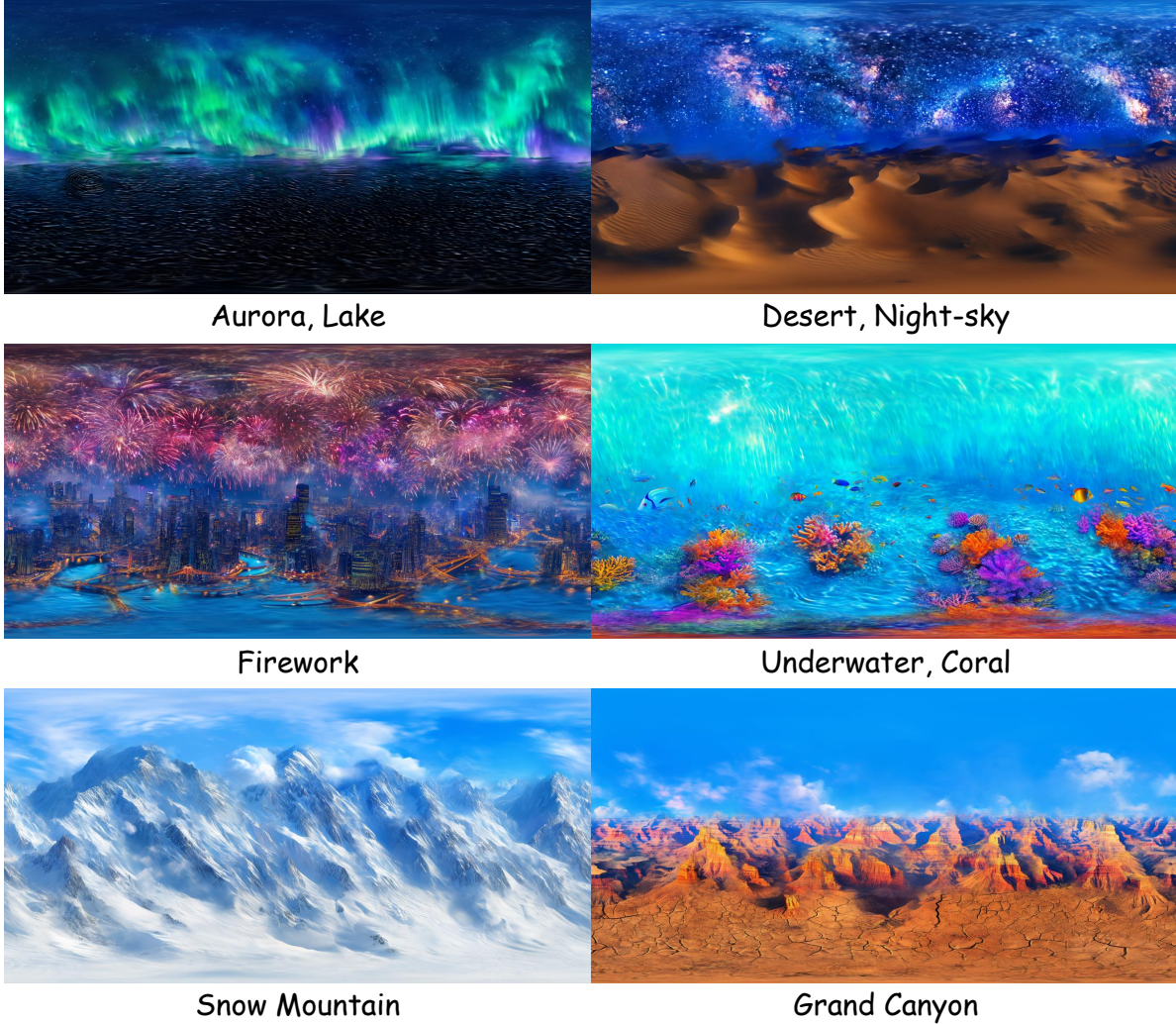


Figure 11. Additional Qualitative Results (Images)

points, which disrupt multi-view stitching and prevent effective information exchange between different views, leading to independent generation of each segment. In particular, applying the Offset Shifting Denoiser (OSD) from DynamicScaler can cause certain points to remain unprocessed during denoising steps, leading to blurry artifacts, as shown in Fig. 10. In contrast, our spherical representation ensures a uniform latent distribution across the sphere, effectively eliminating these issues. As a result, our method enables seamless tuning-free image and video generation.

E. Implementation Details

We used 89 view directions with an 80° FoV and 2,600 points for all experiments. Inference time per sample is approximately 30 seconds for image generation and 20 minutes for video generation. All experiments were conducted

on an A100-40GB. Since our method does not require additional memory beyond what is needed by the models used in our experiments, SANA and LTX Video, we verified that it can also run on an RTX 3090 with 24GB VRAM. For text prompts, we use three descriptions corresponding to the top, middle, and bottom regions of the scene.

E.1. Text Prompt Examples

We conducted evaluations across 20 different scene concepts, ranging from urban to natural landscapes. Since the prompt generation rules differ between SANA [26] and LTX Video [9], slight variations in results may occur. Specifically, LTX Video struggles with very short prompts, so minor modifications were made to improve generation quality. Nevertheless, all prompts will be publicly released. Here are some example text prompts.



Figure 12. Additional Qualitative Results (Videos). The animated results are available on our project page.

Prompts for generating panoramic images in the main manuscript.

• **Underwater**

- An upward view from underwater, looking towards the surface where sunlight beams penetrate the clear ocean. Gentle ripples create a shimmering effect, and the water transitions from deep blue to a lighter, almost turquoise hue near the surface. The light refracts beautifully, creating a dreamlike underwater glow.
- A stunning underwater scene filled with vibrant tropical fish swimming gracefully through the crystal-clear water. Various species, from small neon-colored fish to larger, elegant ones, move in harmony. The environment is serene, with the water gently flowing and reflecting the sunlight from above. The clarity of the water highlights the intricate details of the marine life.
- A breathtaking top-down view of a colorful tropical coral reef. The ocean floor is covered with vibrant corals, ranging from bright orange and pink to deep purple and blue. Small fish dart between the coral formations, while gentle waves create a mesmerizing play of light and shadow on the seabed. The details of the marine ecosystem are incredibly vivid, showing the beauty of underwater biodiversity.

• **Snow Village**

- An upward view of the night sky. Soft moonlight filters through wispy clouds, casting a serene glow over the winter landscape.

- From the snowy fields, a view toward a peaceful village nestled among snow-covered hills. Warm lights glow from the windows of small wooden cabins, contrasting with the crisp, cold air under the moonlit sky.
- A high-angle view of snow-covered paths winding through the landscape. The fresh snow glistens under the moonlight, while the warm glow of lanterns and fireplaces reflects off the frosty roads, creating a cozy contrast against the cold night.

• **Green grass**

- Dark clouds churned in slow, twisting spirals overhead, their shifting forms casting fleeting shadows below. The clouds thickened, their edges curling like ink dissolving in water, deepening into shades of charcoal and silver. Occasionally, bright patches pierced through the dense formations, creating stark contrasts of light and darkness in the sky.
- Dark clouds churned in slow, twisting spirals above the rolling expanse of vibrant green grass, their shifting forms casting fleeting shadows across the land. A soft breeze rustled the blades, carrying the crisp scent of damp earth. The clouds thickened, their edges curling like ink dissolving in water, deepening in shades of charcoal and silver. Patches of sunlight briefly pierced through, creating shifting patterns of light and shadow that danced across the swaying field.
- A rolling expanse of vibrant green grass stretched endlessly, each blade rustling softly in the gentle breeze.

The crisp scent of damp earth lingered in the air as the grass swayed rhythmically, creating subtle waves across the field. The surface shimmered with varying shades of green, highlighting the texture and movement of the landscape.

Prompts for generating panoramic videos in the main manuscript.

• Fireworks

- Vibrant fireworks burst across the night sky, painting the heavens with shimmering trails of vivid colors. Some fireworks transform into heart shapes before fading, adding a touch of elegance to the display. The camera focuses on explosive arcs and sparkling embers, capturing every brilliant flash against an infinite, celestial canvas.
- The city's skyline stretches below, clearly visible as vibrant fireworks light up the night sky. The fireworks burst in various colors, scattering across the air, while their reflections shimmer on the glass windows of skyscrapers. The camera smoothly pans across the city, capturing the river and bridges, with distant car lights creating a flowing effect.
- A breathtaking city skyline stretches below, illuminated by countless lights reflecting off towering skyscrapers. The camera smoothly pans across the landscape, revealing a river winding through the metropolis and bridges glowing under streetlights. Distant car headlights flow like streams of light, adding a dynamic rhythm to the urban nightscape.

• Storm

- Dark storm clouds swirl overhead as multiple lightning bolts strike at different moments, briefly illuminating the chaotic sky. The jagged bolts cut through the darkness, revealing shifting cloud formations in flashes of electric blue and white. The perspective is an upward view, emphasizing the storm's immense scale. The scene is dynamic, with each lightning strike casting sharp contrasts of light and shadow across the turbulent sky.
- A mid-angle view reveals the vast ocean meeting the storm-filled sky, where towering waves rise and fall beneath the relentless tempest. Lightning bolts crack through the heavy clouds, their electric glow reflecting off the turbulent water. The horizon is barely visible, obscured by mist and rain as gusting winds whip across the sea's surface. Each flash of light briefly exposes the chaos, illuminating the swirling storm above and the restless waves below.
- A high-angle view captures the vast, turbulent deep blue ocean as powerful waves crash and churn beneath the storm's force. The water's surface ripples with energy, each undulating motion reflecting the storm's in-

tensity. White foam swirls atop the restless sea, contrasting against the dark depths. Gusts of wind carve patterns into the waves, while distant flashes of lightning momentarily reveal the chaotic movement below.

• Firefly

- An upward view of the vast night sky above an open field, with scattered fireflies drifting gently, their faint glows blending with distant stars.
- A serene nighttime meadow, where countless fireflies flicker softly, casting a warm, golden light that dances over the swaying grass and wildflowers.
- A dramatic top-down view of a vast open field grass, where waves of grass ripple in the night breeze, dotted with countless fireflies drifting just above, their soft glow flickering across the landscape.

Table 4. Evaluation prompt used for assessing 360-degree panoramic generation quality based on four criteria.

Evaluation Prompt
<p>You are an evaluator assessing an image generation model based on a single image at a time. Your evaluation is based on the following four criteria:</p> <ol style="list-style-type: none"> 1. Image Quality: Assess the overall quality of the image. 2. Aesthetic Appeal: Evaluate how visually pleasing the image is. 3. Distortion Level: Determine whether the image appears distorted. If it does not resemble a photo taken with a normal camera, it will receive a lower score. 4. Connectivity: Check if the middle of the image appears disconnected. If there is a noticeable break, the score will be lower. <p>Each criterion is rated on a five-point scale: Excellent (5), Good (4), Fair (3), Poor (2), and Awful (1).</p> <p>You will receive one image at a time. For each criterion, provide a concise reason for the score before listing the rating.</p> <p>Format your response as follows:</p> <ul style="list-style-type: none"> - Image Quality: (Brief reason) → Score - Aesthetic Appeal: (Brief reason) → Score - Distortion Level: (Brief reason) → Score - Connectivity: (Brief reason) → Score <p>{image}</p> <p>Please evaluate the image with the given criteria.</p>