

M2SVid: End-to-End Inpainting and Refinement for Monocular-to-Stereo Video Conversion

Nina Shvetsova^{1,2,3,4*}, Goutam Bhat¹, Prune Truong¹, Hilde Kuehne^{2,3}, Federico Tombari^{1,5}

¹Google, ²Tuebingen AI Center/University of Tuebingen, ³Goethe University Frankfurt,

⁴MPI for Informatics, Saarland Informatics Campus, ⁵Technical University of Munich

Project webpage: <https://m2svid.github.io/>

Abstract

We tackle the problem of monocular-to-stereo video conversion and propose a novel architecture for inpainting and refinement of the warped right view obtained by depth-based reprojection of the input left view. We extend the Stable Video Diffusion (SVD) model to utilize the input left video, the warped right video, and the disocclusion masks as conditioning input to generate a high-quality right camera view. In order to effectively exploit information from neighboring frames for inpainting, we modify the attention layers in SVD to compute full attention for disoccluded pixels. Our model is trained to generate the right view video in an end-to-end manner without iterative diffusion steps by minimizing image space losses to ensure high-quality generation. Our approach outperforms previous state-of-the-art methods, being ranked best $2.6\times$ more often than the second-place method in a user study, while being $6\times$ faster.

1. Introduction

Emerging technologies such as Mixed-Reality headsets and glasses allows users to easily experience immersive content, thanks to the use of separate displays for left and right eyes. Rendering videos from left and right eye viewpoints on these displays creates a stereoscopic 3D effect, giving viewers a sense of depth. However, capturing such stereoscopic videos usually requires specialized cameras that can capture both left and right eye perspectives simultaneously. While manual monocular-to-stereo video conversion is possible, it is costly and time-consuming. This drives the need for automated, fast, and high-quality conversion methods that can enable large scale conversion of videos.

Recent advancements in video generation [4, 5] as well as monocular depth estimation [29, 78] have led to significant improvements in the monocular-to-stereo conversion

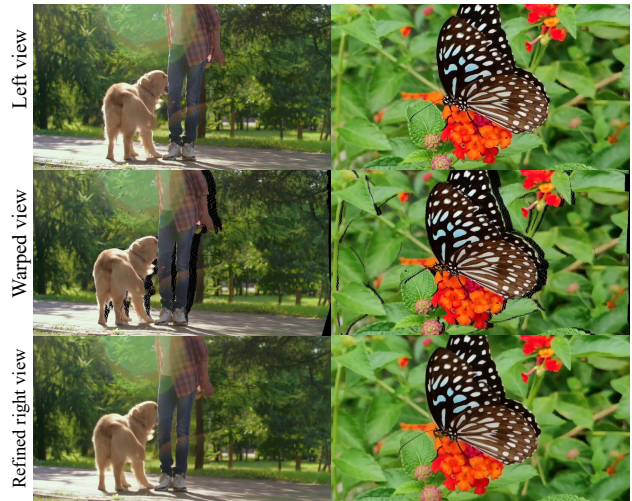


Figure 1. We present M2SVid, an end-to-end video inpainting and refinement approach for monocular-to-stereo video conversion task. Given an initial right view, obtained via *e.g.* depth based warping, our method inpaints the missing region and refines the artifacts introduced by warping in a feed-forward manner.

quality. A common approach is to use depth maps generated by a monocular depth model to warp the input video to a virtual right camera and then inpaint the disoccluded regions [9, 44, 63]. SVG [9] repurposes a pre-trained video diffusion model for the inpainting task, while [44] trains a custom model. Rather than only modifying disoccluded regions, the recent concurrent work StereoCrafter [63] fine-tunes a video diffusion model [3] to inpaint and refine the full warped video, using only the warped view and disocclusion masks as input. While obtaining improved results, StereoCrafter still leverages an architecture primarily designed for monocular video inpainting task, without exploiting the constraints specific to stereoscopic video-refinement. Moreover, the inference requires many diffusion steps leading to high computational cost.

In this work, we propose a novel architecture for efficient inpainting and refinement in the monocular-to-stereo conversion task. Our approach is designed to leverage im-

*This work was conducted during an internship at Google.

portant additional cues available in this task. Firstly, in the areas that are not disoccluded, the right view largely resembles the left one, where the image content is just horizontally shifted. Thus, even if depth-based warping introduces artifacts in the generated right view, we conjecture that most of them can be fixed by relying on left view information. Secondly, the disoccluded regions are generally small in most cases. Furthermore in the presence of camera or scene motion, the disoccluded regions in one frame are often visible at nearby frames, thus simplifying the inpainting task into a spatial-temporal matching problem. The aforementioned aspects motivate us to design M2SVid (**Monocular-to-Stereo Video** conversion), a specialized architecture which performs inpainting and refinement in an feed-forward manner, *without any iterative diffusion steps*.

Our main contributions are the following. First, we extend the Stable Video Diffusion (SVD) architecture by employing *the input left view*, in addition to the warped right view and disocclusion mask as conditioning for the refinement task. This makes it easier for the model to propagate information from the uncorrupted left view to the target right view, thus preserving high frequency details such as fine structures. Secondly, rather than performing temporal attention only over the same spatial locations as in SVD [3], we compute *full cross-attention for disoccluded pixels*. This provides greater flexibility to propagate useful information from either nearby temporal frames or the associated left view, with a limited computational overhead. Thirdly, we train our model in an *end-to-end manner* using perceptual and fidelity based losses, enabling high quality inpainting and refinement (See Fig. 1). Finally, unlike other methods, we train on publicly available datasets, making our approach reproducible and enabling fairer comparisons.

We perform quantitative evaluation, qualitative analysis, as well as a human perception user studies to evaluate our method. It shows that our M2SVid outperforms previous state-of-the-art StereoCrafter and SVG, being ranked best $2.6\times$ more often than StereoCrafter and $4.75\times$ more often than SVG, while running $6\times$ and $600\times$ faster, respectively.

2. Related Work

2.1. Novel View Synthesis

Early methods for novel-view synthesis performed 3D reconstruction of the full scene given dense input views [26, 60, 79], which can then be used to render novel views. While Multi-Plane Image (MPI)-based approaches [16, 68, 69] were successfully used to generate multiplane representations from single images, NeRFs [47] marked a milestone in the field, and spawned many efforts to improve quality [1, 2], reduce latency [17, 49, 50, 73, 80], extend to large-scale scenes [70], explicit representation [30], pose-free setup [7, 36, 46, 67, 72, 76], or fewer inputs [10, 21, 31,

51, 54, 67, 75], including just a single image [38–40, 45]. Approaches such as CAT3D [13], NVS-Diffusion [27], MultiDiff [48] directly generate the queried camera view using a diffusion models. A few works also extend these to videos [74]. While obtaining promising results, the generated views can still have artifacts due to the inherent difficulties of rendering views with large viewpoint changes.

2.2. Monocular-to-Stereo Conversion

Unlike general novel-view synthesis, monocular-to-stereo conversion renders a fixed right view from a left view, enabling monocular content to be experienced in 3D on mixed-reality headsets. Most prior works can be grouped into two families. One line of work aims to directly generate the right image/video given the left input [34, 62, 77, 77]. A notable example is Deep3D [77], which leverages an internal disparity representation to directly render the right video given the left one. The second type of approaches employ an explicit depth map to reproject the left image to the right camera, and then inpaint the disoccluded regions [9, 9, 22, 33, 44, 63, 64, 71]. Notably, Wang *et al.* [71] learn a diffusion model for inpainting disoccluded pixels using self-supervision with random cycle rendering. SVG [9] uses a pre-trained video diffusion model for inpainting. While the method can achieve spatio-temporal consistency in the output, the inpainting result can be incorrect due to the lack of task-specific fine-tuning. Mehl *et al.* [44] aim to instead perform a geometry aware inpainting of the disoccluded areas using local background information, rather than generating arbitrarily realistic content. StereoCrafter [63], which is concurrent to our work, fine-tunes a video diffusion model for inpainting the disoccluded areas, on a large stereo dataset (not released). Another recent work SpatialDreamer [42] mitigates the necessity of paired stereoscopic training data by proposing a self-supervised training framework using a forward-backward rendering mechanism. Note that a number of the above methods do not release code or model, and are often trained on private collected data, making a fair comparison difficult.

2.3. Diffusion Models

Diffusion Models [18, 55, 65] are generative models that iteratively denoise an input to produce the output. They have been tremendously successful for the text-to-image [18, 53, 55, 58, 82] and text-to-video [4, 19] tasks. They have also been successfully employed for diverse computer vision tasks, thanks to the strong image priors learned by the models. These include image [8, 37, 41] and video [6, 35, 87] inpainting, novel view synthesis [25, 48], as well as monocular depth estimation [14, 20, 29, 61], semantic segmentation [28], or normal estimation [12]. Conventionally, diffusion models employ multiple denoising steps during inference, leading to large computations cost. A number of

approaches aim to train one-step models to address this, using knowledge distillation [81, 85], GAN training [43], low-rank adaptors [85]. Garcia *et al.* [14] instead finetune a pre-trained diffusion model to perform direct feed-forward monocular depth estimation using end-to-end training. In this work, we show that such an end-to-end training strategy can also be employed for inpainting disoccluded regions in monocular-to-stereo video conversion task.

3. Conditional Latent Diffusion Models

Here, we provide a brief background on diffusion models employed in our method. Denoising Diffusion Probabilistic Models (DDPMs) [18] are generative models trained to map a simple noise distribution p_T to the data distribution p_0 , by reversing a stochastic forward process p_t , $t = 1, \dots, T$. A denoising model $\hat{v}_\theta(x_t, t)$ is trained to generate an image from pure noise by progressively denoising inputs x_t at time stamp t . In order to reduce computational complexity, Latent Diffusion Models (LDMs) [55] operate in a latent space of a Variational Autoencoder (VAE) [32]. Conditional diffusion models [57, 83] extend the denoising model $\hat{v}_\theta(x_t, t, c)$ to be conditioned on additional input c , such as text [57], images [57], *etc.* to control the generation process.

Training and inference with conditional LDMs: During training, given a data sample \mathbf{x} (e.g., an image or video) and its corresponding conditioning input c (e.g., text or another image), the data sample \mathbf{x} is first encoded into a latent representation $\mathbf{z} = E(\mathbf{x})$ using VAE encoder E . The latent representation \mathbf{z} is perturbed through a forward diffusion process: $\mathbf{z}_t = \sqrt{\bar{\alpha}_t}\mathbf{z} + \sqrt{1 - \bar{\alpha}_t}\epsilon$ where $\epsilon \sim \mathcal{N}(0, I)$ is Gaussian noise, and $\bar{\alpha}_t$ controls the noise schedule. Instead of predicting the noise ϵ directly, diffusion models usually leverage v -parameterization [59], where the model $\hat{v}_\theta(\mathbf{z}_t, t, c)$ learns to predict: $\mathbf{v} = \alpha_t\epsilon - \sqrt{1 - \bar{\alpha}_t}\mathbf{z}$. Therefore the model \hat{v}_θ , typically a U-Net [56], is trained to reconstruct \mathbf{v} , by optimizing the objective:

$$\mathcal{L} = \mathbb{E}_{\mathbf{z}, c, \epsilon \sim \mathcal{N}(0, I), t \sim \mathcal{U}(T)} [\|\mathbf{v} - \hat{v}_\theta(\mathbf{z}_t, t, c)\|^2]. \quad (1)$$

During inference, the denoising process starts from pure noise \mathbf{z}_T and the learned denoiser $\hat{v}_\theta(\mathbf{z}_t, t, c)$ iteratively refines the output over T steps to recover the final sample.

Diffusion as feed-forward models: Garcia *et al.* [14] propose to use pre-trained diffusion U-Net model as a feed-forward models for pixel-to-pixel tasks such as depth and normal estimation. In this scenario, the timestep t is not sampled anymore and fixed as $t = T$ to train the model for single-step prediction. The noise ϵ is additionally replaced by the mean of the noise distribution, *i.e.*, zero. With $t = T$, we get $\bar{\alpha}_T \approx 0$, and thus $\mathbf{z}_T = \sqrt{\bar{\alpha}_T}\mathbf{z} + \sqrt{1 - \bar{\alpha}_T}\epsilon \approx \mathbf{0}$ and $\mathbf{v} = \alpha_T\epsilon - \sqrt{1 - \bar{\alpha}_T}\mathbf{z} \approx -\mathbf{z}$. As such, the model directly learns to predict the clean output $\mathbf{v} \approx -\mathbf{z}$ starting from zero

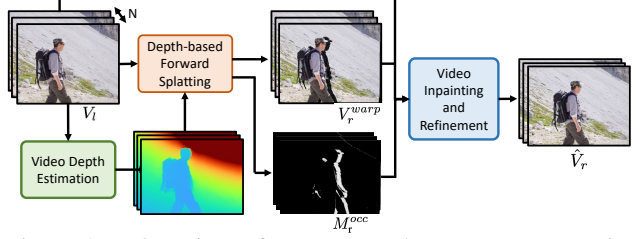


Figure 2. Overview of our monocular-to-stereo conversion pipeline. Given an input monocular video, we first estimate per-pixel depth, which is used to warp the input video to a right camera view. The input video, the warped video, as well as the disocclusion masks are then passed to our video inpainting and refinement module to generate the final right view.

noise vector $\mathbf{z}_T = \mathbf{0}$. With such timestep sampling and assumptions, the standard loss (Eq. (1)) converges to the L_2 loss in the latent space between the ground truth $-\mathbf{z}$ and the predicted latents $\hat{v}_\theta(\mathbf{0}, T, c)$:

$$\mathcal{L}_{latent} = \mathbb{E}_{\mathbf{z}, c} [\|(-\mathbf{z}) - \hat{v}_\theta(\mathbf{0}, T, c)\|^2]. \quad (2)$$

However, [14] proposes decoding the predicted latents, $\hat{\mathbf{z}} = -\hat{v}_\theta(\mathbf{0}, T, c)$, using a VAE decoder D to obtain the reconstructed output $\hat{\mathbf{x}} = D(\hat{\mathbf{z}})$. Instead of minimizing the latent space error Eq. (2), they propose to optimize a task-specific loss directly in image space with $\mathcal{L} = \mathcal{L}_{task}(\mathbf{x}, \hat{\mathbf{x}})$.

4. Monocular-to-Stereo Conversion Pipeline

In this work, we tackle the problem of converting any 2D video to a 3D stereoscopic video. Given a monocular video $V \in \mathbb{R}^{N \times H \times W \times 3}$, containing N frames of resolution $H \times W$, the goal is to transform it into a *stereoscopic video pair* (V_l, V_r) , where $V_l, V_r \in \mathbb{R}^{N \times H \times W \times 3}$ represent the views corresponding to the left and right eyes, respectively. Following prior work, we assume that the input video V corresponds to the left-eye view, *i.e.*, $V_l = V$, and therefore, the goal is to only generate the right-eye view V_r from V . Note that multiple possible right-eye views V_r exist depending on the baseline, *i.e.* the distance between cameras. Hence we condition the generation on a maximum desired disparity D_{max} between V_l and V_r , that can be set by the user to control the stereoscopic effect in practical applications.

An overview of the pipeline is shown in Fig. 2. We follow the commonly employed strategy [44, 63] of performing depth-based reprojection to generate an estimate of the right view, as described next.

Video-depth estimation: We use a monocular video depth estimation model to obtain per-frame depths. Our method is agnostic to the choice of depth estimation method.

Left-to-right video warping: We use the user-provided maximum disparity D_{max} to convert the depth maps to disparities. The disparity maps are then used to forward-splat the input left-view video to a virtual right view denoted as

V_r^{warp} . Note that the warped right frames inevitably contain missing values (*i.e.* holes) introduced by disocclusions. These missing pixels are indicated by the disocclusion mask M_r^{occ} , obtained as a by-product of the splatting.

Video inpainting and refinement: The final step in our pipeline is to generate a high quality right view given the warped video V_r^{warp} . As mentioned before, the warped right video V_r^{warp} contains holes introduced by disocclusions. Furthermore, the warped right video V_r^{warp} can also contain artifacts introduced due to reprojection and depth errors. In this work, we mainly focus on this final “video inpainting and refinement” task. To this end, we introduce an efficient end-to-end model to inpaint the disoccluded regions and fix reprojection artifacts. Our architecture is described in detail in the next section.

5. End-to-End Stereoscopic Video Refinement

We introduce M2SVid, an efficient end-to-end model for generating a high-quality *right video* from the *warped right video* V_r^{warp} . While this problem is similar to the standard video inpainting task, it comes with a set of unique properties, as described next.

(1) In our setup for rendering a right camera view, inpainting holes appear in regions where depth increases from left to right. These regions correspond to background, so unlike general-purpose inpainting (e.g., object removal), our model can rely on context from only one side of the hole.

(2) In a general inpainting setup, the model is only required to generate the missing pixels, while keeping the rest of the image unaltered. In our case, however, regions outside the inpainting mask may also contain artifacts from depth-map errors or interpolation during forward splatting. These artifacts are especially common in presence of thin structures such as fences, where the estimated depth can have substantial errors. Fortunately, we can leverage the original left video to correct the errors introduced by the reprojection.

(3) In the image inpainting task, the model *must hallucinate* content in missing regions due to the lack of additional information. A video inpainting model can instead use information from multiple frames, though large holes still require substantial hallucination. In our task, however, the inpainting regions are much thinner (*i.e.*, as controlled by the maximum disparity D_{max}). As a result, the inpainting problem is greatly simplified and the model *can copy* information from other temporal frames to avoid hallucination. We develop a custom architecture for the right view generation task aiming to exploit the aforementioned properties.

5.1. Overview

An overview of our method is shown in Fig. 3. We base our model on a strong pre-trained video diffusion model to benefit from learned video priors. In particular, we utilize

Stable Video Diffusion (SVD) [3], a latent video diffusion model trained to generate videos from an input image. We customize the SVD architecture for the stereoscopic video generation task as follows. First, we condition each latent generation on the input left video V_l , reprojected right video V_r^{warp} , as well as the disocclusion masks M_r^{occ} . Secondly, only for the inpainted pixels, we extend the spatial attention in SVD to operate over all pixels in all frames, instead of just the pixels in the same spatial location. This gives greater flexibility for the model to copy visible pixels from neighboring frames in order to consistently inpaint. Finally, we train our model *end-to-end* on public stereoscopic datasets to generate the refined right view in a *feed-forward manner*, without requiring multiple denoising steps.

5.2. High-Frequency Details Preservation

Conventional video inpainting methods utilize masked video to condition inpainting models, a strategy widely adopted in stereo conversion methods [44, 63, 71], where the inpainting model is conditioned solely on the warped right video V_r^{warp} and the disocclusion mask M_r^{occ} . However, as discussed in Sec. 5 (2), the warped right video is likely to contain artifacts even in non-inpainting regions. Thus we propose to also utilize the original left video V_l as an extra input to the model. In more detail, the standard SVD model takes the noisy latents for a video snippet and the VAE encoding of a conditioning image as input. We modify the SVD architecture to instead take the following inputs as conditioning: the VAE encoding of the left video snippet $E(V_l)$, the VAE encoding of the warped right snippet $E(V_r^{warp})$, and the disocclusion mask $M_r^{occ, resized}$, resized to the same resolution as the VAE encodings. We thus denote the model $\hat{v}_\theta(z_t, t, c)$ with conditioning inputs c as

$$c = [E(V_l), E(V_r^{warp}), M_r^{occ, resized}], \quad (3)$$

where $[\cdot]$ refers to concatenation. This modification is achieved by modifying the first convolution layer in the U-Net, as is the common practise.

Using the original left video as conditioning allows the model to easily infer high frequency details as well as other information which might have been destroyed during the depth based reprojection stage. As shown in Fig. 5, this improves the quality of the generated right views, with high-frequency details from the left view better preserved.

5.3. Spatio-Temporal Aggregation for Inpainting

The SVD model takes multiple video frames (*i.e.* a snippet) as input and jointly denoise them to produce a temporally consistent output. Ideally, one would compute a full attention over all spatial tokens in all frames for maximal information flow. However, this would be prohibitively expensive, resulting in a complexity of $H = N^2 \times h^2 \times w^2$ for each attention layer, where h and w are the size of the

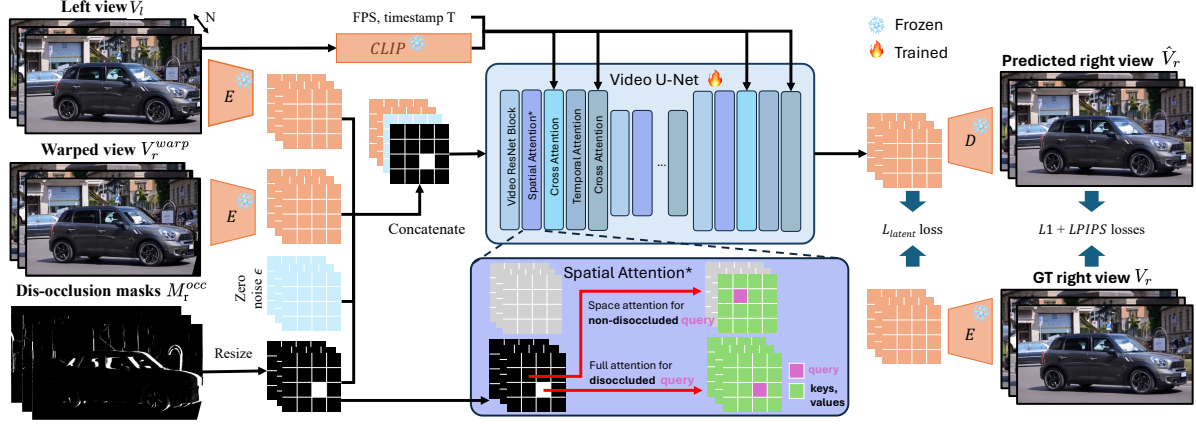


Figure 3. An overview of our proposed stereoscopic video refinement method. Our model inpaints the disoccluded regions in the warped right view, and corrects possible artifacts introduced by warping errors. The model takes the VAE encodings of the input left view, reprojected right view, and the disocclusion mask as conditioning to the U-Net. The latent encodings of the refined right view are then generated in a single denoising step, and then decoded by the VAE Decoder to generate the output right video. In order to effectively utilize the information from neighboring frames for inpainting, we extend the spatial attention layer in SVD to compute full attention for the disoccluded tokens. The model is training end-to-end by minimizing image space and latent space losses.

latent representations. For this reason, the standard video diffusion models such as SVD and VideoCrafter [5] factorize the full attention computation into interleaved spatial and temporal attention layers. In the spatial attention layers, only tokens from the same timestamp attend to each other, while in the temporal attention layers, tokens from the same spatial location attend to each other across different timestamps. This drastically reduces the complexity to $N \times h^2 \times w^2 + N^2 \times h \times w$, which is more palatable.

However, factorizing full attention can reduce modeling capacity for our task, especially in dynamic scenes with camera motion. As discussed in Sec. 5(3), an occluded pixel in one frame is often visible at a different spatial location in another, simplifying inpainting problem. In such cases, we would prefer full spatio-temporal attention. Fortunately, the number of inpainted tokens constitute only a small fraction of all tokens (usually $< 5\%$). We exploit this property and modify the spatial attention layer in SVD to allow tokens corresponding to the disocclusion mask M_r^{occ} to attend to all other tokens, while using the spatial attention only for other tokens. This improves the inpainting capability of our model (Fig. 6) without significant computational overhead.

5.4. Efficient Feed-Forward Prediction

The diffusion models for text-to-image/video generation and inpainting commonly utilize multiple denoising steps to generate the final image [18, 55, 57]. This is important due to the ill-posed nature of the problem wherein multiple solutions exists for a single input prompt. Utilizing fewer denoising steps usually results in blurry output. However, in our task of right view generation, the inpainting regions are generally small. Furthermore, as described in Sec. 5, the neighboring temporal frames as well as the input left frames

contain a significant amount of information needed to generate a high-quality, temporally consistent right video. The model therefore does not need to perform significant hallucination in most cases, but rather needs to fetch the relevant information from other frames. The relatively constrained nature of our problem thus motivates us to generate the right view in a direct feed-forward manner, rather than the multi-step denoising used in prior work [9, 63], similar to the approach employed by Garcia *et al.* [14] for depth estimation.

Given the input left video V_l , warped right video V_r^{warp} , and the disocclusion masks M_r^{occ} , we generate the conditioning input c to the model by obtaining the VAE encodings of left and reprojected videos, as shown in Eq. (3). The initial latent is set to the mean noise $\mathbf{0}$ as in [14] and the timestamp is set to the highest noise $t = T$. As discussed in Sec. 3, this choice of initial latent and timestamp leads the Video U-Net to directly reconstruct a clean latent due to the use of v -parameterization [59]. Therefore, the conditioning input, together with the initial latent are then passed through the Video U-Net $\hat{\mathbf{v}}_\theta$ to obtain the latent encoding $\hat{z} = -\hat{\mathbf{v}}_\theta(\mathbf{0}, T, c)$ of the right video. This is then passed through the VAE decoder D to obtain the predicted right view \hat{V}_r . The inference pipeline can thus be denoted as,

$$\hat{V}_r = D(-\hat{\mathbf{v}}_\theta(\mathbf{0}, T, c)) \quad (4)$$

Our feed-forward prediction strategy significantly reduces the latency of the right view generation step compared to existing methods, as shown in Table 2. Furthermore, this also allows us to train the model end-to-end using image quality losses w.r.t. to the ground truth right video V_r , in order to maximize the quality of the generated video. In particular, we train the model using a combination of LPIPS and L_1 losses directly in the image space, along with an



Figure 4. Qualitative comparison of our approach with state-of-the-art methods SVG [9] and StereoCrafter [63]. Our approach can effectively preserve the high-frequency information from the input video and generate high-quality right views.

auxiliary loss Eq. (2) in the latent space. Our final training loss is thus defined as,

$$\mathcal{L} = \mathcal{L}_{latent}(z, \hat{z}) + \mathcal{L}_{LI}(V_r, \hat{V}_r) + \mathcal{L}_{LPIPS}(V_r, \hat{V}_r) \quad (5)$$

Here z corresponds to the VAE encodings of the ground truth right video V_r . Note that during training, we keep the VAE decoder frozen and only fine-tune the U-Net model.

5.5. Training Data

To train our model using the supervised objective Eq. (5), we require a dataset stereoscopic videos with their corresponding disparity maps. Due to the lack of large-scale stereoscopic datasets, prior works focused on collecting private datasets that usually have not been publicly released. Instead we aim to train using publicly available datasets.

We utilize **Ego4D** [15] and **Stereo4D** [23] datasets for our training. Ego4D a large-scale egocentric video dataset containing approximately 263 long videos (80 hours in total) collected using a stereo camera, while the recently released **Stereo4D** [23] dataset consists of $\sim 200K$ stereoscopic video clips sourced from $\sim 7K$ online videos. The Stereo4D dataset also provides rectified videos and disparity maps. The Ego4D dataset on the other hand contains unrectified videos, without any precomputed disparity maps.

We thus perform the following steps to preprocess Ego4D dataset. First, we uniformly sample frames and perform dense feature matching with LoFTR [66], using these matches to compute the fundamental matrix [86] via RANSAC [11]. Next, we estimate rectification transformations for both views and rectify the videos. We then compute the LoFTR matches again and shift the left and right views horizontally until all disparities between the matched points are positive and the smallest disparity is zero. This ensures that the videos follow a rectified stereo setup. Finally, we use an off-the-shelf stereo matching method BiDaVideo [24] to obtain disparity maps for all stereo pairs.

6. Experimental Results

In this section, we evaluate the quality of the right view videos generated by our method, both qualitatively as well as quantitatively. Further results, analysis, visualizations and implementation details are provided in the Appendix.

Quantitative evaluation: To quantitatively evaluate our approach, we need datasets containing left and right stereo videos, together with the (pseudo) ground truth disparity maps which are used to generate the warped right views for each method. This allows us to directly compare the generated right views with the ground truth right views, without

Method	Training data	Denoising steps	PSNR \uparrow	MS-SSIM \uparrow	LPIPS \downarrow
SVG [9]	- (training free)	50 steps	25.6	0.926	0.217
StereoCrafter [63]	private dataset	25 steps	24.9	0.909	0.242
StereoCrafter [63]	private dataset	1 step	25.3	0.911	0.262
M2SVid (Ours)	Stereo4D + Ego4D	1 step	26.2	0.915	0.180

Table 1. State-of-the-art comparison on Stereo4D test set. Our approach obtains the best scores in terms of PSNR as well as LPIPS.

Method	Denoising steps	Average rank \downarrow	# Chosen best (rank=1) \uparrow	Latency (s) \downarrow
SVG [9]	50 steps	2.88	16 / 112	1270.4
StereoCrafter [63]	25 steps	<u>2.05</u>	<u>29 / 112</u>	12.2
StereoCrafter [63]	1 step	3.46	4 / 112	<u>2.4</u>
M2SVid (Ours)	1 step	1.43	76 / 112	2.1

Table 2. Desktop human perception study (112 rankings, 21 videos, 13 participants), with methods ranked from 1 (best) to 4 (worst). Our method achieved an average rank of 1.43, significantly outperforming others – M2SVid ranked best $2.6\times$ more often than StereoCrafter (25 steps) and $4.75\times$ more often than SVG – while being $6\times$ and $635\times$ faster, respectively. Runtimes were measured on an A100 GPU for 512×512 16-frame videos.

VR headset comparison	StereoCrafter [63] (25 steps) preferred	Ours (1 step) preferred
StereoCrafter vs. Ours	9 / 105	39 / 105

Table 3. VR headset human perception study with 105 comparisons over 21 videos with 5 distinct users. Our method shows a clear advantage over StereoCrafter with 25 denoising steps.

having to worry about errors introduced by incorrect depth estimation. We rely on **Stereo4D** [23] and **Ego4D** [15] test sets and use the standard image quality metrics PSNR, MS-SSIM, and LPIPS [84] for evaluation. We compute the metrics independently over each video, and average them over the full dataset. In the appendix, we also report metrics inside and outside disoccluded regions. All evaluations are performed on 16-frame videos, sampled at 5 FPS, resized to 512 resolution, and centrally cropped.

Qualitative evaluation: For qualitative analysis and user studies, we only need monocular videos along with per-frame depth maps. We use videos from the **DAVIS** dataset [52] as well as videos from free online stock sources¹. We sample 16 frames per video at 8 FPS and re-size them to 640×1152 resolution. We compute the depth maps using the recent DepthCrafter [20] method.

6.1. Comparison to State-of-the-Art

Baselines: We compare our method to the recent state-of-the-art monocular-to-stereo conversion models **SVG** [9] and **StereoCrafter** [63]. SVG is a training-free method that utilizes a frozen diffusion model to inpaint regions within the disocclusion mask while not modifying the remaining regions. StereoCrafter, which is an unpublished concurrent work, fine-tunes SVD 1.1 [3] to inpaint and refine the warped right views in a diffusion-based denoising setup, using the warped right views along with the disocclusion masks as conditioning. We evaluate StereoCrafter using

the default 25 denoising steps as well as single denoising step. We exclude the Deep3D [77] from the evaluation as it doesn’t allow control over the baseline between the stereo cameras and is shown to have inferior results [9, 63] compared to SVG and StereoCrafter. In order to ensure a fair comparison, we use the same depth maps for reprojection in all methods. However, we use the official warping implementation provided by the authors for each method.

Qualitative Results: We perform a qualitative comparison with the SVG and StereoCrafter methods in Fig. 4. Since SVG only performs inpainting within the disocclusion mask, it fails to fix the artifacts introduced by errors in warping. Furthermore, it can incorrectly extend the ‘foreground’ object to fill the inpainting hole due to lack of task-specific fine-tuning. StereoCrafter with a single denoising step produces blurry inpainting and degrades the quality of the warped areas as well. This is expected since the model was trained using a multi-step denoising objective. When using 25 denoising steps, StereoCrafter inference generally produces sharp results. However it can struggle at times to correctly generate the high-frequency details. This is because the high-frequency information can often get degraded during the warping step. Since StereoCrafter only utilizes warped view and masks as conditioning, it can struggle to recover the details. In contrast, our model is conditioned on the input left view as well, allowing it to leverage the full context for inpainting and refinement. Furthermore, since our method is trained in an end-to-end manner with image space losses, it can learn to minimize the loss of high-frequency information introduced by VAE decoder.

Quantitative Results: We quantitatively compare our approach with SVG and StereoCrafter on Stereo4D in Tab. 1. Despite using only a single-step inference, M2SVid significantly outperforms both baselines on all metrics, except for SVG on MS-SSIM, where SVG benefits in MS-SSIM from preserving non-empty pixels (warped using ground truth disparity), while our method incurs a slight MS-SSIM drop due to VAE compression, but achieves superior visual quality by correcting warping errors, e.g., in thin structures.

User studies: We also perform two user studies. First, in a *desktop* user study, participants viewed a random subset of 21 videos from the DAVIS dataset and public sources. For each video, anonymized right views from each method were shown alongside the input left view and disocclusion mask (as reference). Each of the 13 participants were asked to rank the quality of the generated videos from 1 (best) to 4 (worst), taking into account factors such as temporal consistency, image quality (i.e. sharpness) and lack of artifacts. If methods are indistinguishable, the equal ranking was allowed. In total, 112 rankings per method were collected (Tab. 2). Our method significantly outperforms all others, obtaining an average ranking of 1.43, compared to 2.05 ob-

¹<https://www.pexels.com/>

tained by StereoCrafter (25 denoising steps) and 2.88 by SVG. In fact, our method was ranked first $2.6\times$ more often than StereoCrafter (25 steps), and $4.75\times$ than SVG.

To validate that the enhanced visual quality of our method also improves user experience in stereoscopic viewing, we conducted a second user study using a *VR headset*. Participants were shown anonymized stereo videos generated by our method and StereoCrafter (25 denoising steps) and asked which version they preferred, or if both were equal. We collected 105 comparisons from 5 users (21 each). As shown in Tab. 3, our method was preferred in 39 cases, while StereoCrafter was favored in only 9 (the rest are tied), highlighting the clear advantage of our approach.

Run-time: Furthermore, our efficient refinement achieves a $6\times$ and $635\times$ speed-up compared to StereoCrafter with 25 steps and SVG, respectively, on an A100 GPU (Tab. 2).

6.2. Ablation Study

We ablate the key components of our approach in this section on the Ego4D and Stereo4D datasets.

Left view conditioning (Sec. 5.2): The impact of using the left view as conditioning, in addition to using the warped view and disocclusion mask is shown in Tab. 4. Conditioning on the left video shows a 6.0%, 5.8% and 9.8% improvement in PSNR, MS-SSIM and LPIPS metrics respectively on Ego4D with 25 denoising steps. A similar improvement

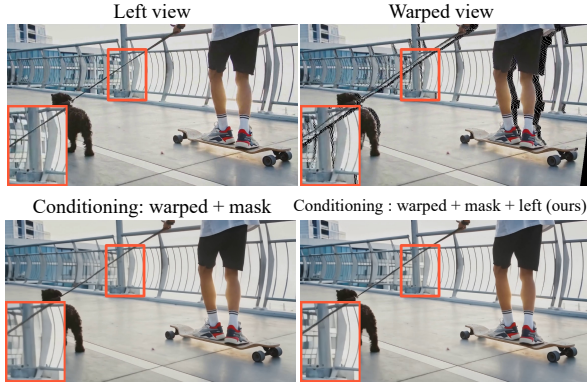


Figure 5. Impact of our left view conditioning (Sec. 5.2).

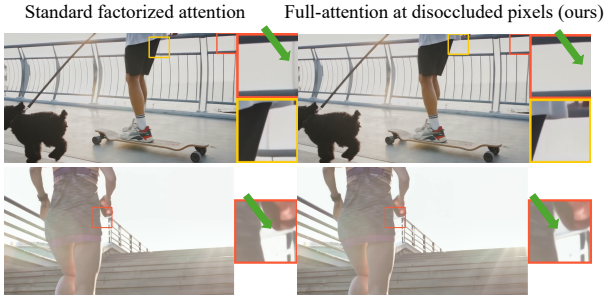


Figure 6. Using full-attention at dis-occluded pixels (Ours, Sec. 5.3) enables the model to exploit information from visible pixels in other frames to improve inpainting (see Appendix).

Model Architecture	feed forward	Inf. steps	Stereo4D			Ego4D		
			PSNR \uparrow	MS-SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	MS-SSIM \uparrow	LPIPS \downarrow
Cond. with $V_r^{warped} + M_r^{occ}$	\times	25	24.5	0.886	0.226	20.8	0.822	0.296
Cond. with $V_r^{warped} + M_r^{occ} + V_l$	\times	25	24.8	0.891	0.215	22.1	0.870	0.267
Cond. with $V_r^{warped} + M_r^{occ}$	\checkmark	1	26.1	0.913	0.187	21.5	0.837	0.276
Cond. with $V_r^{warped} + M_r^{occ} + V_l$	\checkmark	1	26.2	0.915	0.179	22.8	0.886	0.244
+full attention	\checkmark	1	26.2	0.915	0.180	22.7	0.885	0.248

Table 4. Impact of left-view conditioning and full attention at dis-occluded pixels on Stereo4D and Ego4D datasets.

Loss	feed forward	Inf. steps	Stereo4D			Ego4D		
			PSNR \uparrow	MS-SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	MS-SSIM \uparrow	LPIPS \downarrow
A Standard loss (1)	\times	25	24.8	0.891	0.215	22.1	0.870	0.267
B Standard loss (1)	\times	1	25.7	0.901	0.242	23.0	0.877	0.320
C L_{latent} (2)	\checkmark	1	25.6	0.901	0.238	22.9	0.875	0.318
D $L_{latent} + L_{LPIPS} + L_{LI}$ (5)	\checkmark	1	26.2	0.915	0.179	22.8	0.886	0.244

Table 5. Impact of different inference strategies and losses.

is observed in the feed-forward case. This is because the left video conditioning allows the model to recover high-frequency details and correct artifacts introduced during the warping, as seen in Fig. 5.

Full-attention at disoccluded pixels (Sec. 5.3): While this contribution doesn’t impact metrics in Tab. 4 (likely due to the limited number of complex scenes in the test set), our qualitative analysis shows benefits for dynamic scenes with camera motion, where both foreground and background move. In Fig. 6, we observe that our method prevents hallucinations (top example) and enables correct inpainting (bottom) of thin structures. More results are in Appendix.

End-to-end training (Sec. 5.4): In Tab. 5, we ablate our proposed end-to-end training strategy. When comparing the model trained with the standard diffusion loss (Eq. (1)) with 25 (A) or 1 (B) diffusion step at inference, the single step model leads to more blurry results than (A), as evidenced by the 12% and 20% decrease in LPIPS on Stereo4D and Ego4D respectively. End-to-end training with only latent space supervision (Eq. (2)) (C) only slightly improves the LPIPS metric. Our final training loss (Eq. (5)) (D), including LPIPS and L1 losses in image space, largely improves the image sharpness, as evidenced by the 24% and 22% improvement in LPIPS on Stereo4D and Ego4D compared to (C). Notably, we outperform the standard diffusion loss with 25 inference steps (A) on all metrics and datasets.

7. Conclusion

We introduce an end-to-end approach for stereoscopic video inpainting and refinement in this work. First, we extend the SVD model to take the input left video, warped right video, and disocclusion mask as conditioning. Next, we modify the attention layers in SVD to compute full attention for the disoccluded pixels in order to improve inpainting quality. Crucially, we perform the video refinement using a single denoising step, enabling end-to-end training with image space losses. Qualitative and quantitative experiments, as well as user studies, demonstrate that our method clearly outperforms prior state-of-the-art methods for the monocular-to-stereo video conversion task.

References

- [1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. 2
- [2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022. 2
- [3] A. Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, and Dominik Lorenz. Stable video diffusion: Scaling latent video diffusion models to large datasets. *ArXiv*, abs/2311.15127, 2023. 1, 2, 4, 7
- [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. 1, 2
- [5] Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models, 2024. 1, 5
- [6] Nicolas Chérel, Andrés Almansa, Yann Gousseau, and Alasdair Newson. Infusion: Internal diffusion for video inpainting, 2023. 2
- [7] Shin-Fang Chng, Sameera Ramasinghe, Jamie Sherrah, and Simon Lucey. GARF: gaussian activated radiance fields for high fidelity reconstruction and pose estimation. *CORR*, abs/2204.05735, 2022. 2
- [8] Ciprian Adrian Corneanu, Raghudeep Gadde, and Aleix M. Martínez. Latentpaint: Image inpainting in latent space with diffusion models. *WACV*, 2024. 2
- [9] Peng Dai, Feitong Tan, Qiangeng Xu, David Futschik, Ruofei Du, Sean Fanello, Xiaojuan Qi, and Yinda Zhang. Svg: 3d stereoscopic video generation via denoising frame matrix. *arXiv preprint arXiv:2407.00367*, 2024. 1, 2, 5, 6, 7
- [10] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer views and faster training for free. In *CVPR*, 2022. 2
- [11] Martin A. Fischler and Robert C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, 1981. 6
- [12] 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 *ECCV*, 2024. 2
- [13] Ruiqi Gao*, Aleksander Holynski*, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul P. Srinivasan, Jonathan T. Barron, and Ben Poole*. Cat3d: Create anything in 3d with multi-view diffusion models. *NeurIPS*, 2024. 2
- [14] 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, 5
- [15] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *CVPR*, 2022. 6, 7
- [16] Yuxuan Han, Ruicheng Wang, and Jiaolong Yang. Single-view view synthesis in the wild with learned adaptive multiplane images. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–8, 2022. 2
- [17] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul E. Debevec. Baking neural radiance fields for real-time view synthesis. In *ICCV*, 2021. 2
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020. 2, 3, 5
- [19] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey A. Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models. *ArXiv*, abs/2210.02303, 2022. 2
- [20] 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. 2, 7
- [21] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *ICCV*, 2021. 2
- [22] V. Jampani, Huiwen Chang, Kyle Sargent, Abhishek Kar, Richard Tucker, Michael Krainin, Dominik Philemon Kaeser, William T. Freeman, D. Salesin, Brian Curless, and Ce Liu. Slide: Single image 3d photography with soft layering and depth-aware inpainting. *ICCV*, 2021. 2
- [23] Linyi Jin, Richard Tucker, Zhengqi Li, David Fouhey, Noah Snavely, and Aleksander Holynski. Stereo4d: Learning how things move in 3d from internet stereo videos. *arXiv preprint arXiv:2412.09621*, 2024. 6, 7
- [24] Junpeng Jing, Ye Mao, and Krystian Mikolajczyk. Match-stereo-videos: Bidirectional alignment for consistent dynamic stereo matching. 2024. 6
- [25] Nikolai Kalischek, Michael Oechsle, Fabian Manhardt, Philipp Henzler, Konrad Schindler, and Federico Tombari. Cubediff: Repurposing diffusion-based image models for panorama generation, 2025. 2
- [26] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. In *NeurIPS*, 2017. 2
- [27] Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models. In *Proc. CVPR*, 2024. 2
- [28] Yasufumi Kawano and Yoshimitsu Aoki. Maskdiffusion: Exploiting pre-trained diffusion models for semantic segmentation. *IEEE Access*, 12:127283–127293, 2024. 2
- [29] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing diffusion-based image generators for monocular depth estimation. In *CVPR*, 2024. 1, 2
- [30] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 2

- [31] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *CVPR*, 2022. 2
- [32] Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013. 3
- [33] Janusz Konrad, Meng Wang, Prakash Ishwar, Chen Wu, and Debargha Mukherjee. Learning-based, automatic 2d-to-3d image and video conversion. *IEEE TIP*, 22(9):3485–3496, 2013. 2
- [34] Jiyoung Lee, Hyunjoo Jung, Youngjung Kim, and Kwanghoon Sohn. Automatic 2d-to-3d conversion using multi-scale deep neural network. *ICIP*, 2017. 2
- [35] Minhyeok Lee, Suhwan Cho, Chajin Shin, Jungho Lee, Sunghun Yang, and Sangyoun Lee. Video diffusion models are strong video inpainter. In *AAAI*, 2025. 2
- [36] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021. 2
- [37] Anji Liu, Mathias Niepert, and Guy Van den Broeck. Image inpainting via tractable steering of diffusion models. *ArXiv*, abs/2401.03349, 2023. 2
- [38] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T., Zexiang Xu, and Hao Su. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. In *NeurIPS*, 2023. 2
- [39] Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. In *CVPR*, 2024.
- [40] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *ICCV*, 2023. 2
- [41] Andreas Lugmayr, Martin Danelljan, Andrés Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. *CVPR*, 2022. 2
- [42] Zhen Lv, Yangqi Long, Congzhentao Huang, Cao Li, Chengfei Lv, Hao Ren, and Dian Zheng. Spatialdreamer: Self-supervised stereo video synthesis from monocular input. In *CVPR*, 2025. 2
- [43] Xiaofeng Mao, Zhengkai Jiang, Fu-Yun Wang, Wenbing Zhu, Jiangning Zhang, Hao Chen, Mingmin Chi, and Yabiao Wang. Osv: One step is enough for high-quality image to video generation. *ArXiv*, abs/2409.11367, 2024. 3
- [44] Lukas Mehl, Andrés Bruhn, Markus Gross, and Christopher Schroers. Stereo conversion with disparity-aware warping, compositing and inpainting. In *WACV*, 2024. 1, 2, 3, 4
- [45] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion 360° reconstruction of any object from a single image. In *CVPR*, 2023. 2
- [46] Quan Meng, Anpei Chen, Haimin Luo, Minye Wu, Hao Su, Lan Xu, Xuming He, and Jingyi Yu. Gnerf: Gan-based neural radiance field without posed camera. In *JCCV*, 2021. 2
- [47] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2
- [48] Norman Müller, Katja Schwarz, Barbara Rössle, Lorenzo Porzi, Samuel Rota Bulò, Matthias Nießner, and Peter Kontschieder. Multidiff: Consistent novel view synthesis from a single image. In *CVPR*, 2024. 2
- [49] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, 2022. 2
- [50] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty R. Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. Donerf: Towards real-time rendering of compact neural radiance fields using depth oracle networks. *Comput. Graph. Forum*, 40(4):45–59, 2021. 2
- [51] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *CVPR*, 2022. 2
- [52] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, 2016. 7
- [53] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv*, abs/2204.06125, 2022. 2
- [54] Barbara Roessle, Jonathan T. Barron, Ben Mildenhall, Pratul P. Srinivasan, and Matthias Nießner. Dense depth priors for neural radiance fields from sparse input views. In *CVPR*, 2022. 2
- [55] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 3, 5
- [56] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 3
- [57] 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. *NeurIPS*, 2022. 3, 5
- [58] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, Seyedeh Sara Mahdavi, Raphael Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *ArXiv*, abs/2205.11487, 2022. 2
- [59] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022. 3, 5
- [60] Johannes L. Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR 2016, Las Vegas, NV, USA*, pages 4104–4113, 2016. 2
- [61] Jiahao Shao, Yuanbo Yang, Hongyu Zhou, Youmin Zhang, Yujun Shen, Vitor Guizilini, Yue Wang, Matteo Poggi, and

- Yiyi Liao. Learning temporally consistent video depth from video diffusion priors. In *CVPR*, 2025. 2
- [62] Jian Shi, Zhenyu Li, and Peter Wonka. Immersepro: End-to-end stereo video synthesis via implicit disparity learning. *arXiv preprint arXiv:2410.00262*, 2024. 2
- [63] Jian Shi, Qian Wang, Zhenyu Li, and Peter Wonka. Stereocrafter-zero: Zero-shot stereo video generation with noisy restart. *arXiv preprint arXiv:2411.14295*, 2024. 1, 2, 3, 4, 5, 6, 7
- [64] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *CVPR*, 2020. 2
- [65] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2
- [66] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *CVPR*, 2021. 6
- [67] Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt, and Federico Tombari. Sparf: Neural radiance fields from sparse and noisy poses. *CVPR*, 2023. 2
- [68] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *CVPR*, 2020. 2
- [69] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *CVPR*, pages 551–560, 2020. 2
- [70] Haitem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *CVPR*, 2022. 2
- [71] Xiaodong Wang, Chenfei Wu, Shengming Yin, Minheng Ni, Jianfeng Wang, Linjie Li, Zhengyuan Yang, Fan Yang, Lijuan Wang, Zicheng Liu, Yuejian Fang, and Nan Duan. Learning 3d photography videos via self-supervised diffusion on single images. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, 2023. 2, 4
- [72] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. NeRF—: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. 2
- [73] Suttisak Wizatwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *CVPR*, 2021. 2
- [74] Rundi Wu, Ruiqi Gao, Ben Poole, Alex Trevithick, Changxi Zheng, Jonathan T. Barron, and Aleksander Holynski. Cat4d: Create anything in 4d with multi-view video diffusion models. *arxiv/2411.18613*, 2024. 2
- [75] Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P. Srinivasan, Dor Verbin, Jonathan T. Barron, Ben Poole, and Aleksander Holynski. Reconfusion: 3d reconstruction with diffusion priors. In *CVPR*, 2024. 2
- [76] Yitong Xia, Hao Tang, Radu Timofte, and Luc Van Gool. Sinerf: Sinusoidal neural radiance fields for joint pose estimation and scene reconstruction. *CoRR*, abs/2210.04553, 2022. 2
- [77] Junyuan Xie, Ross Girshick, and Ali Farhadi. Deep3d: Fully automatic 2d-to-3d video conversion with deep convolutional neural networks. In *ECCV*, pages 842–857. Springer, 2016. 2, 7
- [78] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv:2406.09414*, 2024. 1
- [79] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *ECCV*, 2018. 2
- [80] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 2
- [81] Xuanwu Yin Yuda Song, Zehao Sun. Sdxs: Real-time one-step latent diffusion models with image conditions. *arxiv*, 2024. 3
- [82] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *ICCV*, 2023. 2
- [83] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 3
- [84] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 7
- [85] Yifan Zhang and Bryan Hooi. Hipa: Enabling one-step text-to-image diffusion models via high-frequency-promoting adaptation. *ArXiv*, abs/2311.18158, 2023. 3
- [86] Zhengyou Zhang. *Eight-Point Algorithm*, pages 370–371. Springer International Publishing, Cham, 2021. 6
- [87] Zhixing Zhang, Bichen Wu, Xiaoyan Wang, Yaqiao Luo, Luxin Zhang, Yanan Zhao, Peter Vajda, Dimitris Metaxas, and Licheng Yu. Avid: Any-length video inpainting with diffusion model. *arXiv preprint arXiv:2312.03816*, 2023. 2