

StereoDiff: Stereo-Diffusion Synergy for Video Depth Estimation

Haodong Li^{1,2} Chen Wang¹ Jiahui Lei¹ Zhiyang Dou^{1,3}

Kostas Daniilidis¹ Jiatao Gu⁴ Lingjie Liu¹

¹University of Pennsylvania; ²HKUST(GZ); ³University of Hong Kong ⁴Apple

{hdli, chenw30, leijh, zydot, lingjie.liu}@seas.upenn.edu;

kostas@cis.upenn.edu; jiatao@apple.com

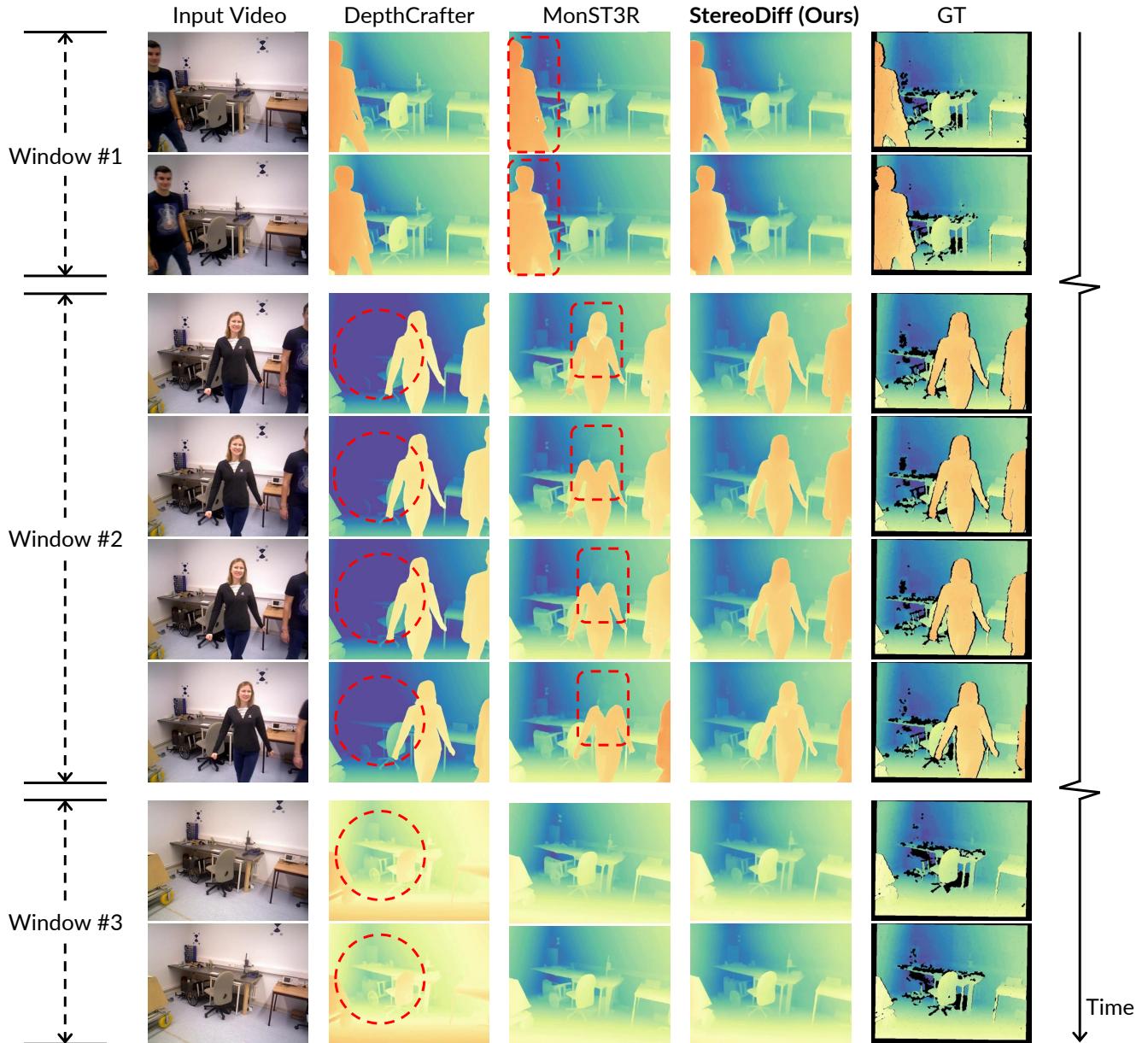


Figure 1. **StereoDiff excels in delivering remarkable global and local consistency for video depth estimation.** In terms of global consistency, StereoDiff achieves highly accurate and stable depth maps on static backgrounds across consecutive windows, leveraging stereo matching to prevent the abrupt depth shifts often seen in DepthCrafter [26], where depth values on static backgrounds can vary significantly between adjacent windows. For local consistency, StereoDiff yields much smoother, flicker-free depth values across consecutive frames, especially in dynamic regions. In contrast, MonST3R [86] suffers from frequent, pronounced flickering and jitters in these areas.

Abstract

Recent video depth estimation methods achieve great performance by following the paradigm of image depth estimation, i.e., typically fine-tuning pre-trained video diffusion models with massive data. However, we argue that video depth estimation is not a naive extension of image depth estimation. The temporal consistency requirements for dynamic and static regions in videos are fundamentally different. Consistent video depth in static regions, typically backgrounds, can be more effectively achieved via stereo matching across all frames, which provides much stronger global 3D cues. While the consistency for dynamic regions still should be learned from large-scale video depth data to ensure smooth transitions, due to the violation of triangulation constraints. Based on these insights, we introduce **StereoDiff**, a two-stage video depth estimator that synergizes stereo matching for mainly the static areas with video depth diffusion for maintaining consistent depth transitions in dynamic areas. We mathematically demonstrate how stereo matching and video depth diffusion offer complementary strengths through frequency domain analysis, highlighting the effectiveness of their synergy in capturing the advantages of both. Experimental results on zero-shot, real-world, dynamic video depth benchmarks, both indoor and outdoor, demonstrate StereoDiff’s SOTA performance, showcasing its superior consistency and accuracy in video depth estimation.

1. Introduction

Monocular video depth estimation is a foundational task in 3D computer vision. Particularly after the hot trend of leveraging pre-trained Stable Diffusion (SD) [46] for image depth prediction [17, 21, 27, 36, 72], e.g., Marigold [27] and Lotus [21], we have witnessed emerging attentions on video depth estimation in the community [26, 30, 53, 64, 73, 77, 86]. Many of them fine-tune the Stable Video Diffusion (SVD) [4] using large-scale video depth data, e.g., ChronoDepth [53] and DepthCrafter [26]. However, most previous methods [26, 53, 67, 73, 76, 77] consider the video depth estimation merely as a video version of image depth estimator, directly modeling a mapping function from the RGB video distribution to the video depth distribution, similar to previous image depth methods that fit a mapping function directly from image distribution to depth.

In this paper, we argue that *video depth estimator is not simply a video version of image depth estimator*. The core attribute of video depth estimation is *consistency*. The consistency for dynamic and static parts of the scene is essentially different and should be handled separately.

① Static regions involve only the camera motion, allowing the 3D structure to be analytically inferred from pairwise correspondences obtained through stereo match-

ing [18, 38, 50, 51, 63, 64, 68, 86] on a sequence of RGB frames, providing strong global 3D cues. The consistency of these areas, primarily about static backgrounds and across all video frames, is termed *global consistency*. Since static elements often occupy a large portion of the scene (e.g., roads, trees, buildings outdoors, or walls, tables, and floors indoors), a strong and robust global consistency is the foundation for achieving consistent and accurate video depth estimation. ② Dynamic parts contain both object motions and camera motion. It is infeasible to achieve analytical 4D reconstruction from RGB sequence alone, as it requires solving unknown object shapes, poses, and motion trajectories simultaneously, which is highly ill-posed. For example, imagine a scene where a person is waving his/her hand from left to right. The predicted depth maps are expected to not only strictly correspond to the RGB inputs in image composition, but also more importantly, maintain consistent, smooth depth changes for the moving hand across consecutive frames, without abrupt fluctuations or flickering. This temporal consistency across short sequences and particularly in dynamic areas, is termed *local consistency*, which should be learned by seeing large amount of video depth data.

Motivated by these analysis, we propose **StereoDiff**, a novel two-stage video depth estimator that synergizes both the stereo matching [30, 64, 86] for accurate global consistency and a video depth diffusion model [26, 53, 73] fine-tuned on large-scale video depth datasets for smooth local consistency. In the first stage of StereoDiff (Sec. 3.2), all video frames are processed in pairs through a stereo matching pipeline and then merged to establish strong global consistency. However, for dynamic objects, depth predictions are limited to pairwise frames (equivalent to a window size of 2), leading to clear inconsistencies (Fig. 1, middle column). Potential camera motion errors can also cause depth jitters across consecutive frames, resulting in suboptimal local consistency. To tackle this issue, in the second stage of StereoDiff (Sec. 3.3), a one-step video depth diffusion process is employed, in order to greatly improve the local consistency of stereo matching-based depth maps while preserving their original strong global consistency, resulting in video depth maps with both high-quality global and local consistency. Leveraging the priors of pre-trained video diffusion models, e.g., SVD, and fine-tuning them with extensive video depth data, video depth diffusion models achieve exceptionally smooth local consistency across neighboring frames. However, it is typically impossible for video diffusion-based video depth estimators to process all video frames simultaneously, which inherently limits their global consistency, as illustrated in the second column of Fig. 1.

We validate StereoDiff on two well-acknowledged, zero-shot, dynamic, and real-world video depth benchmarks (Tab. 1): Bonn [39] for indoor scenes and KITTI [19] for outdoor scenes. We also report the performance on static and

dynamic regions (Tab. 3); and the performance on different frequency domains (Tab. 2) to assess on global and local consistency, respectively. The results demonstrate that StereoDiff achieves SoTA performance on both dynamic video depth benchmarks. Furthermore, StereoDiff effectively retains the strong global consistency established in the first stage while significantly enhancing the local consistency in the second. Additionally, as shown in Tab. 5, thanks to the one-step policy in the second stage, StereoDiff is ~ 2.1 times faster than DepthCrafter.

In summary, our key contributions are as follows:

- We emphasize that achieving consistent video depth estimation requires distinct treatment for static (background) and dynamic (foreground) regions. Specifically, global consistency is better achieved through stereo matching on static regions, while local consistency for dynamic objects should be learned from large-scale video depth data.
- Based on these insights, we introduce **StereoDiff**, a novel two-stage video depth estimator that synergizes stereo matching for strong global consistency and video depth diffusion for smooth local consistency, delivering reliable video depth estimations. StereoDiff is training-free and does not require test-time optimization.
- Experimental results on dynamic, zero-shot, real-world video depth benchmarks (Tab. 1), both indoor and outdoor, demonstrate StereoDiff’s SoTA performance. In addition, analysis across frequency domains (Tab. 2 and Fig. 3) and in dynamic and static regions (Tab. 3) further shows that StereoDiff effectively integrates the strengths of both stereo matching and video depth diffusion models.

2. Related Works

2.1. Image Depth Estimation

Monocular image depth estimation has advanced significantly from early CNN-based approaches [14, 16, 29, 44, 79, 83] to vision transformer-based [13, 45, 80, 84]. To build powerful and generalizable depth estimators, DepthAnything [74, 75] and Metric3D [25, 81] series leveraged extensive training data comprising millions of samples, achieving SoTA performance. Additionally, some methods [1, 5, 40], *e.g.*, DepthPro [5] focus on accurately estimating the metric depth. Recent SD-based depth predictor, *e.g.*, Marigold [27] and GeoWizard [17] incorporated pre-trained diffusion priors for monocular depth estimation, achieved remarkable zero-shot generalizability. More recent studies [21, 36, 72], *e.g.*, GenPercept [72], Lotus [21], have further shown that single-step diffusion delivers even superior performance.

¹For clearer visualization, we filtered out low-confidence 3D points from the full point cloud, like those representing the moving yellow balloon.

2.2. Video Depth Estimation

SfM for Video Depth. Traditional Structure-from-Motion (SfM) methods [50, 51, 54, 61, 69, 90] can estimate only static 3D structure and camera positions, as dynamic objects violate triangulation constraints. Neither can those real-time visual SLAM systems [15, 48, 49, 57, 60], *e.g.*, NeuralRecon [57] and DoubleTake [49]. Earlier approaches [18, 38] adapted SfM for motions with strong assumptions, *e.g.*, rigidity. Recently, self-supervised methods [2, 3, 8, 11, 28, 33, 34, 58, 82, 85, 88] have tackled this via jointly estimating of video depth, camera poses, and motion residuals, *e.g.*, GeoNet [82], CasualSAM [88], and Robust-CVD [28, 34]. However, these methods require resource-intensive test-time optimization (or fine-tuning). More recent advancements, *e.g.*, DUS3R [64], MAST3R [30], and MonST3R [86], deliver more accurate and robust SfM results given monocular videos in an inference-based manner, even with large motions [86]. All video frames are pairwise processed and then merged, which brings global consistency. Nonetheless, due to their pairwise input mechanism, jitters and flickering between consecutive frames still persist, particularly on dynamic objects.

End-to-end Video Depth Estimators. The performance of traditional end-to-end methods [31, 59, 62, 65, 67, 76, 77, 85], *e.g.*, DeepV2D [59], NVDS [67], and FutureDepth [77], are inevitable constrained due to limited training data and model capacity. Recently, benefiting from web-scale image datasets [52], diffusion models [10, 20, 22, 37, 42, 43, 46, 47, 55, 56, 87] have achieved exceptional image generation capability, leading to significant progress in video generation [4, 7, 9, 23, 24, 66, 71, 89], *e.g.*, SVD [4] and Sora [7]. More recently, following the advancements of image depth estimation [17, 21, 27, 72], fine-tuning pre-trained video diffusion models using large-scale video depth data has gained traction [26, 53, 73], *e.g.*, ChronoDepth [53] and DepthCrafter [26], producing exceptionally smooth video depth predictions. However, input videos are typically divided into windows (of continuous or interpolated frames) and processed sequentially, which can lead to cross-window consistencies due to the absence of global 3D constraints.

Motivated by these methods, StereoDiff synergizes the strengths of both SfM and end-to-end video depth diffusion models, aiming to deliver video depth estimations with both strong global consistency and smooth local consistency.

3. Method

Given a monocular video with a sequence of RGB images $\mathcal{I} = \{I_t\}_{t=0}^{T-1}$, the goal of StereoDiff is to predict consistent depth maps across all video frames. As shown in Fig. 2, StereoDiff is a two-stage video depth estimator designed to achieve both global and local consistency. In the first stage, stereo matching [30, 64, 86] is applied

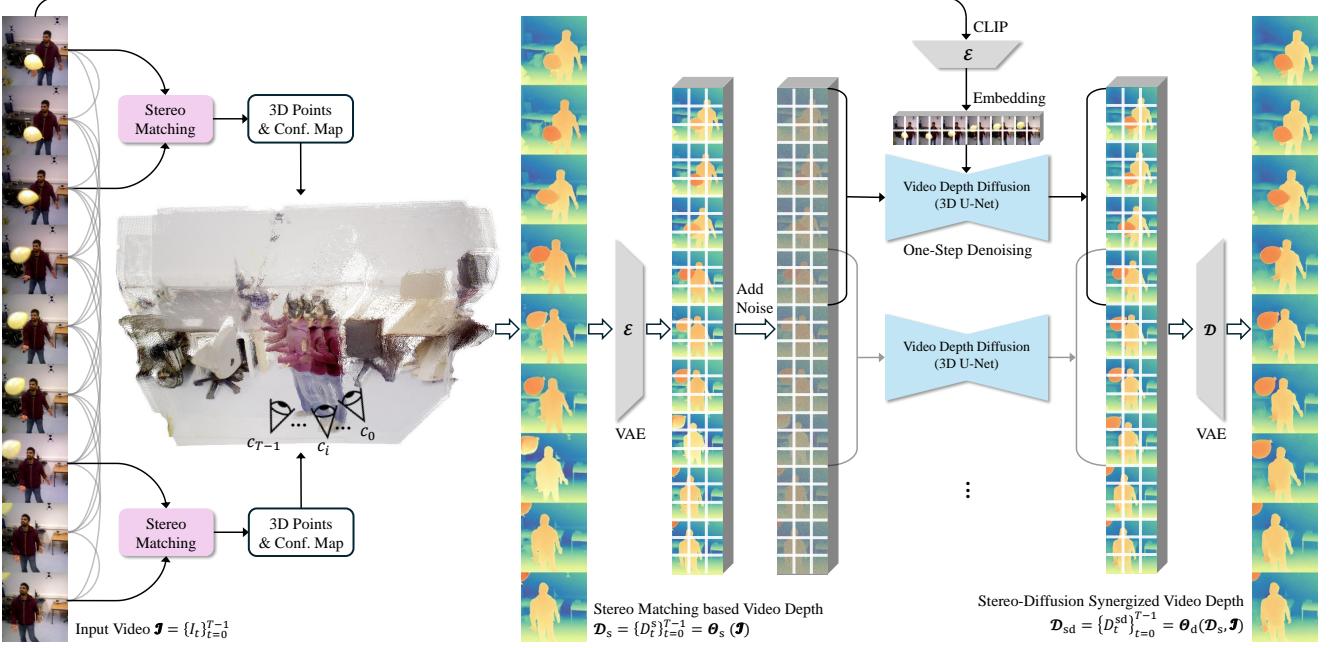


Figure 2. **Pipeline of StereoDiff.** ① All video frames are paired for stereo matching in the first stage, primarily focusing on static backgrounds, in order to achieve a strong global consistency¹. ② Using the stereo matching-based video depth from the first stage, the second stage of StereoDiff applies a video depth diffusion for significantly improving the local consistency without sacrificing its original global consistency, resulting in video depth estimations with both strong global consistency and smooth local consistency.

across all frames to establish strong global consistency, *i.e.*, $\mathcal{D}_s = \{D_t^s\}_{t=0}^{T-1} = \Theta_s(\mathcal{I})$. In the second stage, we use a video depth diffusion model [26, 53, 73] to enhance local consistency, particularly for dynamic objects, while preserving the global coherence achieved in the first stage, *i.e.*, $\mathcal{D}_{sd} = \{D_t^{sd}\}_{t=0}^{T-1} = \Theta_d(\mathcal{D}_s, \mathcal{I})$. This two-stage approach enables StereoDiff to deliver high-quality video depth that maintain coherence across both static and dynamic regions throughout the video. In Sec. 3.1, we formalize global and local consistency from the perspective of frequency domain analysis. Subsequently, Sec. 3.2 and Sec. 3.3 provide detailed descriptions of each stage.

3.1. Formulation of Consistency

Given a video depth estimation $\hat{\mathcal{D}} = \{\hat{D}_t\}_{t=0}^{T-1}$ and the corresponding GT depth \mathcal{D}^* , along with a metric function $f_\epsilon(\cdot)$ to measure the errors between them, we can calculate the sequence of error values:

$$\mathcal{E} = \{\epsilon_t\}_{t=0}^{T-1} = f_\epsilon(\mathcal{D}^*, \hat{\mathcal{D}}) \quad (1)$$

This error sequence can be represented as a sum of orthogonal waves with different frequencies. In this paper, we use fast Fourier transform (FFT) to compute the Discrete Fourier Transform (DFT) of error sequence \mathcal{E} , decomposing it into several frequency components:

$$\mathcal{F}(\epsilon_k) = \sum_{t=0}^{T-1} \epsilon_t \cdot e^{-i2\pi \frac{k}{T} t}, \quad k = 0, 1, \dots, T-1 \quad (2)$$

where $\mathcal{F}(\epsilon_k)$ represents the frequency component at the k -th frequency domain; T is the total number of frames; and i is the imaginary unit. The error sequence can further be reconstructed by Inverse DFT:

$$\epsilon_t = \frac{1}{T} \sum_{k=0}^{T-1} \mathcal{F}(\epsilon_k) \cdot e^{i2\pi \frac{k}{T} t}, \quad t = 0, 1, \dots, T-1 \quad (3)$$

Applying FFT to the error sequence, we can efficiently compute $\mathcal{F}(\epsilon_k)$ for all k frequency domains. This decomposition allows us to analyze the contribution of different frequency bands to the overall error, distinguishing between low-frequency and high-frequency components.

Global consistency refers to the overall stability of depth predictions across the entire video, especially in static backgrounds. For static or minimally dynamic objects, depth changes over time are primarily due to camera motion. Most real-world videos typically have a frame rate much higher than 1 FPS ($\ll 1\text{Hz}$), causing these depth variations to exhibit very low-frequency characteristics, sometimes appearing nearly linear. Global inconsistency often refers to persistent, significant depth deviations that remain stable over long sequences of consecutive frames, which strongly affects the low-frequency components of error sequence \mathcal{E} .

Local consistency focuses on stability between neighboring frames, particularly in dynamic areas with significant motion. Depth variations in these regions are influenced by both camera motion and object motion. Local inconsistencies can

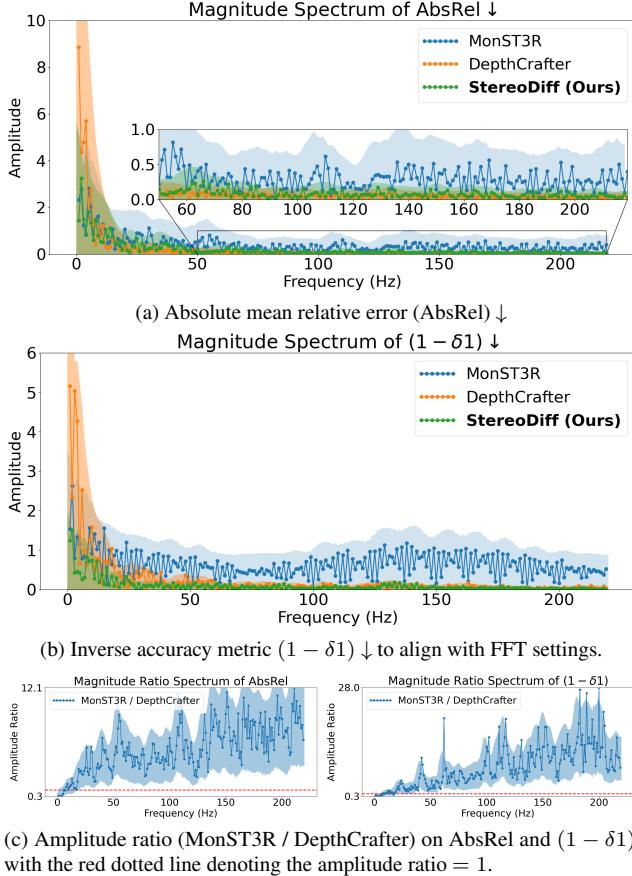


Figure 3. Magnitude spectrum of the error sequence on Bonn [39] dataset. The first scene of Bonn, “balloon”, containing 438 frames, is used as an example here. Due to symmetry, only the second half of the frequency spectrum is shown. Please refer to the supplementary material for additional visualizations of other error metrics, including RMSE ↓ and $(1 - \delta_2) \downarrow$, as well as the error sequence across different frequency domains for a more comprehensive and in-depth understanding.

arise from: 1) errors in camera motion estimation (common in stereo matching-based methods), causing sudden shifts and depth fluctuations in certain frames; and 2) limited window size, which inevitably prevents consistent and accurate depth tracking of moving objects, resulting in jitters and flickering. Although these local inconsistencies may not be clearly reflected on the overall metrics due to the limited number of affected frames, they can significantly increase the high-frequency amplitudes of the error sequence \mathcal{E} .

3.2. Stereo Matching for Global Consistency

Given the input RGB frames \mathcal{I} , the first stage of StereoDiff pairs each frame with the subsequent n frames, forming a total of $nT - (n + 1)n/2$ image pairs. Each pair is then processed through a stereo matching pipeline, resulting in coarse 3D point clouds that ensure the strong global consistency in video depth estimation.

Thanks to the advances of SfM [18, 38, 50, 51, 54, 63, 64, 68, 69, 86], we are fortunate to have works like DUS3R [64], MAST3R [30], and MonST3R [86] that offer highly accurate and robust stereo matching correspondences even without per-scene optimization. In this work, we adopt MonST3R [86] as the stereo matching pipeline, which fine-tunes DUS3R [64] with extensive dynamic video data. Compared to DUS3R, MonST3R more accurately assigns zero confidence to potential low-quality correspondences (*e.g.*, dynamic, blurry) and applies SfM only to static, clear correspondences, significantly enhancing the performance and robustness in dynamic scenes. Typically, an optimization-based post-processing step is applied for improved global alignment after obtaining stereo matching results. However, we exclude this step for three reasons: 1) video depth estimation is a perception task, which is better to be inference-based; 2) the optimization step is both resource-intensive² and time-consuming³; and 3) Similar to DUS3R [64] and MAST3R [30], MonST3R [86] inherently maintains global consistency through its closed-form global point cloud initialization, which uses a Minimum Spanning Tree (MST) to find the optimal path in the pairwise stereo matching graph with maximum confidence, followed by rigid point cloud registration [6, 35] to construct the final coarse 3D point clouds. As a result, StereoDiff is not only training-free but also fully inference-based⁴.

We denote the depth maps estimated only based on stereo matching as $\mathcal{D}_s = \{D_{t+j=0}^s\}_{t=0}^{T-1} = \Theta_s(\mathcal{I})$ and those only generated by video depth diffusion as $\mathcal{D}_d = \{D_t^d\}_{t=0}^{T-1} = \Theta_d(x \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathcal{I})$. As illustrated in Fig. 3, the magnitude spectrum two error sequences measured using AbsRel and $(1 - \delta_1)$ (please see Sec. 4.1.3 for specific definitions) are visualized. It is evident that \mathcal{D}_s exhibits significantly lower low-frequency errors compared to \mathcal{D}_d , indicating strong global consistency. Conversely, \mathcal{D}_d performs much better in high-frequency domains, which primarily represent the local consistency. These findings demonstrate the promising potential of leveraging the priors from video depth diffusion models to greatly enhance the local consistency of \mathcal{D}_s while maintaining its original high-quality global consistency.

3.3. Video Depth Diffusion for Local Consistency

Formally, taking \mathcal{D}_s as input, the video depth diffusion model produces the final video depth prediction, expressed as: $\mathcal{D}_{sd} = \{D_t^{sd}\}_{t=0}^{T-1} = \Theta_d(\mathcal{D}_s, \mathcal{I})$. In this paper, we adopt

²It requires > 80GB of graphics memory for videos with ≥ 300 frames at a resolution of 512×384 , making it impractical for long videos.

³Processing a 200-frame video at 512×384 resolution with a 300-iteration optimization takes over 15 minutes on an NVIDIA A800 GPU.

⁴We omit the Weiszfeld algorithm [41] for focal length estimation, as it requires only 10 iterations and back-propagates gradients into a minimal $T \times 1$ matrix, where T is the number of frames.

⁵Comparisons are conducted in disparity space rather than true-depth space, because both DepthCrafter and StereoDiff represent their video depth estimations using disparity maps.

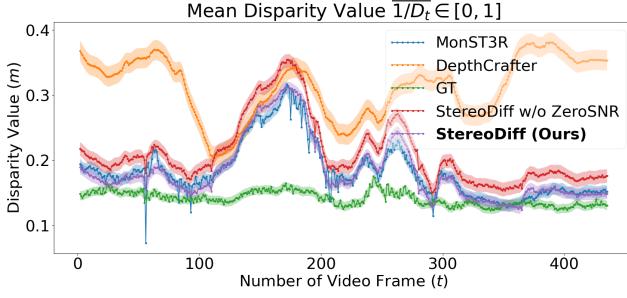


Figure 4. **Comparison of mean disparity⁵ value $\overline{1/D_t}$ tested on Bonn [39] dataset** for MonST3R [86], DepthCrafter [26], and StereoDiff. All disparity maps are normalized to $[0, 1]$ on a per-scene basis before comparison. Incorporating ZeroSNR drags the mean value of StereoDiff’s disparity maps closer to the GT, resulting in improved performance (Tab. 4).

DepthCrafter [26], a fine-tuned SVD model using $\sim 20K$ video sequences, to perform a one-step denoising of \mathcal{D}_s . Unlike SfM-based video depth estimation, which adheres to the “first principle”, video depth diffusion models take a purely “data-driven” approach. These models are fine-tuned from pre-trained video generative models on large-scale video depth data, mapping the RGB video directly to video depth.

As shown in Fig. 3, the depth maps produced by video depth diffusion models \mathcal{D}_d significantly outperform those based on stereo matching \mathcal{D}_s in high-frequency domains. Particularly, Fig. 3c depicts the amplitude ratio of the error sequences calculated on \mathcal{D}_s and \mathcal{D}_d for clearer demonstration. This suggests that the components in higher frequency domains of \mathcal{D}_s , which much more significantly differ from the GT distribution learned by the video depth diffusion models, are more likely treated as noise and effectively denoised. Conversely, the low-frequency characteristics of \mathcal{D}_s align much more closely with the GT video depth distribution, drawing less attention during denoising and thus being better preserved. This results in strong retention of low-frequency features and targeted denoising of high-frequency components, significantly reducing the high-frequency errors in \mathcal{D}_s .

Mathematically, substituting \mathcal{D}_s into Eq. 1 yields the corresponding error sequence $\mathcal{E}_s = \{\epsilon_t^s\}_{t=0}^{T-1}$. This temporal signal can then be transformed into the frequency domain $\mathcal{F}(\epsilon_k^s)$, $k \in [0, T - 1]$ using FFT (Eq. 2). Similarly, we denote the error sequence of \mathcal{D}_{sd} as $\mathcal{E}_{sd} = \{\epsilon_t^{sd}\}_{t=0}^{T-1}$. $\forall t \in [0, T - 1]$, $\epsilon_t^s \geq 0$ and $\epsilon_t^{sd} \geq 0$. The average of error sequence yields the final metric: $(1/T) \sum_{t=0}^{T-1} \epsilon_t$. As discussed above and demonstrated in Fig. 3, during the second stage of StereoDiff, the video depth diffusion model acts as a “low-pass filter” on $\mathcal{F}(\epsilon_k^s)$. Assuming a threshold K_{thr} , for simplicity, we approximate that after the video depth diffusion process, the magnitudes of all frequency components $> K_{thr}$ are re-scaled by a factor $\alpha \in (0, 1)$:

$$\mathcal{F}(\epsilon_k^{sd}) \approx \begin{cases} \mathcal{F}(\epsilon_k^s), & k \leq K_{thr} \\ \alpha \cdot \mathcal{F}(\epsilon_k^s), & k > K_{thr} \end{cases} \quad (4)$$

Following Parseval’s energy theorem, which states that the total energy of the signal in the time domain and frequency domain are equal, we can derive:

$$\begin{aligned} \sum_{k=0}^{T-1} |\mathcal{F}(\epsilon_k^{sd})|^2 &\leq \sum_{k=0}^{T-1} |\mathcal{F}(\epsilon_k^s)|^2 \\ \Rightarrow \sum_{t=0}^{T-1} |\epsilon_t^{sd}|^2 &\leq \sum_{t=0}^{T-1} |\epsilon_t^s|^2 \Rightarrow \frac{1}{T} \sum_{t=0}^{T-1} \epsilon_t^{sd} \leq \frac{1}{T} \sum_{t=0}^{T-1} \epsilon_t^s \end{aligned} \quad (5)$$

This derivation shows that maintaining the low-frequency characteristics of \mathcal{D}_s , while reducing the high-frequency components of its error sequence \mathcal{E}_s , leads to improved performance. In practice, as illustrated in Fig. 3, StereoDiff’s low-frequency error magnitudes \mathcal{D}_{sd} largely inherit those of \mathcal{D}_s , while high-frequency components are significantly reduced by leveraging the video depth diffusion, leading to improved performance (Tab. 1 and 2) and greatly smoothed prediction (Fig. 1), aligning well with our analysis.

ZeroSNR. In diffusion models, the forward process progressively adds Gaussian noise to clean samples according to a pre-defined variance schedule, *i.e.*, β_1, \dots, β_T :

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t \mathbf{I}) \quad (6)$$

Let $\alpha_t = 1 - \beta_t$ and $\bar{a}_t = \prod_{s=1}^t \alpha_s$, x_t can be sampled as:

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{a}_t} x_0, (1 - \bar{a}_t) \mathbf{I}) \quad (7)$$

Equivalently:

$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (8)$$

The SNR is defined as: $\text{SNR}(t) = \bar{a}_t / (1 - \bar{a}_t)$. Specifically, in DepthCrafter [26] and the standard SVD [4] scheduler⁶, the variance sequence is $\beta_0 = 0.00085$ and $\beta_T = 0.012$ with linear scaling, we derive: $x_T \approx 0.0016x_0 + 0.9992\epsilon$. This indicates the input, *i.e.*, x_T , always contains a small amount of signal during training. The leaked signal contains the lowest frequency information, *e.g.*, the mean value. The model learns to denoise with this signal. However, during inference, pure Gaussian noise is used, prompting the model to generate outputs with medium value [32, 36].

As illustrated in Fig. 4, DepthCrafter’s video disparity maps have a mean value closer to 0.5 compared to other methods. Although StereoDiff achieves relatively accurate mean disparity values without ZeroSNR due to its first stage (stereo matching), incorporating ZeroSNR further aligns the mean value of StereoDiff’s disparity maps closer to GT, resulting in improved performance (Tab. 4).

⁶<https://huggingface.co/docs/diffusers/en/api/schedulers/euler>

Method	Bonn [39] (Indoor)				KITTI [19] (Outdoor)				Average	Average
	AbsRel ↓	RMSE ↓	$\delta 1 \uparrow$	$\delta 2 \uparrow$	AbsRel ↓	RMSE ↓	$\delta 1 \uparrow$	$\delta 2 \uparrow$	Rank ⁷ ↓	Rank ⁸ ↓
DepthAnything V2 [75]	0.1250	1.7765	0.8297	0.9912	0.1758	<u>4.2583</u>	0.6872	<u>0.9664</u>	6.25	4.75
DepthAnything [74]	0.1112	1.5191	0.8860	0.9956	0.1755	<u>4.3756</u>	0.6875	0.9678	4.75	3.25
DUSt3R [64]	0.1757	2.3618	0.7798	0.9602	0.3343	7.0966	0.5065	0.7497	7.75	7.75
MAS3R [30]	0.1748	2.2829	0.7698	0.9125	0.2250	5.0800	0.6460	0.8664	6.75	7.13
MonST3R [86]	<u>0.0818</u>	<u>1.2412</u>	<u>0.9542</u>	0.9916	0.1661	4.1881	0.7387	0.9662	2.75	<u>2.38</u>
ChronoDepth [53]	0.1248	1.6918	0.8501	0.9823	0.1749	4.4265	0.7288	0.9334	4.25	5.00
DepthCrafter [26]	0.1104	1.6817	0.8955	<u>0.9945</u>	<u>0.1617</u>	5.3883	<u>0.7695</u>	0.9518	<u>2.50</u>	3.50
StereoDiff (Ours)	0.0799	1.2257	0.9549	0.9870	0.1469	4.4183	0.7764	0.9654	1.00	2.25

Table 1. **Quantitative comparison of StereoDiff with SoTA methods on zero-shot, real-world, dynamic video depth benchmarks.** The four sections from top to bottom represent: image depth estimators, stereo matching-based estimators, video depth diffusion models, and StereoDiff. To ensure comprehensive evaluation, we used two datasets: Bonn [39] for indoor scenes and KITTI [19] for outdoor scenes. We report the mean metric value of StereoDiff across 10 independent runs. Best results are **bolded** and the second best are underlined.

4. Experiments

4.1. Experimental Settings

4.1.1. Implementation Details

In the first stage, we set $n = 2$ for forming image pairs, symmetrizing them before feeding them into the stereo matching pipeline. The Weiszfeld algorithm [41] is adopted for camera intrinsics, and Procrustes alignment [35] is used for solving camera poses. The maximum resolution is limited to 512. In the second stage, following [26], we set the window size to 110 frames with a 25-frame overlap. The ZeroSNR trick is implemented by setting the trailing [32, 36] mode for the timestep spacing in schedulers. Depth maps obtained from the first stage \mathcal{D}_s are resized to the original frame size using nearest interpolation before the one-step denoising process, which is performed from denoising timestep $t = 2$ to $t = 1$ with a total number of denoising timesteps $T = 4$.

4.1.2. Evaluation Datasets.

We use two zero-shot, real-world, dynamic benchmark datasets for evaluation: Bonn [39] (indoor) and KITTI [19] (outdoor). Six dynamic indoor videos from Bonn, with $332 \sim 580$ frames each, and twelve dynamic outdoor videos from KITTI’s validation set, with $17 \sim 251$ frames each, are included. Video frames from Bonn are sized 640×480 and 1216×352 in KITTI. Note that we omit static video depth benchmark datasets such as ScanNet [12] and ScanNet++ [78], as they can be trivially handled by methods like DUSt3R [64], MAS3R [30], and MonST3R [86].

About Zero-Shot. Zero-shot evaluation is more challenging and more close to the practical unknown scenarios during application. Note that for KITTI [19] dataset, it is one of the training datasets used by DepthCrafter [26], whereas DUSt3R [64] and MonST3R [86] were not trained on KITTI.

⁷Rankings on AbsRel and $\delta 1$, the two most recognized metrics in the realm of depth estimation, following [17, 25–27, 53, 72–75, 81, 86].

⁸Rankings on AbsRel, RMSE, $\delta 1$, and $\delta 2$.

For Bonn [39] dataset, it is zero-shot for DepthCrafter [26], DUSt3R [64] and MonST3R [86].

About Real-World. StereoDiff is primarily validated on real-world datasets for several reasons: ① Synthetic datasets are inherently artificial. Dense, high-quality GTs that align perfectly with designed scenarios are inherently accessible. ② In real-world settings, true GTs are inaccessible, and sensor-based approximations often face challenges like missing depth values. Also, real-world situations better reflects the scenarios where the video depth estimators will be potentially applied, *e.g.*, robotics.

4.1.3. Evaluation Metrics.

Following the affine-invariant evaluation protocols from [21, 26, 27, 53, 73, 86], we firstly align the estimated video depth maps with GT using least-squares fitting, and resize all estimations to match the original size of input video in nearest mode. Note that during the least-squares fitting, all frames in a video depth sequence share identical scaling and shifting factors. Specifically, given GT $\mathcal{D}^* = \{D_t^*\}_{t=0}^{T-1}$ and fitted predictions $\hat{\mathcal{D}} = \{\hat{D}_t\}_{t=0}^{T-1}$, we report two error metrics: 1) absolute mean relative error (AbsRel) and 2) root-mean-square deviation (RMSE), *i.e.*:

$$\text{AbsRel}(\mathcal{D}^*, \hat{\mathcal{D}}) = \frac{1}{T} \sum_{t=0}^{T-1} \left[\frac{1}{N} \sum_{j=0}^{N-1} \frac{|D_{tj}^* - \hat{D}_{tj}|}{\hat{D}_{tj}} \right] \quad (9)$$

$$\text{RMSE}(\mathcal{D}^*, \hat{\mathcal{D}}) = \frac{1}{T} \sum_{t=0}^{T-1} \left[\frac{1}{N} \sqrt{\sum_{j=0}^{N-1} (D_{tj}^* - \hat{D}_{tj})^2} \right]$$

where $N = H \times W$, indicating the total number of pixels. We also report two accuracy metrics: $\delta 1$ and $\delta 2$, denoting the proportion of pixels satisfying $\text{Max}(D_{tj}^*/\hat{D}_{tj}, \hat{D}_{tj}/D_{tj}^*) < 1.25$ and 1.25^2 , respectively.

Metrics	Method	Low Freq. $\xleftarrow{\quad}$ High Freq. $\xrightarrow{\quad}$										
		\mathcal{F}_0	\mathcal{F}_1	\mathcal{F}_2	\mathcal{F}_3	\mathcal{F}_4	\mathcal{F}_5	\mathcal{F}_6	\mathcal{F}_7	\mathcal{F}_8	\mathcal{F}_9	\mathcal{F}_{10}
AbsRel \downarrow	DepthCrafter	0.1104	0.0152	0.0215	0.0238	0.0286	0.0206	0.0112	0.0062	0.0023	0.0012	0.0009
	MonST3R	0.0822	0.0130	<u>0.0149</u>	<u>0.0142</u>	0.0149	0.0142	0.0144	0.0116	0.0077	0.0062	0.0067
	StereoDiff (Ours)	0.0806	0.0159	0.0128	0.0132	0.0157	0.0143	0.0135	0.0098	0.0067	0.0043	0.0032
$(1 - \delta_1) \downarrow$	DepthCrafter	0.1046	0.0380	0.0655	0.0696	0.0835	0.0619	0.0331	0.0198	0.0100	0.0046	0.0027
	MonST3R	<u>0.0481</u>	0.0207	<u>0.0247</u>	0.0313	0.0408	0.0411	<u>0.0335</u>	<u>0.0258</u>	0.0180	0.0134	0.0150
	StereoDiff (Ours)	0.0478	0.0246	0.0241	<u>0.0325</u>	0.0428	0.0442	0.0371	0.0261	<u>0.0173</u>	0.0101	0.0069

Table 2. **Quantitative comparisons of MonST3R, DepthCrafter, and StereoDiff on different frequency domains.** We use DFT and Inverse DFT to disentangle the components of the metric sequences calculated on Bonn [39] dataset in various frequency domains. For simplicity, the entire frequency range is divided into 11 discrete groups: $\mathcal{F}_0, \dots, \mathcal{F}_{10}$, representing low to high frequencies. We report the results on two most recognized metrics, AbsRel \downarrow and $(1 - \delta_1) \downarrow$. The results on other metrics are provided in the supplementary material.

Region	AbsRel \downarrow	RMSE \downarrow	$\delta_1 \uparrow$	$\delta_2 \uparrow$
Dynamic	-0.0069	-0.0844	+0.0140	+0.0023
Overall	<u>-0.0020</u>	<u>-0.0150</u>	<u>+0.0013</u>	<u>-0.0042</u>
Static	+0.0009	0	-0.0004	-0.0049

(a) Performance improvement of StereoDiff over MonST3R. For example, $\text{AbsRel} = \text{AbsRel}_{\text{StereoDiff}} - \text{AbsRel}_{\text{MonST3R}}$.

Region	AbsRel \downarrow	RMSE \downarrow	$\delta_1 \uparrow$	$\delta_2 \uparrow$
Dynamic	-0.0178	-0.2575	+0.0413	-0.0055
Overall	<u>-0.0306</u>	<u>-0.4555</u>	<u>+0.0600</u>	<u>-0.0071</u>
Static	-0.0335	-0.4990	+0.0641	-0.0069

(b) Performance improvement of StereoDiff over DepthCrafter.

Table 3. **Quantitative comparisons on dynamic and static regions of the scene** among MonST3R, DepthCrafter and StereoDiff. We use FlowSAM [70] for masking moving objects.

Method	AbsRel \downarrow	RMSE \downarrow	$\delta_1 \uparrow$	$\delta_2 \uparrow$
Naive Solution				
w/ Latent Sharing	0.1245	1.7807	0.8503	0.9719
w/o ZeroSNR	± 0.0002	± 0.0016	± 0.0018	± 0.0006
w/o Latent Sharing	<u>0.0809</u>	<u>1.2383</u>	<u>0.9544</u>	<u>0.9867</u>
w/o ZeroSNR	± 0.0003	± 0.0039	± 0.0006	± 0.0003
StereoDiff (Ours)	0.0799	1.2257	0.9549	0.9870
w/ Latent Sharing	± 0.0001	± 0.0028	± 0.0006	± 0.0004
w/ ZeroSNR				

Table 4. **Ablation studies.** Removing latent sharing strategy and adding the ZeroSNR trick both yield effective performance gains. Here we report the results on Bonn dataset [39].

Method	DepthCrafter	MonST3R	StereoDiff (Ours)
Inf. Time (s)	1.1708	0.4100	0.4100+0.1569

Table 5. **Inference time per frame** tested on the first scene of Bonn [39] dataset (“balloon”), using an NVIDIA A800 GPU. We set $n = 2$ for both MonST3R and StereoDiff.

4.2. Quantitative & Qualitative Comparisons

As shown in Tab. 1, StereoDiff achieves SoTA performance across both dynamic, zero-shot, real-world video depth

benchmarks. Furthermore, the results of frequency domain analysis (Tab. 2) demonstrate that StereoDiff effectively maintains the strong low-frequency global consistency achieved via stereo matching, while significantly enhancing the high-frequency local consistency. This enhancement greatly reduces local jitters and flickering across neighboring frames particularly in dynamic areas (Fig. 1), as high-frequency characteristics of \mathcal{D}_s differ much more significantly from the GT distribution learned by the video depth diffusion models, and are more likely treated as noise and effectively denoised. Additionally, Tab. 3 clearly shows that StereoDiff outperforms MonST3R mainly in high-frequency dynamic regions and outperforms DepthCrafter mainly in low-frequency static regions. These results align well with our analysis in Sec. 3.2 and 3.3. Part of the qualitative comparisons on Bonn [39] dataset are shown in Fig. 1. Please refer to the supplementary material for the qualitative comparisons on KITTI [19] dataset and further results on Bonn [39].

Inference Speed. The inference time comparison among MonST3R [86], DepthCrafter [26] and StereoDiff is reported in Tab. 5. Thanks to efficient stereo matching and MST alignment, especially the one-step denoising policy of the video depth diffusion model in the second stage, StereoDiff is ~ 2.1 times faster than DepthCrafter.

4.3. Ablation Study

As discussed in Sec. 2.2, for video diffusion-based video depth estimators, input videos are typically divided into windows and processed sequentially. In DepthCrafter, this is performed by dividing the video into overlapped windows and sharing the latents of overlapped frames. While this strategy improves continuity, it can still fall short in maintaining consistency between windows, especially on static backgrounds (Fig. 1). As illustrated in Tab. 4, the removal of latent sharing strategy leads to significant performance gains. This is primarily because: 1) the strict spatial correspondence between the diffusion’s latent space and the RGB space, making latent sharing ineffective for scenes with moving cameras or objects, which may lead to harmful feature distortions, especially as the timestep $t \rightarrow 0$; and 2)

in DepthCrafter’s original multi-step denoising process, the latent is progressively refined from Gaussian noise, where sharing latents across overlapping frames can not only aids consistency at early timesteps ($t \rightarrow T$) but also allows the distortions of latent feature to be gradually refined as $t \rightarrow 0$. Additionally, incorporating ZeroSNR aligns the mean value of StereoDiff’s disparity maps more closely with the GT (Fig. 4), further enhancing the performance.

5. Conclusion

In this paper, we emphasize the need for distinct strategies to achieve consistent video depth estimation across static and dynamic regions. Motivated by these insights, we introduce StereoDiff, a novel two-stage video depth estimator that combines stereo matching for strong global consistency provided by the global 3D constraints, and video depth diffusion for significantly enhanced local consistency. Experimental results on two well-acknowledged video depth benchmarks (Tab. 1), including the frequency domain analysis (Tab. 2 and Fig. 3), demonstrate StereoDiff’s effectiveness in synergizing the strengths of both, achieving SoTA performance in dynamic, zero-shot, real-world video depth estimation.

6. Acknowledgments

We sincerely thank Tingyang Zhang and Peng-Shuai Wang from Peking University, Jiahao Shao from Zhejiang University, and Vlas Zyrianov from the University of Illinois Urbana-Champaign for their insightful discussions.

References

- [1] Shariq Farooq Bhat, Reiner Birk, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 3
- [2] Jiawang Bian, Zhichao Li, Naiyan Wang, Huangying Zhan, Chunhua Shen, Ming-Ming Cheng, and Ian Reid. Unsupervised scale-consistent depth and ego-motion learning from monocular video. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 3
- [3] Jia-Wang Bian, Huangying Zhan, Naiyan Wang, Tat-Jin Chin, Chunhua Shen, and Ian Reid. Auto-rectify network for unsupervised indoor depth estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2021. 3
- [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. 2, 3, 6
- [5] Aleksei Bochkovskii, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R Richter, and Vladlen Koltun. Depth pro: Sharp monocular metric depth in less than a second. *arXiv preprint arXiv:2410.02073*, 2024. 3
- [6] Romain Brégier. Deep regression on manifolds: a 3D rotation case study. In *2021 International Conference on 3D Vision (3DV)*, 2021. 5
- [7] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>, 3, 2024. 3
- [8] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Depth prediction without the sensors: Leveraging structure for unsupervised learning from monocular videos. In *Proceedings of the AAAI conference on artificial intelligence*, pages 8001–8008, 2019. 3
- [9] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023. 3
- [10] 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. 3
- [11] Yuhua Chen, Cordelia Schmid, and Cristian Sminchisescu. Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7063–7072, 2019. 3
- [12] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 7
- [13] Ainaz Eftekhari, Alexander Sax, Jitendra Malik, and Amir Zamir. Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10786–10796, 2021. 3
- [14] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014. 3
- [15] Jakob Engel, Vladlen Koltun, and Daniel Cremers. Direct sparse odometry. *IEEE transactions on pattern analysis and machine intelligence*, 40(3):611–625, 2017. 3
- [16] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2002–2011, 2018. 3
- [17] Xiao Fu, Wei Yin, Mu Hu, Kaixuan Wang, Yuexin Ma, Ping Tan, Shaojie Shen, Dahua Lin, and Xiaoxiao Long. Geowizard: Unleashing the diffusion priors for 3d geometry estimation from a single image. In *European Conference on Computer Vision*, pages 241–258. Springer, 2025. 2, 3, 7
- [18] Ravi Garg, Anastasios Roussos, and Lourdes Agapito. Dense variational reconstruction of non-rigid surfaces from monocular video. In *Proceedings of the IEEE Conference on computer*

- vision and pattern recognition*, pages 1272–1279, 2013. 2, 3, 5, 23
- [19] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013. 2, 7, 8, 13, 14, 17, 18, 21, 22, 23, 24, 25
- [20] Jing He, Haodong Li, Yongzhe Hu, Guibao Shen, Yingjie Cai, Weichao Qiu, and Ying-Cong Chen. Disenvisoner: Disentangled and enriched visual prompt for customized image generation. *arXiv preprint arXiv:2410.02067*, 2024. 3
- [21] Jing He, Haodong Li, Wei Yin, Yixun Liang, Leheng Li, Kaiqiang Zhou, Hongbo Liu, Bingbing Liu, and Ying-Cong Chen. Lotus: Diffusion-based visual foundation model for high-quality dense prediction. *arXiv preprint arXiv:2409.18124*, 2024. 2, 3, 7
- [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3
- [23] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 3
- [24] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022. 3
- [25] Mu Hu, Wei Yin, Chi Zhang, Zhipeng Cai, Xiaoxiao Long, Hao Chen, Kaixuan Wang, Gang Yu, Chunhua Shen, and Shaojie Shen. Metric3d v2: A versatile monocular geometric foundation model for zero-shot metric depth and surface normal estimation. *arXiv preprint arXiv:2404.15506*, 2024. 3, 7
- [26] Wenbo Hu, Xiangjun Gao, Xiaoyu Li, Sijie Zhao, Xiaodong Cun, Yong Zhang, Long Quan, and Ying Shan. Depthcrafter: Generating consistent long depth sequences for open-world videos. *arXiv preprint arXiv:2409.02095*, 2024. 1, 2, 3, 4, 6, 7, 8, 14, 20, 22, 23
- [27] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing diffusion-based image generators for monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9492–9502, 2024. 2, 3, 7
- [28] Johannes Kopf, Xuejian Rong, and Jia-Bin Huang. Robust consistent video depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1611–1621, 2021. 3
- [29] Jin Han Lee, Myung-Kyu Han, Dong Wook Ko, and Il Hong Suh. From big to small: Multi-scale local planar guidance for monocular depth estimation. *arXiv preprint arXiv:1907.10326*, 2019. 3
- [30] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. *arXiv preprint arXiv:2406.09756*, 2024. 2, 3, 5, 7
- [31] Zhaoshuo Li, Wei Ye, Dilin Wang, Francis X Creighton, Russell H Taylor, Ganesh Venkatesh, and Mathias Unberath. Temporally consistent online depth estimation in dynamic scenes. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 3018–3027, 2023. 3
- [32] Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and sample steps are flawed. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 5404–5411, 2024. 6, 7
- [33] Chao Liu, Jinwei Gu, Kihwan Kim, Srinivasa G Narasimhan, and Jan Kautz. Neural rgb (r) d sensing: Depth and uncertainty from a video camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10986–10995, 2019. 3
- [34] Xuan Luo, Jia-Bin Huang, Richard Szeliski, Kevin Matzen, and Johannes Kopf. Consistent video depth estimation. *ACM Transactions on Graphics (ToG)*, 39(4):71–1, 2020. 3
- [35] Priyanka Mandikal, KL Navaneet, Mayank Agarwal, and R Venkatesh Babu. 3d-lmnet: Latent embedding matching for accurate and diverse 3d point cloud reconstruction from a single image. *arXiv preprint arXiv:1807.07796*, 2018. 5, 7
- [36] Gonzalo Martin Garcia, Karim Abou Zeid, Christian Schmidt, Daan de Geus, Alexander Hermans, and Bastian Leibe. Fine-tuning image-conditional diffusion models is easier than you think. *arXiv preprint arXiv:2409.11355*, 2024. 2, 3, 6, 7
- [37] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 3
- [38] Kemal E Ozden, Konrad Schindler, and Luc Van Gool. Multi-body structure-from-motion in practice. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(6):1134–1141, 2010. 2, 3, 5, 23
- [39] Emanuele Palazzolo, Jens Behley, Philipp Lottes, Philippe Giguere, and Cyrill Stachniss. Refusion: 3d reconstruction in dynamic environments for rgb-d cameras exploiting residuals. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7855–7862. IEEE, 2019. 2, 5, 6, 7, 8, 13, 14, 16, 19, 20, 23
- [40] Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. Unidepth: Universal monocular metric depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10106–10116, 2024. 3
- [41] Frank Plastria. The weiszfeld algorithm: proof, amendments, and extensions. *Foundations of location analysis*, pages 357–389, 2011. 5, 7
- [42] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 3
- [43] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 3
- [44] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular

- depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020. 3
- [45] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–12188, 2021. 3
- [46] 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, 3
- [47] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 3
- [48] Mohamed Sayed, John Gibson, Jamie Watson, Victor Prisacariu, Michael Firman, and Clément Godard. Simplerecon: 3d reconstruction without 3d convolutions. In *European Conference on Computer Vision*, pages 1–19. Springer, 2022. 3
- [49] Mohamed Sayed, Filippo Aleotti, Jamie Watson, Zawar Qureshi, Guillermo Garcia-Hernando, Gabriel Brostow, Sara Vicente, and Michael Firman. Doubletake: Geometry guided depth estimation. *arXiv preprint arXiv:2406.18387*, 2024. 3
- [50] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016. 2, 3, 5, 23
- [51] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision*, 2016. 2, 3, 5, 23
- [52] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 3
- [53] Jiahao Shao, Yuanbo Yang, Hongyu Zhou, Youmin Zhang, Yujun Shen, Matteo Poggi, and Yiyi Liao. Learning temporally consistent video depth from video diffusion priors. *arXiv preprint arXiv:2406.01493*, 2024. 2, 3, 4, 7
- [54] Noah Snavely, Steven M Seitz, and Richard Szeliski. Photo tourism: exploring photo collections in 3d. In *ACM siggraph 2006 papers*, pages 835–846. 2006. 3, 5, 23
- [55] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. 3
- [56] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3
- [57] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. Neuralrecon: Real-time coherent 3d reconstruction from monocular video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15598–15607, 2021. 3
- [58] Libo Sun, Jia-Wang Bian, Huangying Zhan, Wei Yin, Ian Reid, and Chunhua Shen. Sc-depthv3: Robust self-supervised monocular depth estimation for dynamic scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2023. 3
- [59] Zachary Teed and Jia Deng. Deepv2d: Video to depth with differentiable structure from motion. *arXiv preprint arXiv:1812.04605*, 2018. 3
- [60] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. *Advances in neural information processing systems*, 34:16558–16569, 2021. 3
- [61] George Vogiatzis and Carlos Hernández. Video-based, real-time multi-view stereo. *Image and vision computing*, 29(7):434–441, 2011. 3
- [62] Chaoyang Wang, Simon Lucey, Federico Perazzi, and Oliver Wang. Web stereo video supervision for depth prediction from dynamic scenes. In *2019 International Conference on 3D Vision (3DV)*, pages 348–357. IEEE, 2019. 3
- [63] Hengyi Wang and Lourdes Agapito. 3d reconstruction with spatial memory. *arXiv preprint arXiv:2408.16061*, 2024. 2, 5, 23
- [64] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 2, 3, 5, 7, 23
- [65] Yiran Wang, Zhiyu Pan, Xingyi Li, Zhiguo Cao, Ke Xian, and Jianming Zhang. Less is more: Consistent video depth estimation with masked frames modeling. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 6347–6358, 2022. 3
- [66] Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *arXiv preprint arXiv:2309.15103*, 2023. 3
- [67] Yiran Wang, Min Shi, Jiaqi Li, Zihao Huang, Zhiguo Cao, Jianming Zhang, Ke Xian, and Guosheng Lin. Neural video depth stabilizer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9466–9476, 2023. 2, 3
- [68] Philippe Weinzaepfel, Vincent Leroy, Thomas Lucas, Romain Brégier, Yohann Cabon, Vaibhav Arora, Leonid Antsfeld, Boris Chidlovskii, Gabriela Csurka, and Jérôme Revaud. Croco: Self-supervised pre-training for 3d vision tasks by cross-view completion. *Advances in Neural Information Processing Systems*, 35:3502–3516, 2022. 2, 5, 23
- [69] Changchang Wu. Towards linear-time incremental structure from motion. In *2013 International Conference on 3D Vision-3DV 2013*, pages 127–134. IEEE, 2013. 3, 5, 23
- [70] Junyu Xie, Charig Yang, Weidi Xie, and Andrew Zisserman. Moving object segmentation: All you need is sam (and flow). *arXiv preprint arXiv:2404.12389*, 2024. 8, 14

- [71] 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, 2025. 3
- [72] Guangkai Xu, Yongtao Ge, Mingyu Liu, Chengxiang Fan, Kangyang Xie, Zhiyue Zhao, Hao Chen, and Chunhua Shen. What matters when repurposing diffusion models for general dense perception tasks? *arXiv preprint arXiv:2403.06090*, 2024. 2, 3, 7
- [73] Honghui Yang, Di Huang, Wei Yin, Chunhua Shen, Haifeng Liu, Xiaofei He, Binbin Lin, Wanli Ouyang, and Tong He. Depth any video with scalable synthetic data. *arXiv preprint arXiv:2410.10815*, 2024. 2, 3, 4, 7
- [74] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10371–10381, 2024. 3, 7
- [75] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv preprint arXiv:2406.09414*, 2024. 3, 7
- [76] Rajeev Yasarla, Hong Cai, Jisoo Jeong, Yunxiao Shi, Risheek Garrepalli, and Fatih Porikli. Mamo: Leveraging memory and attention for monocular video depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8754–8764, 2023. 2, 3
- [77] Rajeev Yasarla, Manish Kumar Singh, Hong Cai, Yunxiao Shi, Jisoo Jeong, Yinhao Zhu, Shizhong Han, Risheek Garrepalli, and Fatih Porikli. Futuredepth: Learning to predict the future improves video depth estimation. In *European Conference on Computer Vision*, pages 440–458. Springer, 2025. 2, 3
- [78] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12–22, 2023. 7
- [79] Wei Yin, Xinlong Wang, Chunhua Shen, Yifan Liu, Zhi Tian, Songcen Xu, Changming Sun, and Dou Renyin. Diversedepth: Affine-invariant depth prediction using diverse data. *arXiv preprint arXiv:2002.00569*, 2020. 3
- [80] Wei Yin, Jianming Zhang, Oliver Wang, Simon Niklaus, Long Mai, Simon Chen, and Chunhua Shen. Learning to recover 3d scene shape from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 204–213, 2021. 3
- [81] Wei Yin, Chi Zhang, Hao Chen, Zhipeng Cai, Gang Yu, Kaixuan Wang, Xiaozihi Chen, and Chunhua Shen. Metric3d: Towards zero-shot metric 3d prediction from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9043–9053, 2023. 3, 7
- [82] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1983–1992, 2018. 3
- [83] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. Neural window fully-connected crfs for monocular depth estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3916–3925, 2022. 3
- [84] Chi Zhang, Wei Yin, Billzb Wang, Gang Yu, Bin Fu, and Chunhua Shen. Hierarchical normalization for robust monocular depth estimation. *Advances in Neural Information Processing Systems*, 35:14128–14139, 2022. 3
- [85] Haokui Zhang, Chunhua Shen, Ying Li, Yuanzhouhan Cao, Yu Liu, and Youliang Yan. Exploiting temporal consistency for real-time video depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1725–1734, 2019. 3
- [86] Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing Sun, and Ming-Hsuan Yang. Monst3r: A simple approach for estimating geometry in the presence of motion. *arXiv preprint arXiv:2410.03825*, 2024. 1, 2, 3, 5, 6, 7, 8, 14, 15, 20, 22, 23
- [87] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 3
- [88] Zhoutong Zhang, Forrester Cole, Zhengqi Li, Michael Rubinstein, Noah Snavely, and William T Freeman. Structure and motion from casual videos. In *European Conference on Computer Vision*, pages 20–37. Springer, 2022. 3
- [89] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022. 3
- [90] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1851–1858, 2017. 3

StereoDiff: Stereo-Diffusion Synergy for Video Depth Estimation

Supplementary Material

7. More Results of Frequency Domain Analysis

Additional magnitude spectrum visualizations of error sequences from the Bonn [39] dataset, including $\text{RMSE}\downarrow$ and $(1 - \delta_2)\downarrow$, are presented in Fig. 5. Similarly, the visualizations of magnitude spectrum on $\text{AbsRel}\downarrow$, $\text{RMSE}\downarrow$, $(1 - \delta_1)\downarrow$, and $(1 - \delta_2)\downarrow$ of KITTI [19] dataset are illustrated in Fig. 6. It is evident that the depth maps generated by video depth diffusion models \mathcal{D}_d significantly outperform those based on stereo matching \mathcal{D}_s in high-frequency domains, reflecting smoother local consistency. Conversely, \mathcal{D}_s performs better in low-frequency domains than \mathcal{D}_d , highlighting robust and strong global consistency. In practice, as illustrated in Fig. 5 and 6, StereoDiff's depth maps \mathcal{D}_{sd} effectively inherit the low-frequency error magnitudes of \mathcal{D}_s while significantly reducing the high-frequency errors by leveraging a one-step denoising process based on the video depth diffusion model. This process effectively achieves greatly smoothed local consistency, while retaining the original high-quality global consistency, aligning well with our analysis.

Furthermore, the supplementary results of frequency domain analysis on Bonn [39] and KITTI [19] dataset is re-

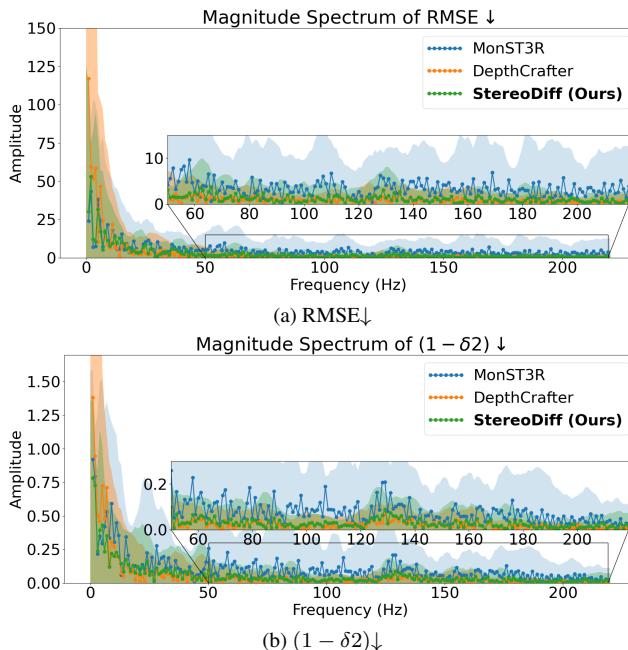


Figure 5. **Magnitude spectrum of the error sequence on the Bonn [39] dataset.** The first scene of Bonn, “balloon”, containing 438 frames, is used as an example here. Due to symmetry, only the second half of the frequency spectrum is shown.

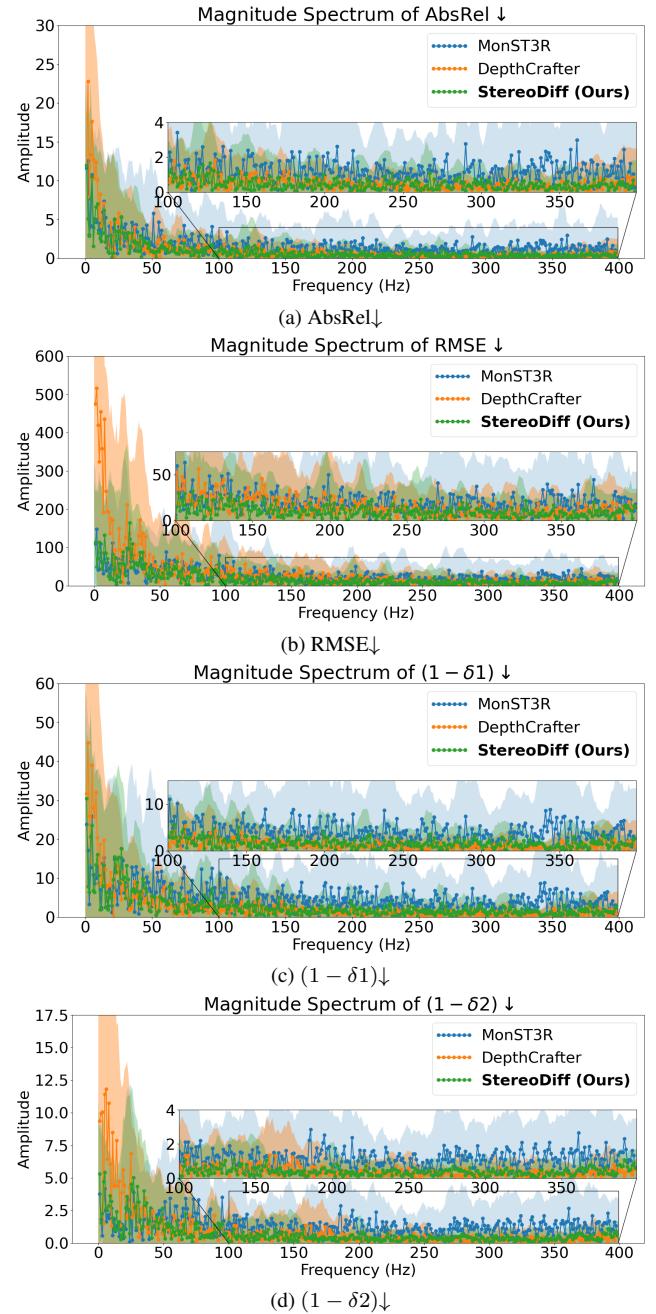


Figure 6. **Magnitude spectrum of the error sequence on the KITTI [19] dataset.** The spectrum is computed across all samples used for validation, as many KITTI outdoor videos contain limited number (< 100) of frames. Similar to Fig. 5, only the second half of the frequency spectrum is shown due to symmetry.

Metrics	Method	Low Freq. $\xleftarrow{\quad}$ High Freq. $\xrightarrow{\quad}$										
		\mathcal{F}_0	\mathcal{F}_1	\mathcal{F}_2	\mathcal{F}_3	\mathcal{F}_4	\mathcal{F}_5	\mathcal{F}_6	\mathcal{F}_7	\mathcal{F}_8	\mathcal{F}_9	\mathcal{F}_{10}
RMSE \downarrow	DepthCrafter	1.6823	0.1783	0.3221	0.2269	0.3125	0.2567	0.1448	0.0884	0.0355	0.0191	0.0144
	MonST3R	1.2427	0.0949	0.1075	0.1633	0.1503	0.1579	0.1604	0.1356	0.0848	0.0678	0.0726
	StereoDiff (Ours)	1.2294	0.1349	0.1065	0.1657	0.1659	0.1565	0.1469	0.1187	0.0786	0.0517	0.0421
$(1 - \delta^2) \downarrow$	DepthCrafter	0.0055	0.0008	0.0019	0.0028	0.0035	0.0030	0.0025	0.0019	0.0007	0.0004	0.0004
	MonST3R	0.0084	0.0013	0.0014	0.0054	0.0054	0.0058	0.0065	0.0057	0.0072	0.0041	0.0022
	StereoDiff (Ours)	0.0133	0.0064	0.0050	0.0136	0.0141	0.0174	0.0186	0.0156	0.0033	0.0023	0.0018

Table 6. **Quantitative comparisons of different frequency domains on Bonn [39] dataset**, among MonST3R [86], DepthCrafter [26], and StereoDiff. We use DFT and Inverse DFT to disentangle the components of the metric sequences into various frequency domains. For clearer visualization, the entire frequency range is divided into 11 discrete groups: $\mathcal{F}_0 \sim \mathcal{F}_{10}$, representing low to high frequencies. Here we only report supplementary results calculated on RMSE \downarrow and $(1 - \delta^2) \downarrow$.

Metrics	Method	Low Freq. $\xleftarrow{\quad}$ High Freq. $\xrightarrow{\quad}$								
		\mathcal{F}_0	\mathcal{F}_1	\mathcal{F}_2	\mathcal{F}_3	\mathcal{F}_4	\mathcal{F}_5	\mathcal{F}_6	\mathcal{F}_7	
AbsRel \downarrow	DepthCrafter	0.1620	0.0306	0.0324	0.0363	0.0272	0.0169	0.0129	0.0103	0.0076
	MonST3R	0.1666	0.0258	0.0221	0.0277	0.0279	0.0208	0.0190	0.0135	0.0135
	StereoDiff (Ours)	0.1476	0.0209	0.0155	0.0285	0.0247	0.0171	0.0136	0.0106	0.0078
RMSE \downarrow	DepthCrafter	5.4048	0.7941	0.8940	1.0056	0.8343	0.4651	0.3548	0.2641	0.1965
	MonST3R	4.1926	0.4247	0.3956	0.4656	0.5366	0.5599	0.5215	0.3529	0.2526
	StereoDiff (Ours)	4.4291	0.2985	0.3678	0.5270	0.5345	0.4628	0.3496	0.2690	0.2293
$(1 - \delta^1) \downarrow$	DepthCrafter	0.2322	0.0635	0.0671	0.0821	0.0674	0.0482	0.0352	0.0269	0.0204
	MonST3R	0.2647	0.0679	0.0506	0.0977	0.0853	0.0605	0.0555	0.0428	0.0427
	StereoDiff (Ours)	0.2304	0.0777	0.0557	0.0930	0.0744	0.0618	0.0403	0.0325	0.0262
$(1 - \delta^2) \downarrow$	DepthCrafter	0.0489	0.0155	0.0201	0.0279	0.0257	0.0178	0.0105	0.0072	0.0055
	MonST3R	0.0344	0.0100	0.0073	0.0159	0.0127	0.0111	0.0113	0.0100	0.0092
	StereoDiff (Ours)	0.0367	0.0136	0.0133	0.0216	0.0171	0.0147	0.0118	0.0086	0.0068

Table 7. **Quantitative comparisons across different frequency domains on KITTI [19] dataset**, among MonST3R [86], DepthCrafter [26], and StereoDiff. Following the settings in Tab. 6, we apply DFT and Inverse DFT to decompose the metric sequences into various frequency domains. For clearer visualization, the entire frequency range is grouped into 9 discrete bands, \mathcal{F}_0 to \mathcal{F}_8 , representing low to high frequencies. Results are reported on AbsRel \downarrow , RMSE \downarrow , $(1 - \delta^1) \downarrow$, and $(1 - \delta^2) \downarrow$.

Region	AbsRel \downarrow	RMSE \downarrow	$\delta 1 \uparrow$	$\delta 2 \uparrow$
Dynamic	-0.0463	-0.5809	+0.0982	+0.0294
Overall	<u>-0.0191</u>	<u>+0.2364</u>	<u>+0.0375</u>	<u>-0.0017</u>
Static	-0.0171	+0.2968	+0.0326	-0.0042

(a) Performance improvement of StereoDiff over MonST3R. For example, $\text{AbsRel} = \text{AbsRel}_{\text{StereoDiff}} - \text{AbsRel}_{\text{MonST3R}}$.

Region	AbsRel \downarrow	RMSE \downarrow	$\delta 1 \uparrow$	$\delta 2 \uparrow$
Dynamic	+0.0110	-0.4692	-0.0344	-0.0131
Overall	-0.0147	-0.9638	+0.0067	+0.0128
Static	-0.0184	-1.0070	+0.0126	+0.0163

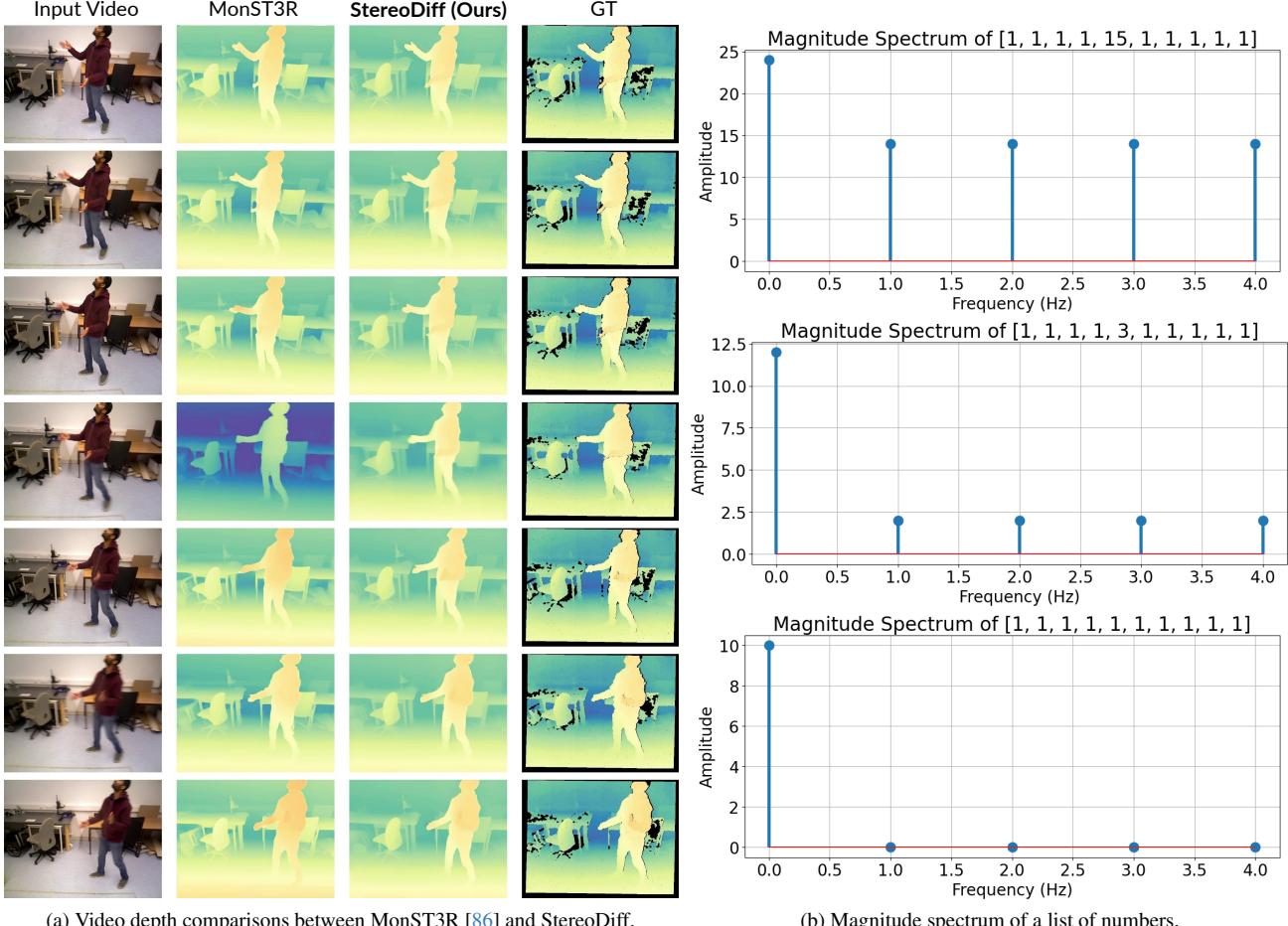
(b) Performance improvement of StereoDiff over DepthCrafter.

Table 8. **Quantitative comparisons on dynamic and static regions of the scenes in the KITTI [19] dataset**, among MonST3R [86], DepthCrafter [26] and StereoDiff. FlowSAM [70] is utilized for masking moving objects.

ported in Tab. 6 and 7. These results also demonstrate that StereoDiff effectively maintains the strong low-frequency global consistency achieved via stereo matching, while significantly enhancing the high-frequency local consistency.

For a more comprehensive and in-depth understanding, the error sequences calculated on AbsRel \downarrow , RMSE \downarrow , $(1 - \delta^1) \downarrow$, and $(1 - \delta^2) \downarrow$, across different frequency domains are visualized, as shown in Fig. 8 for Bonn [39] dataset and Fig. 9, 10 for KITTI [19] dataset. The frequency domains are grouped exponentially: the first range, \mathcal{F}_0 , includes the first 2^0 discrete frequencies, \mathcal{F}_1 covers the next 2^1 , \mathcal{F}_2 includes the next 2^2 , and so on. For clearer visualization, only the first scene of Bonn dataset is illustrated, which contains 438 frames and 219 discrete frequency domains, resulting in 8 groups. However, for quantitative comparisons (Tab. 6 and main paper's Tab. 2), all 6 scenes with a total of 2744 frames are used, resulting in 11 groups. For KITTI, all evaluation scenes are included in both the visualizations (Fig. 6, 9, and 10) and quantitative comparisons (Tab. 7), with 797 frames in total, resulting in 9 frequency domain groups.

Additionally, Tab. 8 clearly shows that StereoDiff outperforms MonST3R mainly in high-frequency dynamic regions and outperforms DepthCrafter mainly in low-frequency static regions. These results also align well with our analysis.



(a) Video depth comparisons between MonST3R [86] and StereoDiff.
(b) Magnitude spectrum of a list of numbers.

Figure 7. Smoothing abrupt depth jitters and flickering also improves the low-frequency performance. On the left (a) we display a sequence of video depth comparisons between MonST3R [86] and StereoDiff, demonstrating that StereoDiff effectively eliminate the harmful and sudden depth jitters caused by the significantly inaccurate camera pose estimations of certain frames. On the right (b), we present the magnitude spectrum of three different number lists, which are served as toy examples representing the error sequences of (from top to bottom): stereo matching-based video depth \mathcal{D}_s , StereoDiff’s video depth \mathcal{D}_{sd} that *partially* mitigates harmful abrupt depth shifts, and the ideal StereoDiff output that *completely* eliminates these abrupt depth shifts, respectively.

8. Analysis on Performance in Low Frequency

Generally, we expect StereoDiff’s depth maps \mathcal{D}_{sd} can basically inherit the global consistency obtained via the stereo matching stage, with significantly enhanced local consistency. However, as shown in Tab. 6 and 7, and main paper’s Tab. 1 and 2, the low-frequency error components of stereo matching-based depth maps \mathcal{D}_s are often slightly reduced, which may seem counterintuitive. This section provides supplementary analysis to explain this phenomenon.

As discussed in Sec. 1 and 3.2 of the main paper, local consistency, particularly in dynamic areas, relates to the temporal stability of depth values across short sequences and exhibits high-frequency characteristics. As shown in Fig. 7 (a), besides the inconsistencies on dynamic objects caused by limited window-size, stereo matching-based depth

maps \mathcal{D}_s can also suffer from abrupt errors in camera motion estimation, causing sudden shifts and depth fluctuations in certain frames, which is one of the reasons that bring local inconsistencies (Sec. 3.2 of the main paper). Additionally, Fig. 7 (b) displays the magnitude spectrum of three number lists representing the error sequences of stereo matching-based video depth \mathcal{D}_s and StereoDiff’s video depth \mathcal{D}_{sd} that *partially* and *completely* eliminates these abrupt shifts. The results reveal that while StereoDiff significantly reduces high-frequency errors in \mathcal{D}_s by approximately 6 times ($\sim 15.0 \rightarrow \sim 2.5$), the low-frequency errors are also slightly reduced ($\sim 25.0 \rightarrow \sim 12.5$, approximately 2 times smaller). These findings indicate that the second stage of StereoDiff, while primarily targeting high-frequency errors to enhance local consistency, can also work on low-frequency errors, slightly improving the global consistency.

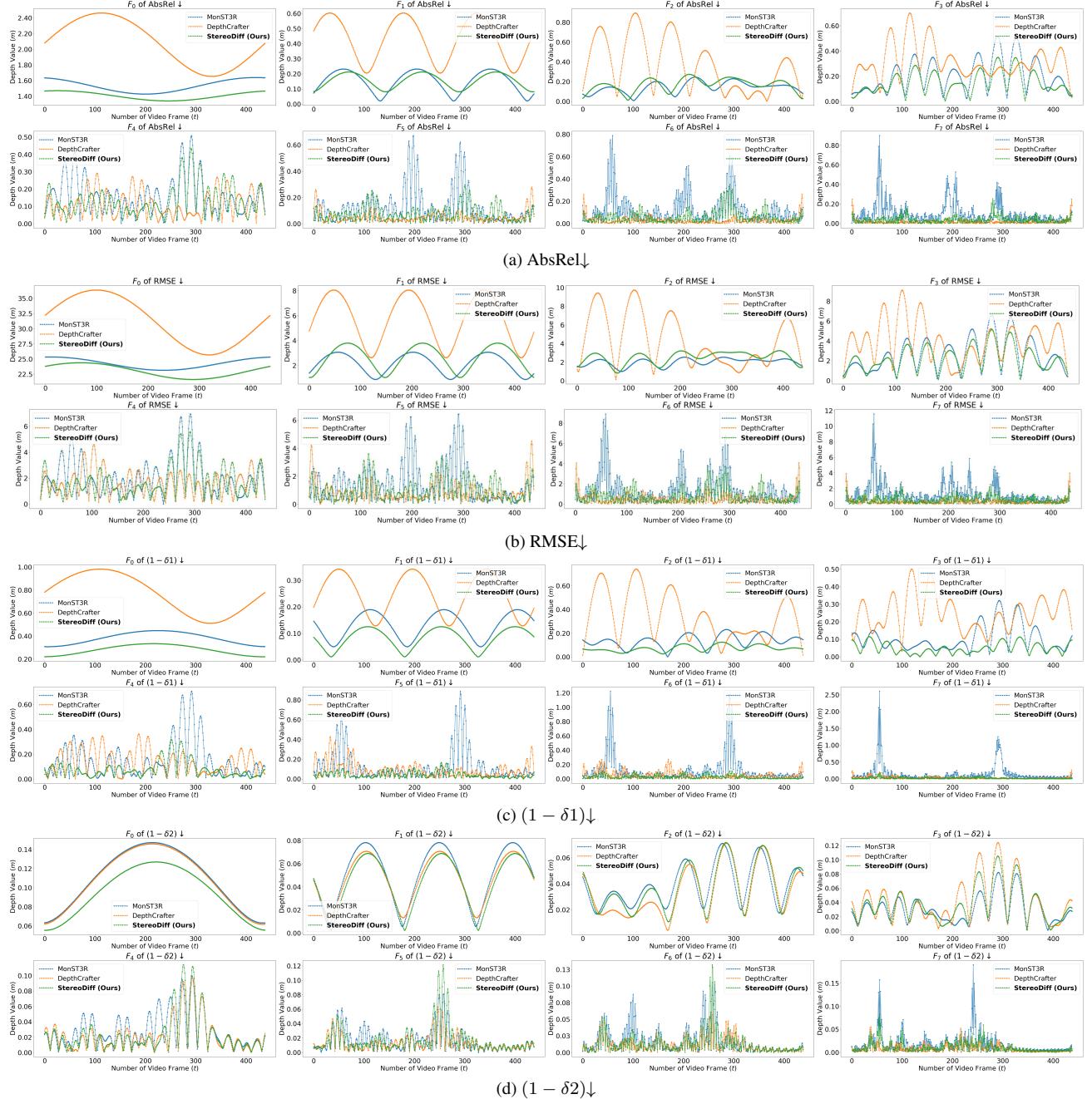


Figure 8. The error sequence across different frequency domains on Bonn [39] dataset. The first scene of Bonn, “balloon”, containing 438 frames, is used as an example here. DFT and Inverse DFT are utilized to disentangle the components of the metric sequences into various frequency domains. For clearer visualization, only the modulus of the obtained complex numbers after Inverse DFT is plotted.

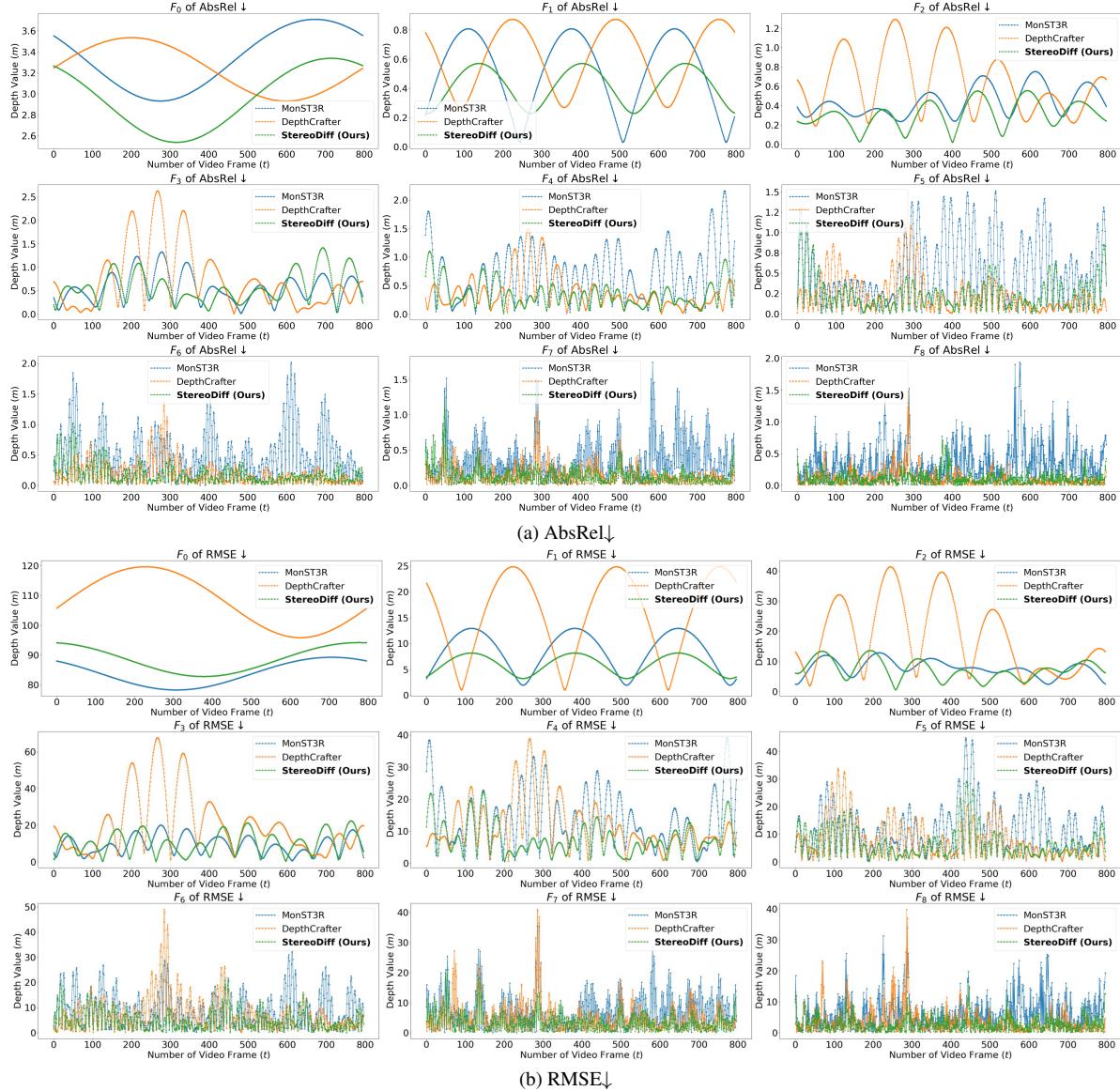


Figure 9. The error sequence of $\text{AbsRel}\downarrow$ and $\text{RMSE}\downarrow$ across different frequency domains on KITTI [19] dataset. Following Tab. 6, DFT and Inverse DFT are utilized to disentangle the components of the metric sequences into various frequency domains. Also for clearer visualization, only the modulus of the obtained complex numbers after Inverse DFT is plotted.



Figure 10. The error sequence of $(1 - \delta 1) \downarrow$ and $(1 - \delta 2) \downarrow$ across different frequency domains on KITTI [19] dataset. The settings for this visualization are exactly the same as Fig. 9.

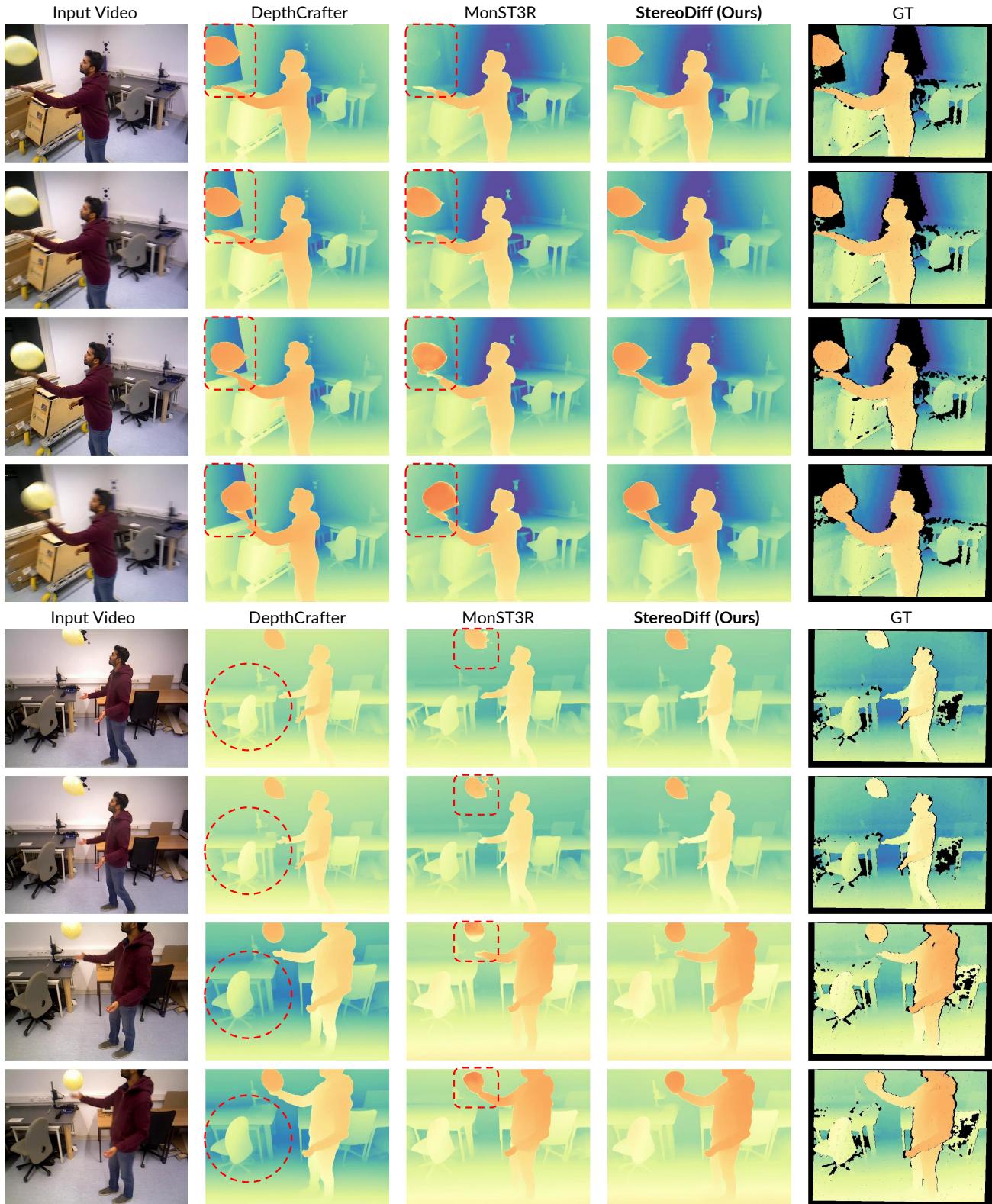


Figure 11. Qualitative Comparisons on Bonn [39] dataset. Please see the specific caption on the next page.

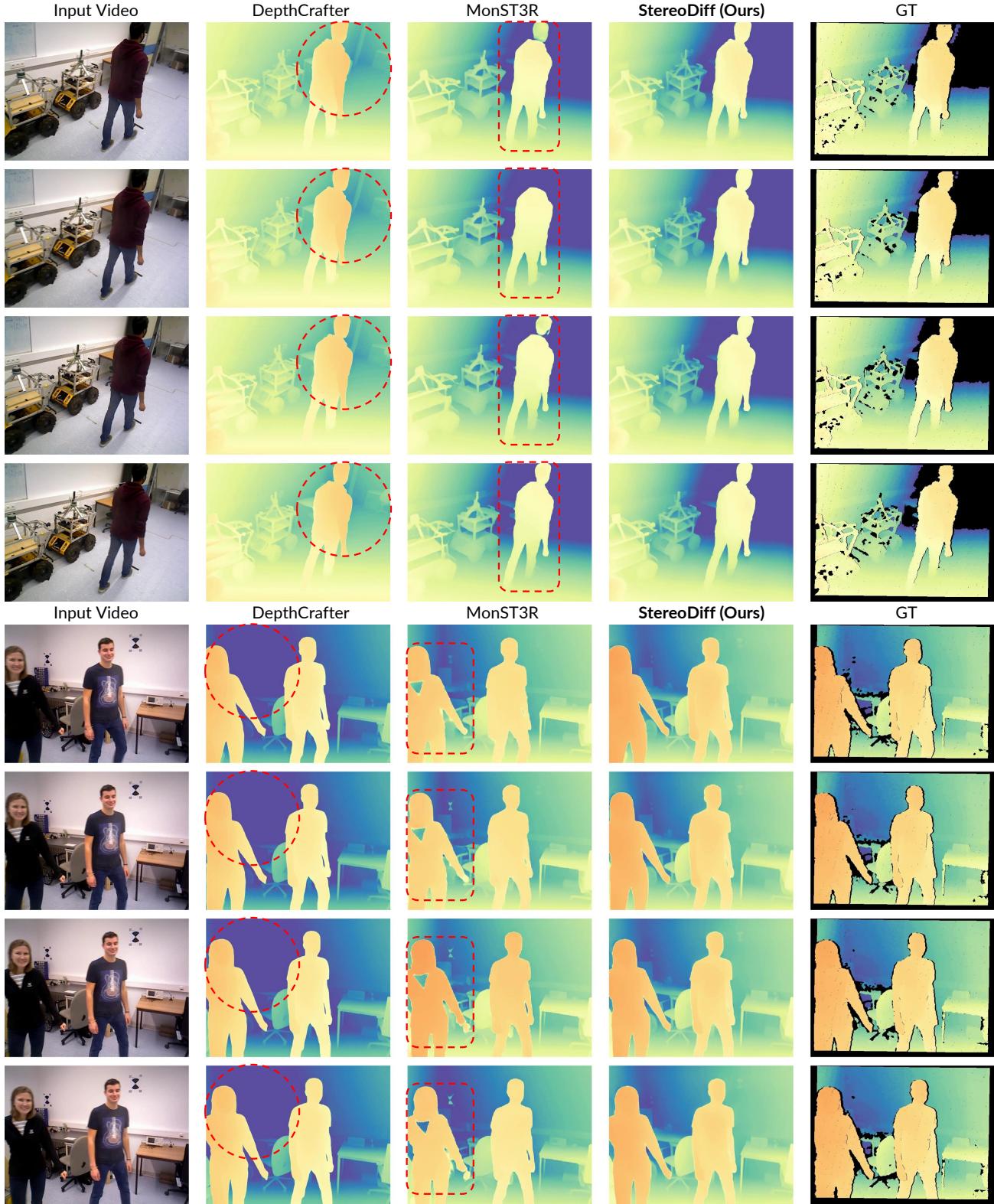


Figure 12. Qualitative Comparisons on Bonn [39] dataset, conducted among two SoTA video depth estimators, MonST3R [86] and DepthCrafter [26], alongside StereoDiff. Four frames are sampled from a video depth sequence to form a complete comparison set.

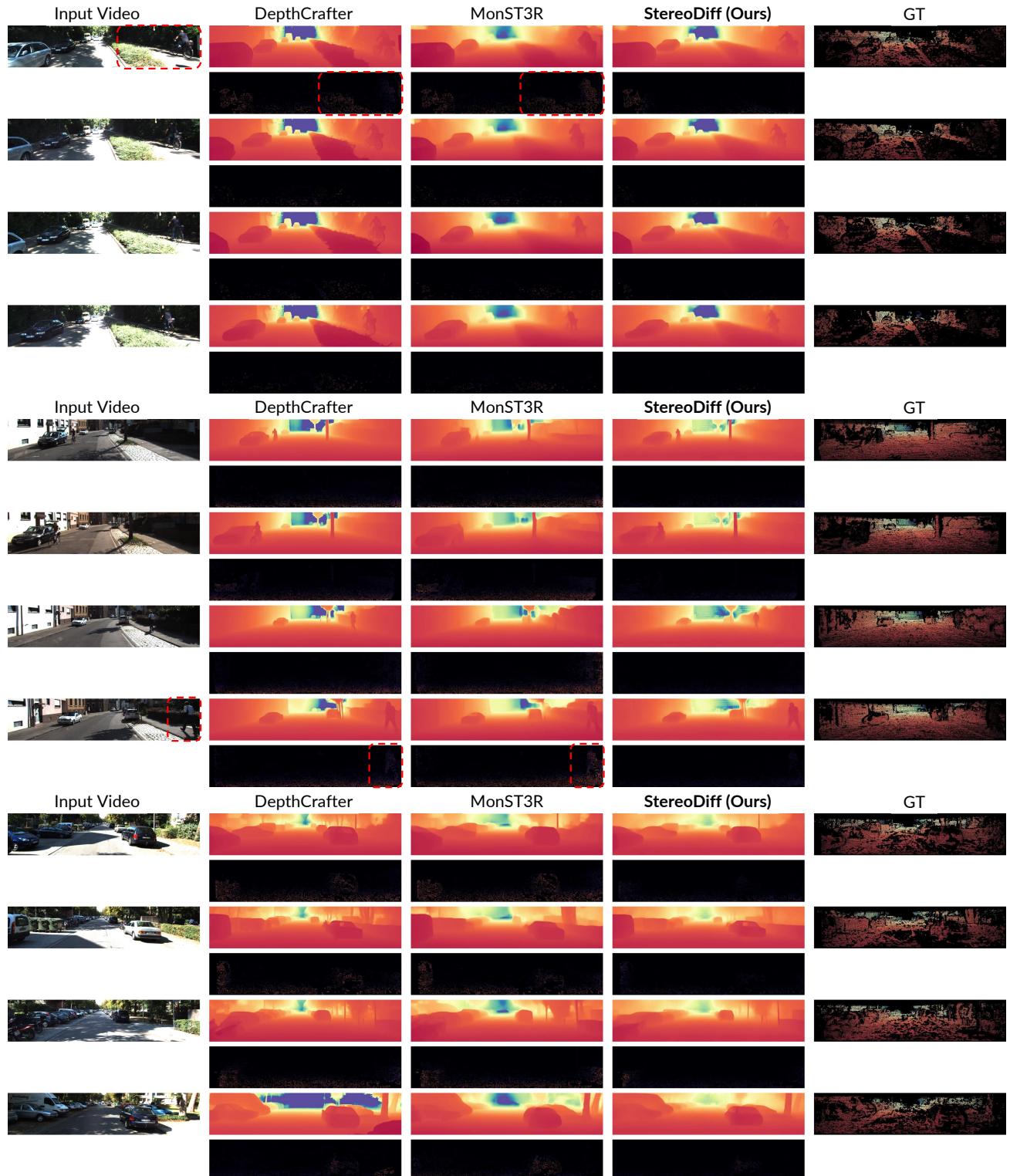


Figure 13. **Qualitative Comparisons on KITTI [19] dataset.** Please see the specific caption on the next page.

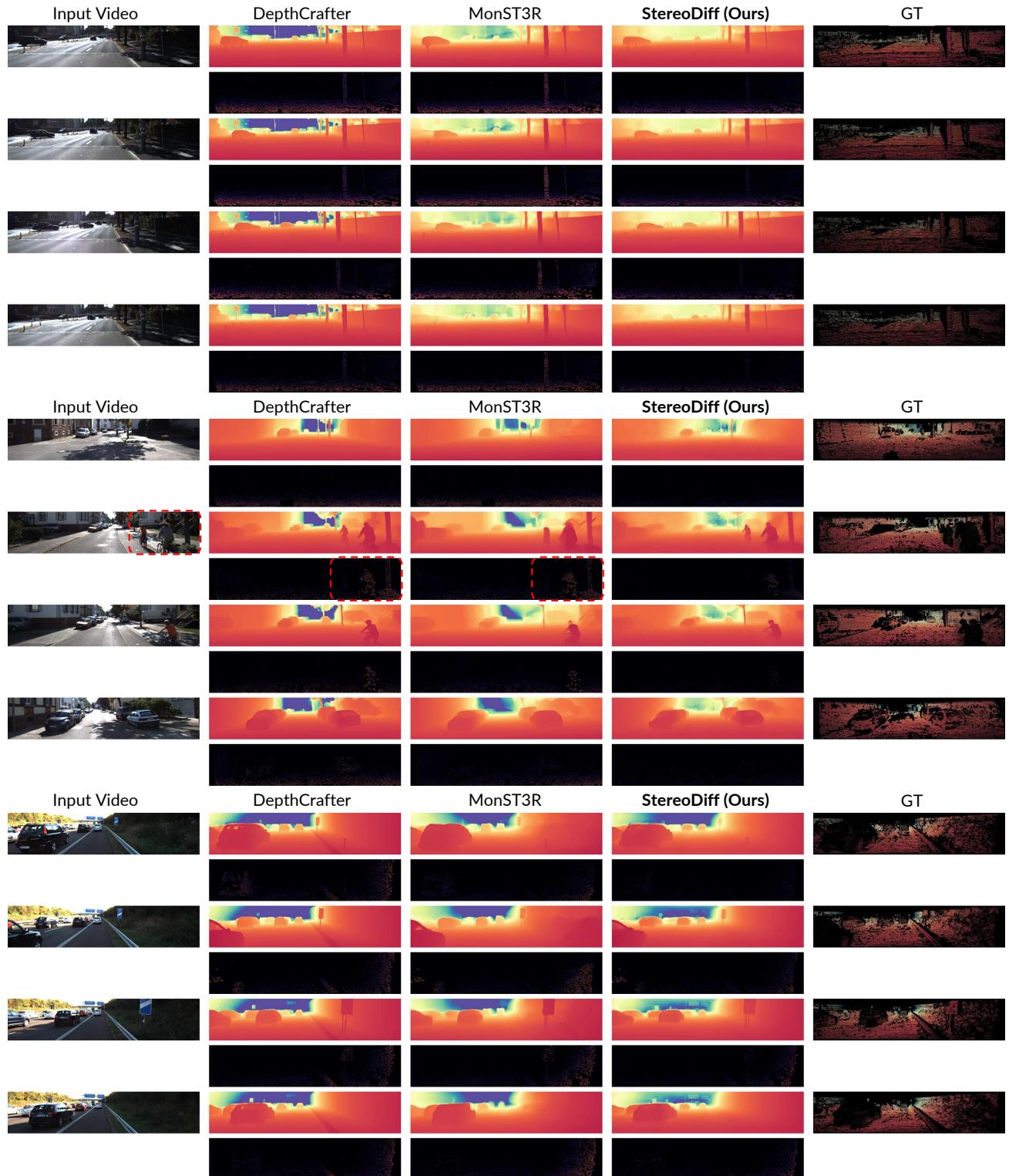


Figure 14. **Qualitative Comparisons on KITTI [19] dataset.** Comparisons among MonST3R [86], DepthCrafter [26], and StereoDiff are presented, with four frames sampled from a video depth sequence as a comparison set, similar to Fig. 8. For better clarity, the corresponding error maps are provided below each estimated depth map. Please zoom in for detailed views.

9. Qualitative Comparisons

Qualitative comparisons on two zero-shot, dynamic, and real-world video depth benchmarks, among DepthCrafter [26], MonST3R [86], and StereoDiff are illustrated in Fig. 11 and 12 for Bonn [39] dataset and Fig. 13 and 14 for KITTI [19] dataset. For better clarity, Fig. 13 and 14 include the corresponding error map below each estimated depth map. In static regions, StereoDiff effectively utilizes stereo matching to deliver highly robust and stable video depth estimations. This approach minimizes large depth shifts commonly observed in DepthCrafter [26], where depth values on static backgrounds can largely vary between adjacent windows. This demonstrates the advantage of stereo matching in enhancing global consistency. In dynamic regions, StereoDiff excels in maintaining smooth local consistency across consecutive frames, addressing challenges posed by both the object motion and camera movement. In contrast, MonST3R [86] suffers from pronounced flickering and depth jitters in these areas. StereoDiff’s one-step denoising process (the second stage), performed by the video depth diffusion model, ensures smoother and flicker-free depth predictions.

StereoDiff’s two-stage pipeline synergizes the strengths of stereo matching and video depth diffusion. For static regions, stereo matching ensures robust and strong global consistency thanks to the global 3D constraints. Meanwhile, the video depth diffusion stage greatly improves the smoothness of depth maps especially in dynamic areas, without sacrificing the obtained global consistency. This two-stage approach effectively addresses the distinct challenges of static and dynamic regions, delivering a comprehensive solution for consistent and accurate video depth estimation.

As for visualization, all depth maps use the same colormap (Spectral) from the `matplotlib` library. Before visualization, both predicted and GT depth maps are normalized by the maximum depth value of the GT. Notably, the color tones of depth maps differ largely between Bonn and KITTI. In KITTI, the minimum depth values for both GT and predictions are close to 0, resulting in normalized values covering the full range of $0 \sim 1$ and basically showing the entire colormap. In contrast, for Bonn, the minimum depth values are typically in the range of $3 \sim 5m$. After normalization, the depth maps may span merely $0.3 \sim 1$, leading to inadequate red tones in the visualizations.

10. Video Results for Better Comparisons

Please see the attached `*.mp4` files for the video results comparisons on Bonn [39] dataset, among DepthCrafter [26], MonST3R [86], and StereoDiff. We eliminate the video results for KITTI [19] because of KITTI’s video length is much shorter than Bonn. The average video length of KITTI is ~ 100 frames while for Bonn the average length is ~ 500

frames. Longer video depth results can more clearly show the superior global and local consistency of StereoDiff over DepthCrafter [26] and MonST3R [86].

11. Limitations and Benchmarks Selection

Limitations. The limitation of StereoDiff mainly stems from its first stage, which is a stereo matching process designed to achieve robust and strong global consistency through global 3D constraints. SfM methods [18, 38, 50, 51, 54, 63, 64, 68, 69, 86] inevitably face failure cases due to various limitations. These include challenges with textureless or repetitive surfaces, constantly changing lighting conditions, and computational challenges in large-scale scenarios. While improvements can reduce failures, the various limitations cannot be entirely avoided. In the KITTI dataset [19], the quality of 3D structures varies significantly across scenes due to two main factors: 1) the extremely limited resolution, as images are down-scaled from 1216×352 to 512×144 ; and 2) large proportions of textureless surfaces, such as roads. We display both the good cases and failure cases in Fig. 15 for clearer understanding.

Benchmarks Selection. As discussed in Sec. 1 and 4.1 of the main paper, StereoDiff is validated on two well-acknowledged, zero-shot, dynamic, and real-world video depth benchmarks: Bonn [39] for indoor scenes and KITTI [19] for outdoor scenes. The complete Bonn dataset comprises 24 dynamic scenes and 2 static scenes, totaling $\sim 12,000$ frames. The motions of dynamic scenes can be classified into 3 categories: 1) one moving object and one moving person, 2) only one moving person, and 3) two moving persons. For the diversity of motions and evaluation efficiency, 6 dynamic scenes are selected⁹, covering all 3 motions categories:

Scene Name	Motion Category
balloon	1
balloon2	1
person_tracking	2
person_tracking2	2
synchronous	3
synchronous2	3

For KITTI [19], following prior works [26, 86], the validation set is employed for evaluation. The validation set contains 13 outdoor dynamic scenes. One scene, `2011_09_26_drive_0036`, is excluded due to extreme failures, as illustrated in Fig. 15 (a). The remaining 12 scenes are used for validation. While some scenes demonstrate good 3D structures, as seen in (b, c) of Fig. 15, others can be less optimal. Please refer to Fig. 16 for additional examples of 3D structures in KITTI.

⁹For `balloon_tracking` and `balloon_tracking2`, both a dynamic object and a moving person are included.



(a) Failure case (KITTI's 2011_09_26_drive_0036)



(b) Good case (2011_09_26_drive_0095) (right view)

(c) Good case (left view)

Figure 15. **3D structures of scenes in KITTI [19] dataset.** We display both the failure cases (a) and good cases from different views (b, c).



(a) 2011_09_26_drive_0005



(b) 2011_09_26_drive_0013



(c) 2011_09_26_drive_0023



(d) 2011_09_29_drive_0026

Figure 16. Additional 3D scene structures in KITTI [19] dataset, for a more comprehensive understanding.