# PanoWan: Lifting Diffusion Video Generation Models to 360° with Latitude/Longitude-aware Mechanisms

**Yifei Xia**[1,2,3]*  **Shuchen Weng**[4]*  **Siqi Yang**[5]  **Jingqi Liu**[1,2]  **Chengxuan Zhu**[6]

**Minggui Teng**[1,2]  **Zijian Jia**[7]  **Han Jiang**[3]  **Boxin Shi**[1,2]†

[1]State Key Lab of Multimedia Info. Processing, School of Computer Science, Peking University
[2]Nat'l Eng. Research Ctr. of Visual Tech., School of Computer Science, Peking University
[3]OpenBayes Information Technology Co., Ltd.  [4]Beijing Academy of Artificial Intelligence
[5]Institute for Artificial Intelligence, Peking University
[6]Nat'l Key Lab of General AI, School of Intelligence Science and Technology, Peking University
[7]School of Artificial Intelligence, Beijing University of Posts and Telecommunications
`{yfxia,shuchenweng,yousiki,peterzhu,minggui_teng,shiboxin}@pku.edu.cn`
`liujingqi@stu.pku.edu.cn    jiazijian@bupt.edu.cn    hahn@openbayes.com`

## Abstract

Panoramic video generation enables immersive 360° content creation, valuable in applications that demand scene-consistent world exploration. However, existing panoramic video generation models struggle to leverage pre-trained generative priors from conventional text-to-video models for high-quality and diverse panoramic videos generation, due to limited dataset scale and the gap in spatial feature representations. In this paper, we introduce PanoWan to effectively lift pre-trained text-to-video models to the panoramic domain, equipped with minimal modules. PanoWan employs latitude-aware sampling to avoid latitudinal distortion, while its rotated semantic denoising and padded pixel-wise decoding ensure seamless transitions at longitude boundaries. To provide sufficient panoramic videos for learning these lifted representations, we contribute PANOVID, a high-quality panoramic video dataset with captions and diverse scenarios. Consequently, PanoWan achieves state-of-the-art performance in panoramic video generation and demonstrates robustness for zero-shot downstream tasks. Our project page is available at `https://panowan.variantconst.com`.

## 1 Introduction

Text-based panoramic video generation aims to produce a complete 360° view, ensuring coherent spatial and visual relationships between elements within the scene. Such inherent property is highly valuable for conventional VR content, the construction of interactive game worlds [35, 8], and the simulation of environments for embodied AI [19].

The remarkable capabilities of conventional text-to-video models [29, 33, 4] motivate researchers to leverage their generative priors to panoramic video generation. One intuitive strategy generates local perspective conventional videos and integrates them during inference. While these training-free methods [12, 22] entirely preserve generative priors, they sacrifice overall consistency as they struggle to establish cross-view long-range dependencies. Alternatively, fine-tuning conventional text-to-video models [32, 40] also faces challenges. On one hand, existing panoramic video datasets are limited in scope and scale compared to conventional ones. On the other hand, the gap in spatial

---

*Equal contribution.

†Corresponding author.

**Text-to-video generation**



*A black hyper-car speeds through cyberpunk highway.*



*Cowboys ride through sunset-lit western town, visitors explore old streets.*



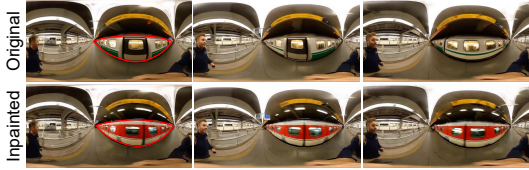*Medieval event with knights jousting, crowds cheering, camps bustling lively.*

**Long video**



*Sunset at a beach.*

**Super-resolution**



*Lab scene with researchers and equipment.*

**Semantic editing**



*Change the color of the train to red.*

**Video outpainting**



*Colorful hot air balloons.*

Figure 1: PanoWan is a text-based panoramic video generation framework. It lifts pre-trained generative priors from a conventional text-to-video model to the panorama, and enables generating diverse scenarios for long videos. Equipped with training-free techniques, PanoWan supports zero-shot editing of panoramic videos, including super-resolution, semantic editing, and video outpainting.

representations (*e.g.*, latitudinal distortions and seam longitudes) between panoramic and conventional videos potentially hinder the effective prior leverage from pre-trained conventional models.

In this paper, we pursue a training-based approach to overcome current bottlenecks. We propose **PanoWan**, a framework for **Pano**ramic video generation based on **Wan** 2.1 [33]. Equipped with minimal modules, PanoWan effectively lifts generative priors from a pre-trained conventional text-to-video model to the panorama. To bridge the gap in spatial feature representations between panoramic and conventional videos, we design latitude-aware sampling to address latitudinal distortions caused by equirectangular projection. Since text-to-video models lack continuity awareness for left and right boundaries, we achieve seamless longitude transitions using the rotated semantic denoising to address semantic inconsistency and the padded pixel-wise decoding to resolve pixel-wise disharmony.

As learning to lift spatial representations from conventional videos to the panorama requires large-scale data, we further introduce PANOVID. This **Pano**ramic **Vid**eo dataset offers diverse scenarios (*e.g.*, landscape, streetscape, and humanscape), includes over 13K captioned video clips totaling 944 hours, and features data processing tailored for panoramic video generation. As shown in Fig. 1, our framework produces panoramic videos from text descriptions, enabling long video generation. Additionally, it enables training-free editing of user-provided panoramic videos, including super-resolution, semantic inpainting, and video outpainting. Extensive experiments demonstrate that

PanoWan achieves state-of-the-art panoramic video generation performance across seven metrics, alongside robust zero-shot capabilities for various downstream tasks.

Our contributions can be summarized as follows:

- We contribute PANOVID, a large-scale and high-quality panoramic video dataset with captioned video clips, tailored for text-based panoramic video generation.
- We propose PanoWan, lifting generative priors from a pre-trained model to show state-of-the-art performance on panoramic video generation and robustness for downstream tasks.
- We integrate the latitude-aware sampling, rotated semantic denoising, and padded pixel-wise decoding to bridge the spatial difference between panoramic and conventional videos.

## 2 Related Works

### 2.1 Video Diffusion Models

Diffusion models have demonstrated impressive results in image generation, but extending them to videos introduces additional challenges (*e.g.*, temporal consistency and computational efficiency). Early models adapt image diffusion models with cascaded architecture [10], temporal layers [26], and attention mechanisms [3]. To further improve efficiency, LVDM [9] introduces a lightweight latent video diffusion model with a 3D latent space and hierarchical structure for long video generation. Meanwhile, SVD [4] improves video diffusion using large and high-quality datasets. Recent Diffusion Transformers (DiTs) [23] effectively model complex spatio-temporal dynamics for video generation, operating on latent space compressed by 3D VAE [39]. After that, more large-scale models (*e.g.*, HunyuanVideo [14], CogVideoX [38], and Seaweed [25]) emerge and demonstrate the benefits of scaling both model and data size. This motivates us to adopt Wan 2.1 [33] as the backbone to leverage its strong generative priors and temporal modeling capabilities for panoramic video generation.

### 2.2 Panoramic Video Generation

**Text conditioned generation.** 360DVD [32] is the pioneer to introduce stable video generation techniques to panoramic video generation by proposing a 360-Adapter and a set of 360 enhancement techniques. Training-free methods (*e.g.*, DynamicScaler [12] and SphereDiff [22]) create panoramic videos by generating local patches and then composing them together into a complete panorama, which inherently break global consistency. With the development of video diffusion models, VideoPanda [36] augments diffusion models with multi-view attention. PanoDiT [40] uses a DiT backbone with global-temporal attention and panoramic-specific losses for coherent long-range generation. Despite these advancements, existing methods still suffer from observable latitude distortions and issues with seam-free longitude transitions. In this work, we introduce PanoWan, a framework that addresses these challenges by lifting generative priors from pre-trained text-to-video models to panorama.

**Image or video conditioned generation.** Imagine360 [27] adopts antipodal-aware motion modeling to convert perspective videos to panoramic views. Building on static panoramas, 4K4DGen [16] lifts them into dynamic 4D scenes via spatial-temporal denoising. HoloTime [41] further leverages Gaussian splatting and a two-stage diffusion process for high-fidelity 4D reconstruction. VidPanos [21] treats panorama generation as a space-time outpainting task from panning video inputs, while Argus [20] integrates motion and geometry cues for enhanced video-to-360° synthesis. These explorations highlight a trend toward unifying spatial, temporal, and geometric reasoning. We demonstrate that PanoWan possesses these capabilities, with robust zero-shot capabilities for downstream tasks.

## 3 PANOVID Dataset

The absence of paired datasets has long been regarded as one of the primary barriers to advancing the performance of panoramic video generation models [32]. Existing text to panoramic video generation methods [32, 36, 40] mainly rely on WEB360 dataset [32], which contains only 2114 video clips of 10 seconds each. Although Argus [20] filters out over 283K video clips from the 360-1M dataset [30], it is not built for the text-based panoramic video generation task, providing no paired captions, and showing significant distribution bias for the scenario semantics.

To address these limitations, we present PANOVID, a large-scale and high-quality dataset with diverse scenarios and balanced semantics, tailored for text-based panoramic video generation. Our data collection process begins by aggregating videos from existing panoramic sources, including 360-1M [30], 360+x [5], Imagine360 [27], WEB360 [32], Panonut360 [37], the Miraikan 360-degree Video Dataset [1], and a public dataset of immersive VR videos [15]. Subsequently, we use Qwen-2.5-VL [2] to generate text descriptions and predict POI (Point-of-Interest) categories for each video. To ensure semantic uniqueness, we perform redundancy removal based on caption similarity. Under each POI category, the videos are further filtered according to optical flow [7] smoothness and aesthetic scores [34], retaining only the top-ranked samples. Thanks to the filtering pipeline, PANOVID features more than 13K high-quality video clips totaling approximately 944 hours, and is semantically diverse and balance.

## 4 Method

Panoramic videos have a different spatial feature representation compared to conventional ones. Inspired by GEN3C [24], we effectively preserve the generative prior of pre-trained models by equipping minimal modules and fine-tuning a small subset of parameters via LoRA [11]. We firstly introduce our video diffusion backbone and formulate the spherical coordinate mapping (Sec. 4.1). Next, we propose the latitude-aware sampling to avoid latitude distortion, along with its corresponding analysis (Sec. 4.2). Finally, we present the rotated semantic denoising and the padded pixel-wise decoding to achieve the seamless longitude transitions (Sec. 4.3).

### 4.1 Preliminaries

**Video diffusion models.** We employ Wan 2.1 [33] as the video generation backbone, with spatial-temporal Variational AutoEncoders (VAEs) to map high-dimensional videos into compact latent codes. The flow matching framework [18] is used to model a unified denoising diffusion process. Specifically, given a clear video $x$, a VAE encoder $\mathrm{E}(\cdot)$ first projects the video into the latent space $z_1 = \mathrm{E}(x)$. During training, a noise $z_0 \sim \mathcal{N}(0, I)$ is sampled, and an intermediate latent code $z_t = tz_1 + (1-t)z_0$ is constructed by linearly interpolating between $z_1$ and $z_0$ at timestep $t \in [0, 1]$. The training goal is to predict the ground truth velocity $v_t = \mathrm{d}z_t/\mathrm{d}t = z_1 - z_0$, and the loss function is formulated as:

$$\mathcal{L} = \mathbb{E}_{z_0, z_1, c_{\text{txt}}, t} ||u(z_t, c_{\text{txt}}, t; \theta) - v_t||^2, \tag{1}$$

where $c_{\text{txt}}$ is the text embedding, $\theta$ is the parameters of the prediction model, and $u(z_t, c_{\text{txt}}, t; \theta)$ is the predicted velocity of the model.

**Spherical coordinate mapping.** Panorama captures a $360°$ view, inherently representing signals in spherical coordinates $(\varphi, \theta)$. To leverage generative priors from conventional images and videos that operate in Cartesian coordinates $(x, y)$, we employ the equirectangular projection (ERP) $\mathcal{P}_{\text{ERP}}$ to map between these coordinate systems for panoramic videos:

$$\mathcal{P}_{\text{ERP}} : [2R] \times [R] \to [0, 2\pi] \times [-\frac{\pi}{2}, \frac{\pi}{2}], \quad (x, y) \mapsto (\varphi, \theta) = \left(\frac{2x+1}{2R}\pi, \frac{2y+1-R}{2R}\pi\right), \tag{2}$$

where $R$ is the radius of the sphere. $\varphi$ and $\theta$ are longitude and latitude respectively. While ERP enables the direct application of pre-trained VAEs to encode panoramic videos into latent codes for diffusion processes, it introduces extreme horizontal stretching in polar regions. This horizontal stretching phenomenon arises from the altered representation of distances during projection, and is recognized by changes in horizontal signal frequency. Let $\mathrm{d}s_\varphi$ and $\mathrm{d}s_\theta$ represent the infinitesimal arc lengths along lines of constant latitude and longitude, respectively. They are formulated as:

$$\mathrm{d}s_\varphi = 2R \arcsin(\cos\theta \cdot \sin\frac{\mathrm{d}\varphi}{2}) = R\cos\theta \, \mathrm{d}\varphi, \quad \mathrm{d}s_\theta = R \, \mathrm{d}\theta, \tag{3}$$

where $\theta$ is the latitude. We further consider the spherical frequency $f_{\text{sph}}$ (cycles per unit physical distance) and the Cartesian frequency $f_{\text{car}}$ (cycles per pixel in the image). Assuming that warping preserves content, their relationship in frequency is scaled by the change in distance:
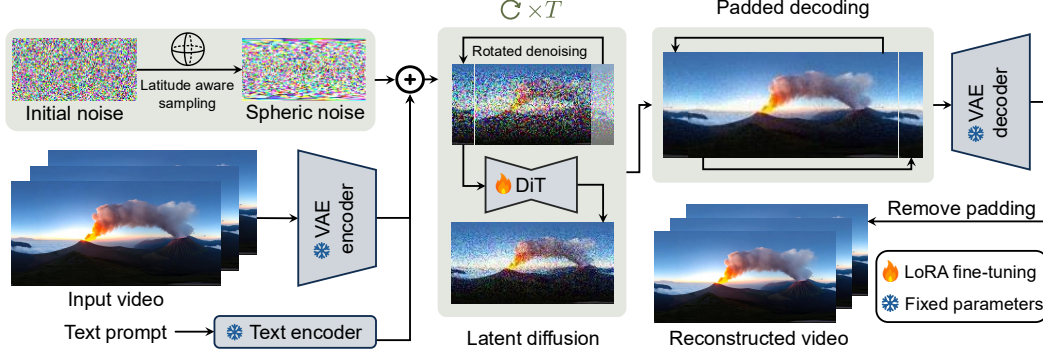
Figure 2: The pipeline of our proposed PanoWan, aware of spherical coordinates. To avoid latitudinal distortion, initial random Gaussian noise is remapped to align with the spherical frequency distribution using the latitude-aware sampling (Sec. 4.2). Next, this remapped noise serves as the latent code within the VAE-encoded latent space. A DiT-based denoising network then iteratively refines this latent representation, where rotated denoising is applied by rolling the latent grid to ensure semantic consistency across longitudinal boundaries. After that, padded pixel-wise decoding provides the VAE decoder with extended context, enabling the mapping of the denoised latent code back into seamless panoramic videos (Sec. 4.3). The DiT backbone within PanoWan is efficiently fine-tuned using LoRA, where most parameters of the pre-trained text-to-video model remain frozen to preserve its strong generative priors.

$$f_{\text{car},y}(x) = f_{\text{sph},\theta}(\varphi)\frac{\mathrm{d}s_\theta}{\mathrm{d}\theta} = R f_{\text{sph},\theta}(\varphi), \quad f_{\text{car},x}(y) = f_{\text{sph},\varphi}(\theta)\frac{\mathrm{d}s_\varphi}{\mathrm{d}\varphi} = R f_{\text{sph},\varphi}(\theta)\cos\theta. \quad (4)$$

Consequently, in polar regions ($|\theta| \approx \frac{\pi}{2}$, namely $y \approx 0$ or $y \approx R$), the horizontal frequency in the Cartesian coordinate becomes near-zero ($f_{\text{car},x}(y) \approx 0$). Such distortion in the horizontal frequency distribution significantly degrades the effectiveness of transferring priors.

## 4.2 Latitude-Aware Mechanisms

**Latitude-aware sampling.** Conventional text-to-video models typically assume independent and identically distributed (i.i.d.) Gaussian noise vectors for each Cartesian coordinate $(x, y)$. To avoid latitudinal distortion in polar regions of ERP, we propose the latitude-aware sampling to better align the initial noise with the spherical frequency distribution for panoramic video generation. As illustrated in the top-left of Fig. 2, our latitude-aware sampling remaps the horizontal sampling coordinates based on latitude to preserve frequency consistency across the sphere. Specifically, after initializing the latent map with i.i.d. Gaussian noise vectors, we calculate the sampling noise by remapping the horizontal sampling coordinate $x$ based on the latitude corresponding to row $y$, and then applying the interpolation:

$$P'(x, y) = \text{Interp}_P\Big(R + (x - R)\cos(\frac{2y + 1 - R}{2R}\pi), y\Big), \quad (5)$$

where $P'(x, y)$ is the interpolated noise vector at coordinate $(x, y)$. Interp$(\cdot)$ is formulated as the interpolation function for normalization:

$$\text{Interp}_P(x, y) = \text{sgn}(\text{BI}(P, x, y))\sqrt{\text{BI}(P^2, x, y)}, \quad (6)$$

where sgn$(\cdot)$ is the sign function, and BI$(P, x, y)$ is the standard bilinear interpolation for vector $P$ at coordinate $(x, y)$. Consequently, the resulting noise vectors sampled by our strategy preserve $\mathbb{E}[P'(x, y)] = 0$ and $\mathbb{E}[\text{Var } P'(x, y)] = 1$, approaching the distribution on which the diffusion models are pre-trained. The proof is given in the supplementary materials.

**Frequency domain analysis.** We aim to prove that the horizontal frequency properties of the proposed sampling correctly represent the inherent properties of the spherical coordinate, following

the methodology of 1-D Discrete Fourier Transform (DFT). Denote the maximal frequency of the original signal in the spherical space as $f_{\max}$, and the maximal frequency in the Cartesian coordinates at the latitude of $\theta$ is determined by:

$$\max f_{\mathrm{car},x}(y) = \max R \cos \theta \cdot f_{\mathrm{sph},\varphi}(\theta) \leq R f_{\max}, \tag{7}$$

where the equation only holds when $\theta = 0$, namely on the equator. For the proposed design, it guarantees $\max f_{\mathrm{car},x} = 2R$ along the equator as $\cos(\theta) = 1$ and the original pixels along $P(\cdot, \frac{R-1}{2})$ is taken. Therefore, the maximal spherical frequency is $f_{\max} = 2$. According to Eq. (5), the warped results only depends on the values between $(R - R\cos(\theta), y)$ and $(R + (R-1)\cos(\theta), y)$ in the original Cartesian grid $P$. According to DFT, the support of the spectrum is reduced to:

$$\lceil R\cos(\theta)\rceil + \lceil (R-1)\cos(\theta)\rceil \approx 2R\cos\theta = f_{\max}R(\theta)\cos\theta = \max f_{\mathrm{car},x}(y), \ \forall \theta \in [-\frac{\pi}{2}, \frac{\pi}{2}]. \tag{8}$$

This matches the inherent property of the panorama in the frequency domain in every latitude.

### 4.3 Longitude Continuity Mechanisms

**Seamless longitude transitions.** Pre-trained conventional text-to-video models lack continuity awareness required between the columns of left and right boundaries. Consequently, applying their generative priors directly for panoramic video generation leads to seam artifacts, resulting in an observable transition where the easternmost and westernmost longitudes meet. To achieve seamless longitude transitions, we recognize that these artifacts arise from both semantic inconsistency and pixel-wise disharmony. This motivates us to propose the rotated denoising and the padded decoding to significantly remove the artifacts.

**Rotated semantic denoising.** The video generation backbone inherently introduces semantic inconsistency at each denoising step. Since the pre-trained generative priors lack the continuity awareness, the semantic in leftmost and rightmost longitudes is typically inconsistent, which are further accumulated during the iterative denoising steps and finally produce an obvious transition.

Our proposed rotated semantic denoising aims to spread the transition error evenly to different longitudes. Let $\mathcal{R}_{s_t}(\cdot)$ be the circular-shift operator and $W$ denote the width of the latent code. As shown in Fig. 2, we horizontally roll the latent code $Z_t$ by $\{s_t = t \bmod W\}$ columns at denoising step $t$ and then undo the shift:

$$Z_{t+1} = \mathcal{R}_{-s_t}\Big( \phi_\theta\big(\mathcal{R}_{s_t}(Z_t)\big)\Big), \tag{9}$$

where $\phi_\theta(\cdot)$ is the noise predictor. As a result, the inherent accumulative error for horizontal coordinate $x$ after $T$ denoising steps is:

$$E_T(x) = \sum_{t=1}^{T} \varepsilon_t\big((x + s_t) \bmod W\big), \tag{10}$$

where $\varepsilon_t(\cdot)$ is prediction error for the transition and step $t$, which would concentrate at a fixed seam if no rotation are applied. Due to the rotation strategy, this error at physical coordinate $x$ at step $t$ is determined by the logical position $\{(x + s_t) \bmod W\}$. Over $T$ steps, these logical coordinates $\{(x + s_t) \bmod W\}_{t=1}^{T}$ ideally approach a uniform permutation for all longitudes. This effectively suppresses seam artifacts by a factor approaching $1/W$.

**Padded pixel-wise decoding.** When decoding latent codes back to the pixel space, the pre-trained VAE decoder D often introduces pixel-wise inconsistencies, as it is trained on conventional videos and lacks awareness of the spatial continuity required across the left-right seam of panoramic videos [20]. To address it, we present the padded pixel-wise decoding. Let $Z_0$ be the denoised latent code. We first create a padded latent code $Z_0' = P_r(Z_0)$, where $P_r(\cdot)$ is a circular padding operator that extends $Z_0$ by $r$ columns of context on the side, and the content at horizontal coordinate $x$ is $\{x \bmod W\}$ in $Z_0$. Finally, we center crop the decoded panoramic videos after the decoding $V = \mathrm{Crop}\big(\mathrm{D}(Z_0')\big)$, as illustrated in Fig. 2. This approach ensures that pixels near the original seam boundaries are decoded with $r$ columns of horizontal panoramic context. Consequently, the VAE decoder can effectively leverage its generative priors learned from conventional videos to avoid the seam artifacts.

6

Table 1: Quantitative comparison results of PanoWan and previous text-based panoramic video generation models. ↑ (↓) means higher (lower) is better. Throughout the paper, best performances are highlighted in **bold**.

| Method | General Metrics | | | Panoramic Metrics | | |
|---|---|---|---|---|---|---|
| | FVD ↓ | VideoCLIP-XL ↑ | Image Quality ↑ | End Continuity ↓ | Motion Pattern ↑ | Scene Richness ↑ |
| 360DVD [32] | 1750.36 | 20.27 | 0.7054 | 0.0323 | 5.8% | 6.6% |
| DynamicScaler [12] | 2146.04 | 21.13 | 0.7188 | 0.0339 | 4.0% | 2.6% |
| Ours (W/o LAS) | 1520.69 | 21.20 | 0.7205 | 0.0278 | 16.2% | 19.4% |
| Ours (W/o RSD) | 1302.48 | 21.76 | 0.7243 | 0.0327 | 15.6% | 18.8% |
| Ours (W/o PPD) | 1294.03 | 21.81 | 0.7239 | 0.0294 | 22.0% | 17.4% |
| **Ours (full)** | **1281.21** | **21.86** | **0.7249** | **0.0270** | **36.4%** | **35.2%** |

## 5 Experiments

### 5.1 Training Details

PanoWan is built on Wan 2.1-1.3B-T2V [33] as the video generation backbone. We train PanoWan at a resolution of $448 \times 896$, closely matching the pre-trained resolution of this backbone model. For parameter-efficient training, LoRA [11] with a rank of 64 is applied to the query, key, value, and output projections of the attention mechanisms, as well as to the feed-forward networks. The model is trained for 200K iterations on our contributed PANOVID dataset. The training process employs the AdamW optimizer [13] with a learning rate of $1 \times 10^{-4}$ and a batch size of 8. Training is conducted on 8 NVIDIA H100 GPUs for approximately 18 hours. During each iteration, clips of 81 consecutive frames are randomly sampled from the videos. Consequently, only 21.9M parameters are adjusted, constituting approximately 1.6% of the base model's total parameters.

### 5.2 Panoramic Video Evaluation Metrics

Existing panoramic video generation methods either directly apply conventional video evaluation metrics [12, 40] or rely on subjective user preferences [32], lacking metrics that comprehensively assess both perceptual quality and spherical consistency critical for panoramic video evaluation. This motivates us to adapt general video quality metrics for panoramic videos and to introduce additional panorama-specific metrics for structural properties of 360° content.

**General metrics.** We apply Frechét Video Distance (FVD) [28] to evaluate overall video quality and VideoCLIP-XL [31] to assess text-video alignment. Following DynamicScalar [12], we also calculate specific metrics for image quality. To adapt these general metrics for panoramic videos, we project each video onto a cube map and compute metric scores separately on each of the six faces. The final reported score for a video $v$ is a weighted average:

$$\overline{\mathbf{f}}(v) = \sum_{f \in \mathcal{F}} \alpha_f \cdot \Phi\big(\mathcal{P}_f(v)\big), \tag{11}$$

where $\mathcal{F}$ denotes the set of cube map faces, $\mathcal{P}_f(v)$ is the projection of video $v$ onto face $f \in \mathcal{F}$, $\Phi$ is the metric function, and $\alpha_f$ is the weight assigned to face $f$. Following OmniFID [6], we assign weights $\alpha_{\text{top}} = \alpha_{\text{bottom}} = \frac{1}{3}$ and $\alpha_{\text{side}} = \frac{1}{12}$ for each of the four lateral faces.

**Panoramic metrics.** Following previous works [32, 12], we evaluate motion patterns and scene richness with user preferences. We additionally introduce a quantitative metric for evaluating the end continuity of generated panoramic videos, tailored to capture artifacts across longitude boundaries. Specifically, this metric computes the mean absolute pixel difference across the left and right boundaries, directly capturing discontinuities at the longitude seam.

### 5.3 Comparison with State-of-the-art Methods

We evaluate PanoWan against existing text-based panoramic video generation methods, including 360DVD [32] and DynamicScaler [12]. Quantitatively, as shown in Tab. 1, PanoWan achieves state-of-the-art performance across both general and panoramic metrics (detailed in Sec. 5.2). Qualitatively,

*Wide-angle panoramic interior capturing an energetic hot pot restaurant bustling with diners enthusiastically cooking ingredients in bubbling pots. Colorful plates arranged invitingly, steam rising, spirited conversations filling warm, inviting atmosphere emphasizing community and culinary enjoyment.*



*Expansive panoramic view capturing a vibrant ski resort nestled among towering snowy mountains, as skiers gracefully descend pristine slopes amid cozy alpine chalets. Chairlifts glide leisurely under crisp blue skies while visitors gather around lively outdoor cafes, soaking in the sunny winter scenery.*

Figure 3: Visual comparison results with existing text-based panoramic video generation methods.

we present visual results to highlight our advantages. For instance, DynamicScaler [12] falls short in complex scenarios (Fig. 3, first sample), and 360DVD [32] exhibits notable distortion in polar regions (Fig. 3, second sample). In contrast, PanoWan effectively maintains global consistency and visual coherence, achieving superior performance in generating high-fidelity panoramic videos.

## 5.4 Ablation Studies

We conduct ablation studies to validate the effectiveness of proposed modules in PanoWan: latitude-aware sampling (LAS), rotated semantic denoising (RSD), and padded pixel-wise decoding (PPD).

**Quantitative results.** As shown in Tab. 1, removing LAS primarily affects general metrics—FVD increases from 1281.21 to 1520.69, and VideoCLIP-XL drops from 21.86 to 21.20—indicating the model struggles to learn panoramic features in high-latitude regions without frequency-aligned noise initialization. In contrast, removing RSD or PPD mainly degrades panoramic metrics (*e.g.*, end continuity increases from 0.0270 to 0.0327 and 0.0294, respectively), confirming their roles in achieving seamless longitude transitions.

**Qualitative results.** We further provide qualitative evaluation in Fig. 4. When generating high-latitude elements like LED panels (which should appear straight in perspective views but are inherently distorted in the equirectangular projection), PanoWan without LAS fails to render them with the correct geometric appearance. When RSD is discarded, semantic inconsistencies become apparent at the longitude seam, due to the lack of mechanism for continuity awareness and the error accumulation during the denoising process. When PPD is removed, observable seam artifacts occur because
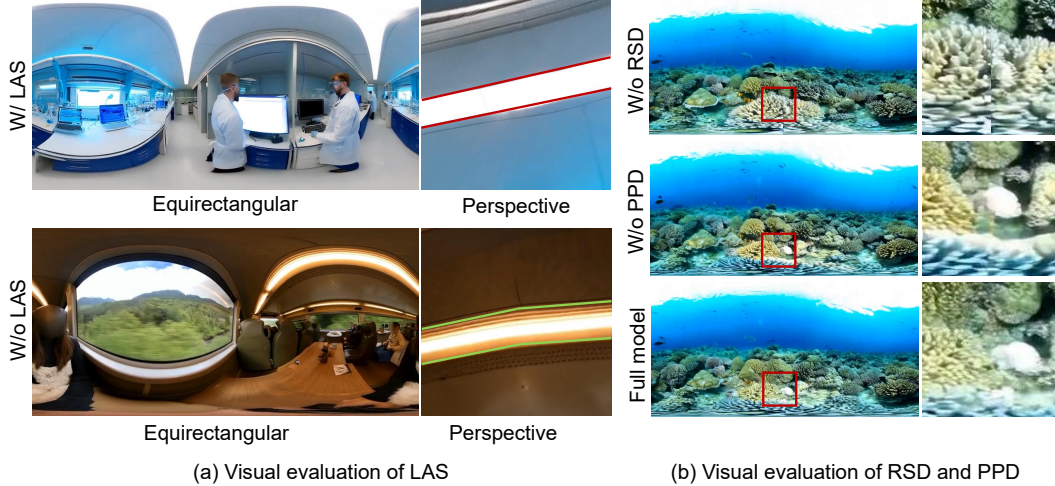
Figure 4: Qualitative evaluation of proposed latitude/longitude-aware mechanisms. (a) With the proposed Latitude-Aware Sampling (LAS), PanoWan ensures that content generated at high latitudes exhibits an accurate geometry when presented in a perspective view. (b) By combining Rotated Semantic Denoising (RSD) and Padded Pixel-wise Decoding (PPD), PanoWan achieves seamless longitude transitions. For visualization, videos are rolled $180°$ to center the seam.

conventional VAE decoder introduces pixel-wise inconsistencies at boundaries. Consequently, the full PanoWan model with all modules enabled achieves the best performance.

## 5.5 Application

As a text-based panoramic video model, PanoWan shows robust zero-shot capabilities across a wide range of downstream tasks. We present representative examples in Fig. 1 and additional examples in supplementary materials due to the space limitation.

**Long video generation.** To generate long panoramic videos, we employ a local windowing strategy within the latent space. Specifically, at each denoising step, we partition the latent code into temporally overlapping chunks. These chunks are processed independently and then seamlessly merged using a linear blending function on the overlapping segments to ensure smooth temporal transitions.

**Super-resolution for panoramic videos.** To generate high-resolution panoramic videos from low-resolution ones, we first encode each low-resolution video into its corresponding latent code. After injected noise, the latent code is denoised based on user-provided text descriptions, producing results with structural consistency and visual fidelity across the spherical representation.

**Inpainting for semantic editing.** Given a panoramic video, we identify and mask regions for modification. Next, we apply the denoising process to these regions, guided by user-provided text descriptions. Leveraging its understanding of spherical representations, the inpainted content naturally exhibits the properties of ERP projection.

**Outpainting for conventional videos.** Similar to the inpainting process, we first map the conventional video to the latent code and then mask the surrounding unseen panoramic regions. With user-provided text descriptions, we denoise the masked regions to generate corresponding content. Our pre-trained model maintains the spatial and temporal consistency for generated panoramic videos.

## 6 Conclusion

We present PanoWan, a text-based panoramic video generation framework that effectively lifts pre-trained diffusion model to the panorama. By integrating latitude-aware sampling, PanoWan addresses latitudinal distortions caused by equirectangular projection. Equipped with the rotated semantic denoising and the padded pixel-wise decoding, PanoWan achieves the seamless longitude transitions. To provide large-scale data for lifting representations from conventional videos to the

panorama, we contribute the PANOVID dataset, offering high-quality and semantically rich 360° video data with annotated text descriptions. Extensive experiments demonstrate that PanoWan achieves state-of-the-art performance on text-based panoramic video generation and strong generalization across diverse zero-shot downstream tasks.

**Limitation.** While PanoWan benefits from the strong priors of its pretrained text-to-video models, it is also inherits a common challenge: the content forgetting problem often seen in such models. This issue is particularly evident when generating long videos due to limited temporal memory. We believe this challenge can be substantially alleviated through future advancements in memory-aware generation techniques (*e.g.*, video caching mechanisms).

# 7   Appendix

## 7.1   Proof of Noise Distribution Preservation

To align the initial noise with the spherical frequency distribution and avoid latitudinal distortion in polar regions of ERP, Sec. 4.2 proposes the *latitude-aware sampling* strategy. For clarity, we recall the noise at each coordinate after remapping the horizontal sample coordinate $x$ based on latitude $y$, as stated in Eq. (5) of the main paper:

$$P'(x,y) = \text{Interp}_P\Big(R + (x-R)\cos(\frac{2y+1-R}{2R}\pi), y\Big), \tag{12}$$

To best exploit the diffusion prior, it is desired to have $\mathbb{E}[P'(x,y)] = 0$ and $\mathbb{E}[\text{Var } P'(x,y)] = 1$. First, we provide $\mathbb{E}[P'(x,y)] = 0$ as follows:

$$\mathbb{E}[P'(x,y)] = \mathbb{E}_{x,y}\Big[\text{sign}(\text{BI}(P,x,y))\sqrt{\text{BI}(P^2,x,y)}\Big] \tag{13}$$

$$= \mathbb{E}_{P_{ij}\sim\mathcal{N}(0,1)}\Big[\text{sign}\Big(\sum w_{ij}P_{ij}\Big)\sqrt{\sum w_{ij}P_{ij}^2}\Big] \tag{14}$$

$$\overset{\tilde{P}_{ij}:=-P_{ij}}{=\!=\!=\!=\!=} \mathbb{E}_{(-\tilde{P}_{ij})\sim\mathcal{N}(0,1)}\Big[\text{sign}\Big(\sum w_{ij}(-\tilde{P}_{ij})\Big)\sqrt{\sum w_{ij}(-\tilde{P}_{ij})^2}\Big] \tag{15}$$

$$= -\mathbb{E}_{\tilde{P}_{ij}\sim\mathcal{N}(0,1)}\Big[\text{sign}\Big(\sum w_{ij}(P_{ij})\Big)\sqrt{\sum w_{ij}P_{ij}^2}\Big] = -\mathbb{E}\left[P'(x,y)\right]. \quad \square \tag{16}$$

After that, we prove $\mathbb{E}[\text{Var } P'(x,y)] = 1$ as follows:

$$\mathbb{E}[\text{Var } P'(x,y)] = \mathbb{E}\left[P'(x,y)^2\right] - \left(\mathbb{E}\left[P'(x,y)\right]\right)^2 = \mathbb{E}\left[P'(x,y)^2\right] \tag{17}$$

$$= \mathbb{E}_{P_{ij}\sim\mathcal{N}(0,1)}\left[\sum_{i,j\in\{0,1\}} w_{ij}P_{ij}^2\right] = \sum_{i,j\in\{0,1\}} w_{ij}\mathbb{E}_{P_{ij}\sim\mathcal{N}(0,1)}\left[P_{ij}^2\right] \tag{18}$$
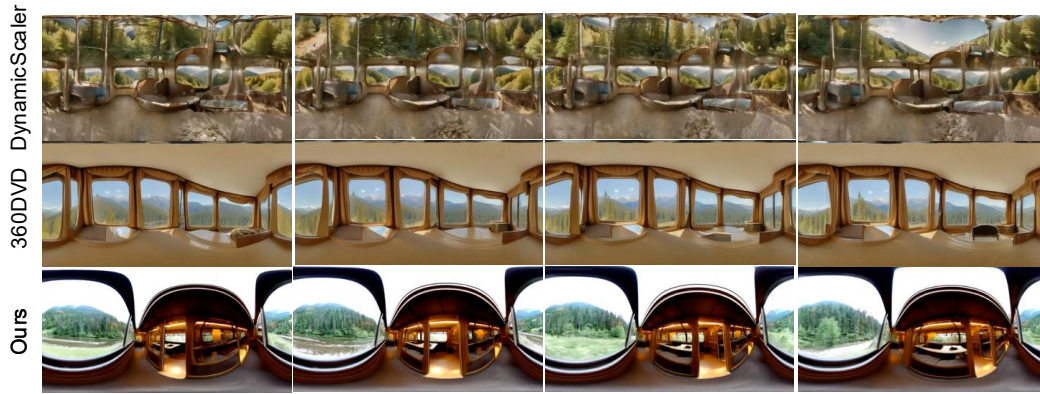
$$= \sum_{i,j\in\{0,1\}} w_{ij} = 1. \quad \square \tag{19}$$

Note that the first equation in Eq. (19) is possible only because $P_{ij}$ is independent and identically distributed.

## 7.2   Additional Experiment Results

In this section, we provide additional comparison results between PanoWan and existing text-based panoramic video generation methods [12, 32], where PanoDiT [40] is omitted as its code is unavailable. We also present additional examples showcasing applications such as long video generation, super-resolution, semantic inpainting, and video outpainting. Finally, we provide a detailed discussion of failure cases.
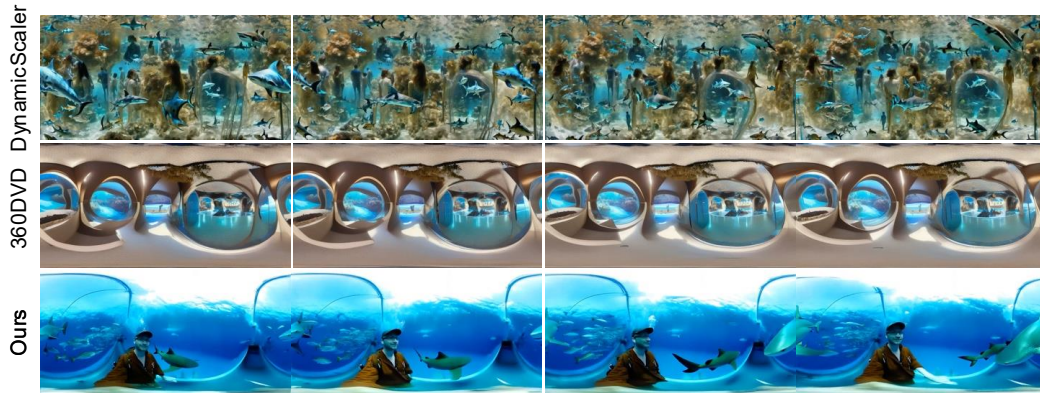
**Additional comparison results.** As shown in Fig. 5, DynamicScaler [12] suffers from a limited local denoising window, resulting in globally inconsistent content with repeated semantic elements appearing in different regions. 360DVD [32] often produces observable artifacts in high-latitude

*Wide-angle panoramic view from inside a luxurious vintage train wagon as it gently traverses through spectacular landscapes. Mountains, forests, and meandering rivers captured clearly through panoramic windows, accentuated by warm interior lighting and cozy compartments, providing a relaxing and nostalgic atmosphere.*
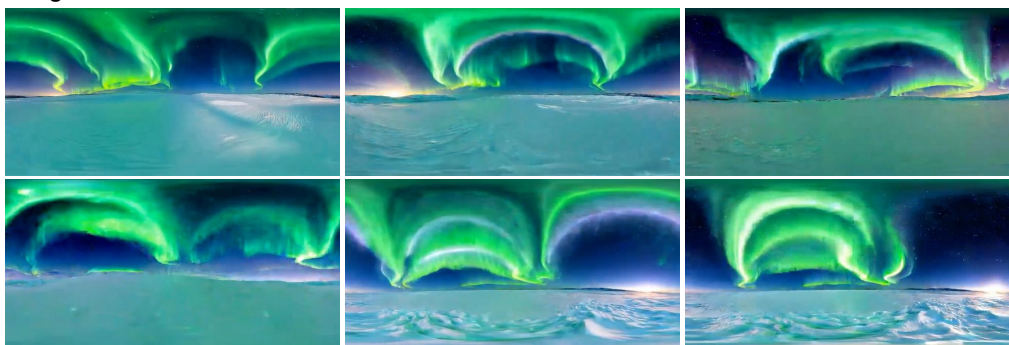


*Stunning panoramic underwater shot of a vibrant coral reef ecosystem brimming with marine life. Colorful fish dart effortlessly among intricate coral formations, soft rays of sunlight filter through the crystal-clear waters, creating mesmerizing patterns on the ocean floor. Wide-angle capturing vivid hues and abundant biodiversity.*



*Immersive panoramic view inside an elongated aquarium tunnel, visitors walking beneath a transparent underwater canopy surrounded by vibrant fish, graceful manta rays, and large sharks moving serenely through clear azure waters, producing a compelling sense of underwater wonder and tranquility.*

Figure 5: Additional comparison results with existing text-based panoramic video generation methods.
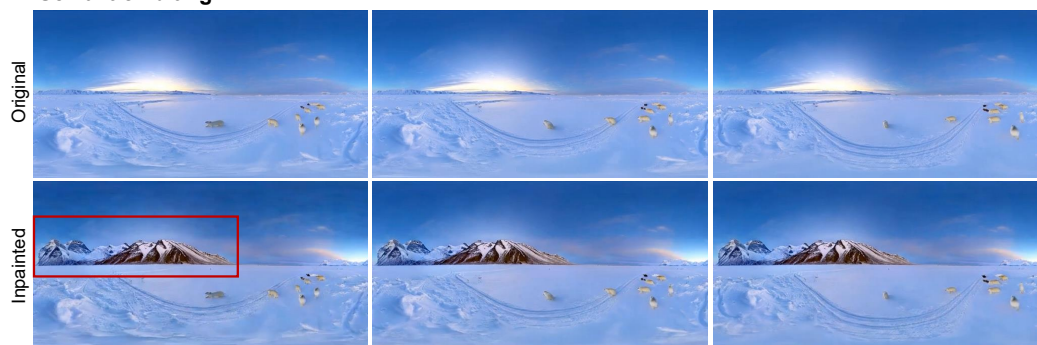
**Long Video**



*Majestic panoramic shot capturing vivid green and violet northern lights gracefully illuminating a quiet, snow-covered tundra beneath a star-studded night sky. Gentle, fluid ribbons of color dance overhead in mesmerizing motion, creating an awe-inspiring spectacle. Ultra-wide landscape shot emphasizing grandeur and serenity.*

**Super-resolution**



*360-degree panoramic interior view inside a charming artisan bakery bustling with activity, bakers carefully preparing handcrafted breads, pastries, and desserts. Shelves stocked with warm baked goods, aromatic scents filling the air, creating feelings of warmth, comfort, and culinary delight.*

**Semantic Editing**
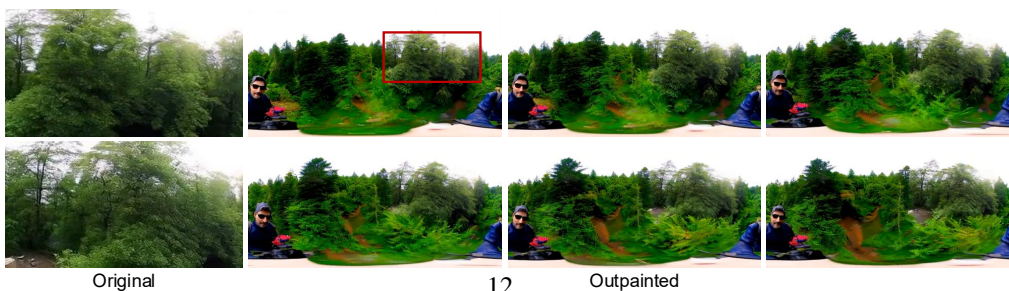


*Add a snowy mountain.*

**Video outpainting**



Original                                          Outpainted

Figure 6: Additional application results, showcasing the zero-shot capabilities for downstream tasks.

*Vibrant panoramic shot inside an animal adoption center where families interact warmly with playful puppies, kittens, and other shelter animals. Staff members guide hopeful adopters amid cheerful play areas, conveying compassion, hope, and care.*



*Inside a bustling Starbucks, a young woman sits by the window, sipping a grande latte, engrossed in a thick novel. Sunlight filters through, casting warm glows on her focused face. Surrounding her are chic wooden interiors, the aroma of freshly brewed coffee, and the chatter of patrons. Medium shot, capturing the vibrant cafe ambiance.*
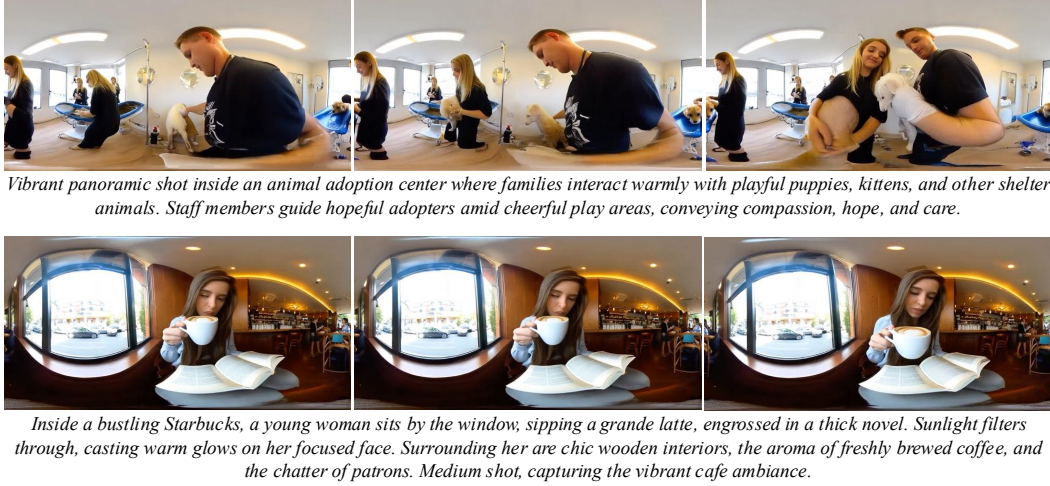
Figure 7: Visualization of failure cases.

regions and exhibits relatively limited scene consistency. In contrast, our PanoWan achieves the most coherent and visually consistent results across diverse scenarios.

**Additional application results.** We provide additional examples across four application scenarios. As shown in Fig. 6, (a) PanoWan enables long video generation while maintaining consistent semantics throughout extended temporal durations. (b) For the super-resolution, directly applying Wan 2.1 [33] leads to severe artifacts in high-latitude regions, whereas PanoWan produces consistent and artifact-free results across all latitudes. (c) and (d) further demonstrate its effectiveness in semantic editing and video outpainting, respectively. These results highlight the strong potential of PanoWan as a versatile model for high-quality panoramic video generation and editing.

**Failure cases.** PanoWan can occasionally exhibit failures in certain scenarios. As illustrated in the top row of Fig. 7, the generated pet dog exhibits inconsistent features. In the bottom row, the book being read by the woman appears with two spines and three pages, deviating from real-world structure and common sense. We believe these failure cases are not primarily related to panoramic properties but are largely inherited from the pre-trained text-to-video model [33]. As the backbone model improves, these failure cases will be improved.

## 7.3 Dataset Details

Existing panoramic video datasets [32, 27, 20] are limited in scale or lack paired text captions. Recognizing the need for a dataset that supports effective and scalable training of text-based panoramic generation models, we collected our PANOVID dataset. In this section, we first present the data sources, followed by an efficient and accurate filtering pipeline designed to support the scalable collection of panoramic videos with diverse and paired captions. After that, we provide an analysis of the semantic distribution and showcase additional examples from the final dataset.

**Data collection.** We incorporate videos from large-scale collections (*e.g.*, 360-1M [30] and Imagine360 [27]), which are primarily composed of user-uploaded YouTube content. These sources offer vast quantity of data but relatively lower visual quality. To compensate for these quality concerns, more curated datasets are incorporated (*e.g.*, WEB360 [32], 360+x [5], Panonut360 [37], and the Miraikan 360-degree Video Dataset [1]), which are fine-grained but limited in volume.

**Data filtering.** To efficiently filter and curate high-quality panoramic video clips from these varied noisy sources, we design a five-stage filtering pipeline based on a vision-language model. This pipeline ensures that only relevant and high-quality content is selected for further processing:

- **Initial filtering by popularity.** From the 360-1M dataset, we retain only videos with at least 1000 views. This heuristic filters out low-engagement content, which is often associated with poor quality or uninformative scenes.
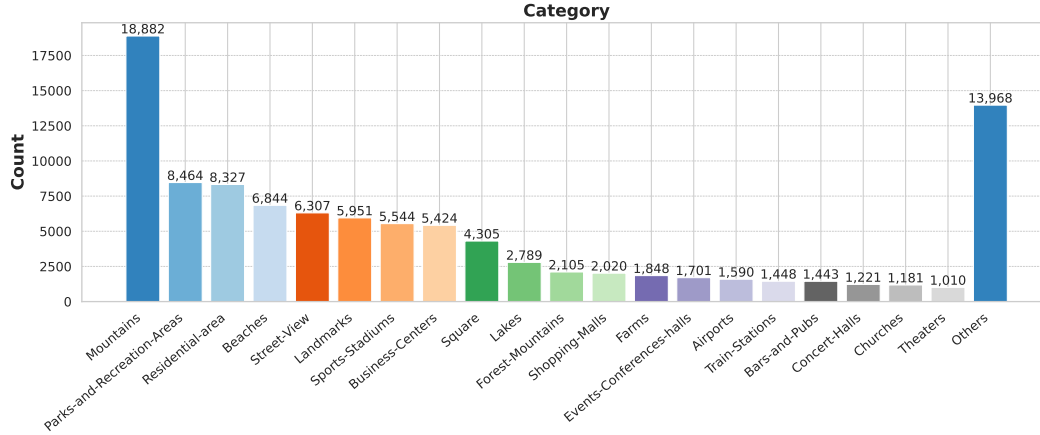
Figure 8: Category distribution of PANOVID dataset before balancing the semantics.

- **Video segmentation.** We segment each video into 10-second clips using the PySceneDetect library, ensuring that each clip contains a single continuous scene without abrupt transitions.
- **Vision-language-based annotation.** For each segmented clip, we employ the Qwen-2.5-VL [2] model to generate structured annotations in JSON format. The model receives the following prompt:

```
Please analyze this video and provide your response in JSON format with
    the following fields:

1. 'caption': A detailed description of what's happening in the video. Do
    not show your analysis, just describe what you see. Do not start with
    "The video shows", describe the video itself as a whole. Include the
    panoramic statement like "Panoramic view of ..." or "360 degree view
    of ..." if it is a panoramic video. Here are some examples:
  - 'Panoramic shot of colorful hot air balloons gracefully ascend,
      floating over lush green fields, their vibrant hues contrasting
      against a vast, cloud-dappled blue sky. Gentle breezes propel them
      in a serene dance, casting dynamic shadows on the verdant landscape
       below. Wide shot from ground level, capturing the expansive scene.'

  - 'Panoramic shot of an active volcano spewing smoky plumes against a
      fiery sunset sky, majestic mountains shrouded in misty clouds in
      the foreground, creating a breathtaking contrast. Camera pans
      slowly, capturing the vastness and awe-inspiring beauty of nature.'
  - 'Aerial perspective of vibrant fireworks blossoming in the ink-black
      sky, casting shimmering lights over a sprawling urban landscape
      below. Mesmerizing pyrotechnics burst in various colors and
      patterns against the starless expanse, illuminating cityscapes with
       transient brilliance. Wide shot from a plane window, capturing the
      nocturnal city alive under the grand firework spectacle.'
2. 'is_panorama': A boolean (true or false) indicating whether this is a
    360-degree ERP projected video.
3. 'poi_category': A list of strings indicating the points of interest in
    the video. If there are no points of interest, set this to an empty
    list. If there are multiple points of interest, describe each of them
    in a string, sorted by their importance. You should use the available
    POI categories. In case none of the provided POI categories can
    describe the video, you may return a succinct word in the similar
    pattern as given categories. For example:
  - ['Coffee-Shop']
  - ['Mountains', 'Lakes']

Your response should be valid JSON string, like this:
```

14

```json
{
  "caption": "Panoramic shot of colorful hot air balloons gracefully ascend
      , floating over lush green fields, their vibrant hues contrasting
      against a vast, cloud-dappled blue sky. Gentle breezes propel them
      in a serene dance, casting dynamic shadows on the verdant landscape
      below. Wide shot from ground level, capturing the expansive scene.",
  "is_panorama": true,
  "poi_category": ["Mountains"]
}
```

```
Note that the video may be compressed with limited fps to reduce uploaded
    file size. Your response in 'caption' should not include any
    description of the video quality or compression. Just focus on the
    content of the video.

Here are the available POI categories: "Restaurant, Coffee-Shop, Bars-and-
    Pubs, Residential-area, Hotels-Motels, Vaccation-Rentals, Hospitals-
    Clinics, Pharmacies, Dentists, School-Universities, Library,
    Supermarkets, Shopping-Malls, Clothing-Stores, Shoe-Stores, Bookstores,
     Flowerstore, Furniture-Stores, Electorical-Store, Pet-Store, Toy-Shop,
     Airports, Train-Stations, Bus-Stops, Gas-Station, Car-Rental-Agencies,
     Theaters, Concert-Halls, Sports-Stadiums, Parks-and-Recreation-Areas,
     Museums, Art-Galleries, Zoos-Aquariums, Botanical-Gardens, Landmarks,
     Cultural-Centers, Post-Offices, Police-Stations, Courthouses, City-
    Halls, Banks-ATMs, Events-Conferences-halls, Beaches, Hiking-Trails,
    Campgrounds, Lakes, Mountains, Forest-Mountains, Farms, Street-View,
    Square, Business-Centers, Tech-Companies, Co-working-Spaces, Gyms-and-
    Fitness-Centers, Sports-Clubs, Swimming-Pools, Tennis-Courts, Auto-
    Repair-Shops, Car-Washes, Parking-Lots, Churches, Mosques, Temples,
    Graveyards."
```

The model's responses are used to filter out clips that are not panoramic videos in ERP format. Note that the POI (Points Of Interest) categories follow the setup of DL3DV [17].

- **Motion score filtering.** We compute a normalized motion score for each clip based on optical flow magnitude (following [7]). Only clips with motion scores above 0.4 are retained, to avoid including static or nearly still scenes that are less informative for training generative models.

- **Aesthetic score filtering.** Each frame is scored using Q-Align [34] for visual aesthetics, and the clip's final score is computed as the average of all frame scores. We discard clips with aesthetic scores below 3, ensuring the training data maintains a minimum level of visual quality.

This scalable filtering pipeline enables the construction of a large, diverse, and high-quality panoramic video dataset from noisy web-scale sources, thereby effectively lifting pre-trained text-to-video generation models to the panorama.

**Semantic distribution.** We observe that the POI labels generated through the proposed pipeline exhibit a highly imbalanced distribution, as shown in Fig. 8. A few categories (*e.g.*, *Mountains*, *Parks-and-Recreation-Areas*, and *Residential-area*) dominate the dataset with thousands of samples, while many others (*e.g.*, *Bookstores*, *Dentists*, and *Tennis-Courts*) have fewer than 50 instances. This imbalanced distribution poses a challenge for training panoramic video generation models that require semantic diversity and balanced representation. To address this issue and ensure a more balanced semantic coverage, we select up to 200 video clips with the highest aesthetic quality for each category. For categories with fewer than 200 clips, all available samples are retained. By combining this filtered set with existing curated datasets, PANOVID contains over 13K high-quality clips with diverse scene types and detailed captions.

**Dataset examples.** PANOVID dataset covers a wide range of scene categories (*e.g.*, natural landscapes, urban environments, and indoor scenarios), with paired detailed captions. As shown in Fig. 9, we provide representative samples from the PANOVID dataset to illustrate its diversity and quality.
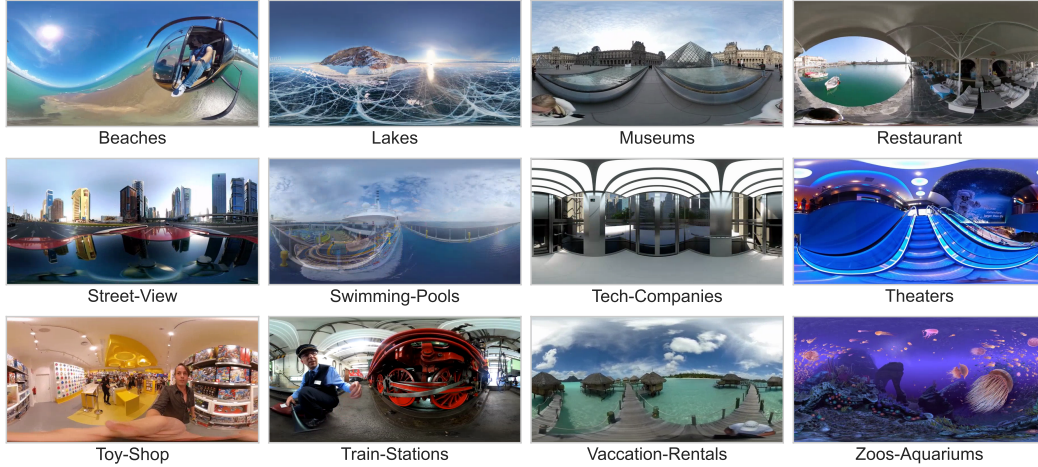
Figure 9: Representative samples from the PANOVID dataset.

# References

[1] Miraikan 360-degree video dataset. `https://www.miraikan.jst.go.jp/en/research/AccessibilityLab/dataset360/`.

[2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2.5-VL technical report. *arXiv preprint arXiv:2502.13923*, 2025.

[3] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, 2021.

[4] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.

[5] Hao Chen, Yuqi Hou, Chenyuan Qu, Irene Testini, Xiaohan Hong, and Jianbo Jiao. 360+x: A panoptic multi-modal scene understanding dataset. In *Proc. of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.

[6] Anders Christensen, Nooshin Mojab, Khushman Patel, Karan Ahuja, Zeynep Akata, Ole Winther, Mar Gonzalez-Franco, and Andrea Colaco. Geometry fidelity for spherical images. In *Proc. of European Conference on Computer Vision*, 2024.

[7] Gunnar Farnebäck. Two-frame motion estimation based on polynomial expansion. In *Image Analysis*, 2003.

[8] Hao He, Ceyuan Yang, Shanchuan Lin, Yinghao Xu, Meng Wei, Liangke Gui, Qi Zhao, Gordon Wetzstein, Lu Jiang, and Hongsheng Li. CameraCtrl II: Dynamic scene exploration via camera-controlled video diffusion models. *arXiv preprint arXiv:2503.10592*, 2025.

[9] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation. *arXiv preprint arXiv:2211.13221*, 2022.

[10] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

[11] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *Proc. of the International Conference on Learning Representations*, 2022.

[12] Liu Jinxiu, Lin Shaoheng, Li Yinxiao, and Yang Ming-Hsuan. DynamicScaler: Seamless and scalable video generation for panoramic scenes. In *Proc. of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.

[13] Diederik P Kingma, J Adam Ba, and J Adam. A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2020.

[14] 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.

[15] Benjamin J. Li, Jeremy N. Bailenson, Adam Pines, Walter J. Greenleaf, and Leanne M. Williams. A public database of immersive vr videos with corresponding ratings of arousal, valence, and correlations between head movements and self report measures. *Frontiers in Psychology*, 2017.

[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.

[17] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin, Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu, et al. DL3DV-10K: A large-scale scene dataset for deep learning-based 3D vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22160–22169, 2024.

[18] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.

[19] Taiming Lu, Tianmin Shu, Junfei Xiao, Luoxin Ye, Jiahao Wang, Cheng Peng, Chen Wei, Daniel Khashabi, Rama Chellappa, Alan Yuille, and Jieneng Chen. GenEx: Generating an explorable world. *Proc. of the International Conference on Learning Representations*, 2025.

[20] Rundong Luo, Matthew Wallingford, Ali Farhadi, Noah Snavely, and Wei-Chiu Ma. Beyond the frame: Generating 360° panoramic videos from perspective videos. *arXiv preprint arXiv:2504.07940*, 2025.

[21] Jingwei Ma, Erika Lu, Roni Paiss, Shiran Zada, Aleksander Holynski, Tali Dekel, Brian Curless, Michael Rubinstein, and Forrester Cole. VidPanos: Generative panoramic videos from casual panning videos. In *Proc. of ACM SIGGRAPH Asia*, 2024.

[22] Minho Park, Taewoong Kang, Jooyeol Yun, Sungwon Hwang, and Jaegul Choo. SphereDiff: Tuning-free omnidirectional panoramic image and video generation via spherical latent representation. *arXiv preprint arXiv:2504.14396*, 2025.

[23] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, 2023.

[24] Xuanchi Ren, Tianchang Shen, Jiahui Huang, Huan Ling, Yifan Lu, Merlin Nimier-David, Thomas Müller, Alexander Keller, Sanja Fidler, and Jun Gao. GEN3C: 3D-informed world-consistent video generation with precise camera control. In *Proc. of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.

[25] SeaweadTeam, Ceyuan Yang, Zhijie Lin, Yang Zhao, Shanchuan Lin, Zhibei Ma, Haoyuan Guo, Hao Chen, Lu Qi, Sen Wang, et al. Seaweed-7B: Cost-effective training of video generation foundation model. *arXiv preprint arXiv:2504.08685*, 2025.

[26] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. In *ICLR*, 2023.

[27] Jing Tan, Shuai Yang, Tong Wu, Jingwen He, Yuwei Guo, Ziwei Liu, and Dahua Lin. Imagine360: Immersive 360 video generation from perspective anchor. *arXiv preprint arXiv:2412.03552*, 2024.

[28] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018.

[29] Veo-Team, :, Agrim Gupta, Ali Razavi, Andeep Toor, Ankush Gupta, Dumitru Erhan, Eleni Shaw, Eric Lau, Frank Belletti, Gabe Barth-Maron, Gregory Shaw, Hakan Erdogan, Hakim Sidahmed, Henna Nandwani, Hernan Moraldo, Hyunjik Kim, Irina Blok, Jeff Donahue, José Lezama, Kory Mathewson, Kurtis David, Matthieu Kim Lorrain, Marc van Zee, Medhini Narasimhan, Miaosen Wang, Mohammad Babaeizadeh, Nelly Papalampidi, Nick Pezzotti, Nilpa Jha, Parker Barnes, Pieter-Jan Kindermans, Rachel Hornung, Ruben Villegas, Ryan Poplin, Salah Zaiem, Sander Dieleman, Sayna Ebrahimi, Scott Wisdom, Serena Zhang, Shlomi Fruchter, Signe Nørly, Weizhe Hua, Xinchen Yan, Yuqing Du, and Yutian Chen. Veo 2. 2024. URL `https://deepmind.google/technologies/veo/veo-2/`.

[30] Matthew Wallingford, Anand Bhattad, Aditya Kusupati, Vivek Ramanujan, Matt Deitke, Aniruddha Kembhavi, Roozbeh Mottaghi, Wei-Chiu Ma, and Ali Farhadi. From an image to a scene: Learning to imagine the world from a million 360° videos. In *Proc. of Neural Information Processing Systems*, 2024.

[31] Jiapeng Wang, Chengyu Wang, Kunzhe Huang, Jun Huang, and Lianwen Jin. VideoCLIP-XL: Advancing long description understanding for video clip models. *arXiv preprint arXiv:2410.00741*, 2024.

[32] Qian Wang, Weiqi Li, Chong Mou, Xinhua Cheng, and Jian Zhang. 360DVD: Controllable panorama video generation with 360-degree video diffusion model. In *Proc. of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.

[33] WanTeam, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.

[34] Haoning Wu, Zicheng Zhang, Weixia Zhang, Chaofeng Chen, Chunyi Li, Liang Liao, Annan Wang, Erli Zhang, Wenxiu Sun, Qiong Yan, Xiongkuo Min, Guangtai Zhai, and Weisi Lin. Q-Align: Teaching LMMs for visual scoring via discrete text-defined levels. *arXiv preprint arXiv:2312.17090*, 2023.

[35] Zeqi Xiao, Yushi Lan, Yifan Zhou, Wenqi Ouyang, Shuai Yang, Yanhong Zeng, and Xingang Pan. WORLDMEM: Long-term consistent world simulation with memory. *arXiv preprint arXiv:2504.12369*, 2025.

[36] Kevin Xie, Amirmojtaba Sabour, Jiahui Huang, Despoina Paschalidou, Greg Klar, Umar Iqbal, Sanja Fidler, and Xiaohui Zeng. VideoPanda: Video panoramic diffusion with multi-view attention. *arXiv preprint arXiv:2504.11389*, 2025.

[37] Yutong Xu, Junhao Du, Jiahe Wang, Yuwei Ning, Sihan Zhou, and Yang Cao. Panonut360: A head and eye tracking dataset for panoramic video. In *Proceedings of the 15th ACM Multimedia Systems Conference*, 2024.

[38] 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. In *ICLR*, 2025.

[39] Lijun Yu, José Lezama, Nitesh B. Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, Alexander G. Hauptmann, Boqing Gong, Ming-Hsuan Yang, Irfan Essa, David A. Ross, and Lu Jiang. Language model beats diffusion – tokenizer is key to visual generation. In *ICLR*, 2024.

[40] Muyang Zhang, Yuzhi Chen, Rongtao Xu, Changwei Wang, JinMing Yang, Weiliang Meng, Jianwei Guo, Huihuang Zhao, and Xiaopeng Zhang. PanoDit: Panoramic videos generation with diffusion transformer. In *Proc. of the AAAI Conference on Artificial Intelligence*, 2025.

[41] Haiyang Zhou, Wangbo Yu, Jiawen Guan, Xinhua Cheng, Yonghong Tian, and Li Yuan. Holotime: Taming video diffusion models for panoramic 4d scene generation, 2025. URL `https://arxiv.org/abs/2504.21650`.