

NormalCrafter: Learning Temporally Consistent Normals from Video Diffusion Priors

Yanrui Bin¹ Wenbo Hu^{2*} Haoyuan Wang³ Xinya Chen⁴ Bing Wang^{1†}

¹Spatial Intelligence Group, The Hong Kong Polytechnic University ²ARC Lab, Tencent PCG

³City University of Hong Kong ⁴Huazhong University of Science and Technology

<https://normalcrafter.github.io/>

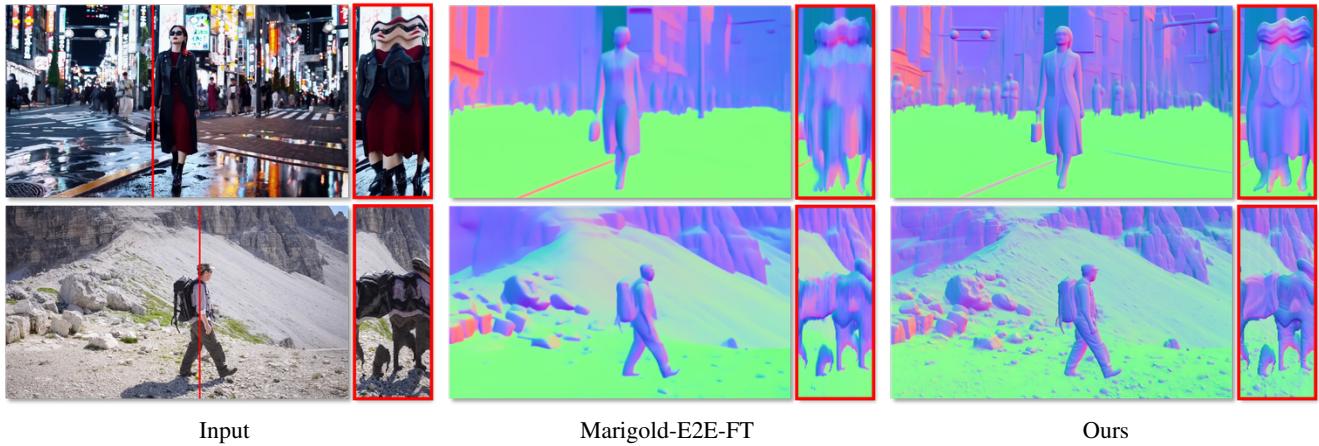


Figure 1. We innovate NormalCrafter, a novel video normal estimation model, that can generate temporally consistent normal sequences with fine-grained details from open-world videos with arbitrary lengths. Compared to results from state-of-the-art image normal estimators, Marigold-E2E-FT [26], our results exhibit both higher spatial fidelity and temporal consistency, as shown in the frame visualizations and temporal profiles (marked by the red lines and rectangles).

Abstract

Surface normal estimation serves as a cornerstone for a spectrum of computer vision applications. While numerous efforts have been devoted to static image scenarios, ensuring temporal coherence in video-based normal estimation remains a formidable challenge. Instead of merely augmenting existing methods with temporal components, we present NormalCrafter to leverage the inherent temporal priors of video diffusion models. To secure high-fidelity normal estimation across sequences, we propose Semantic Feature Regularization (SFR), which aligns diffusion features with semantic cues, encouraging the model to concentrate on the intrinsic semantics of the scene. Moreover, we introduce a two-stage training protocol that leverages both latent and pixel space learning to preserve spatial accuracy while maintaining long temporal context. Extensive evaluations demonstrate the efficacy of our method, showcasing a superior performance in generating temporally consistent

normal sequences with intricate details from diverse videos.

1. Introduction

Surface normals, as pivotal descriptors of 3D scene geometry, underpin a spectrum of applications, including 3D reconstruction, relighting, video editing, and mixed reality. Estimating high-fidelity and temporally consistent normals from diverse, unconstrained videos remains a formidable challenge, owing to variations in scene layouts, illuminations, camera motions, and scene dynamics.

Recent advancements in normal estimation from monocular images have embraced both discriminative [2, 3, 11, 12, 20, 35] and generative paradigms [15, 16, 26, 36, 37]. While discriminative approaches remain hampered by the limitations of training data scale and quality, resulting in suboptimal zero-shot generalization, generative methods harness pre-trained diffusion priors to deliver state-of-the-art performance on open-world images, even when confined to syn-

* Project leader. † Corresponding author.

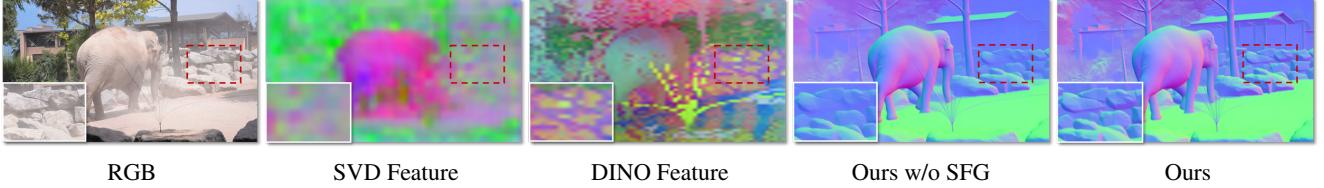


Figure 2. Naively repurposing video diffusion models, *e.g.* SVD [4], for normal estimation (Ours w/o SFG) produces over-smoothed predictions, due to insufficient high-level semantic cues in SVD features. By leveraging Semantic Feature Regularization (SFR) to align diffusion features with DINO [8], our approach yields sharper and more fine-grained normal predictions.

thetic training data. However, these methods are inherently designed for static imagery, neglecting the temporal dynamics of videos and consequently inducing temporal inconsistency or flickering, as demonstrated in Fig. 1.

In this paper, we propose NormalCrafter, a novel video normal estimation model that generates temporally consistent normal sequences exhibiting rich, fine-grained details from unconstrained open-world videos of arbitrary lengths. Rather than incrementally incorporating temporal layers or devising complex stabilization schemes for image-based normal estimators, we harness the potential of video diffusion models for a more robust approach to video normal estimation. Although repurposed video diffusion models have achieved remarkable success in depth estimation, normal estimation presents its own set of challenges, particularly in preserving the high-frequency, semantics-driven details inherent in surface normals. Naively applying video diffusion models to normal estimation often yields suboptimal performance, such as the over-smoothing of normal predictions, as illustrated in Fig. 2. To this end, we introduce a *semantic feature regularization* (SFR) technique that directs the model’s focus toward the semantics by aligning diffusion features with semantic representations extracted from an external encoder, *e.g.*, DINO [8]. Furthermore, recent findings demonstrate that supervising the final output of the variational autoencoder (VAE) in image-based depth or normal estimation, rather than operating solely in the latent space, significantly enhances spatial fidelity. However, this direct supervision considerably increases GPU memory consumption during training, as it requires expanding the compact latent space into the high-dimensional pixel space, thereby restricting training to shorter video clips. To address this issue, we propose a two-stage training strategy: first, training the full model in the latent space to effectively capture long-term temporal context, and then fine-tuning the spatial layers in the pixel space to improve spatial accuracy while preserving the capacity for long sequence inference.

We perform a comprehensive evaluation of our NormalCrafter across a wide range of datasets under zero-shot settings. Both qualitative and quantitative analyses reveal that NormalCrafter attains state-of-the-art performance in open-world video normal estimation, significantly surpassing ex-

isting methodologies. Moreover, our rigorous ablation experiments substantiate the effectiveness of the proposed semantic feature regularization and two-stage training strategy in enhancing both the spatial fidelity and temporal consistency of normal predictions. Our contributions are summarized below:

- We introduce NormalCrafter, a novel framework that generates temporally consistent normal sequences with intricate, fine-grained details for open-world videos of arbitrary lengths, outperforming existing approaches by a substantial margin.
- We propose the semantic feature regularization (SFR) technique, which directs the model’s focus towards meaningful semantics by aligning diffusion features with high-level semantic representations.
- We devise a two-stage training strategy that leverages both latent and pixel-space supervision, enabling the generation of normal sequences with long temporal context while preserving high spatial accuracy.

2. Related Work

Our method relates to two primary research streams: video diffusion models and surface normal estimation. The latter can be categorized into discriminative methods that directly regress normal maps from the input, and more recent diffusion-based approaches that leverage the priors of generative diffusion models for this task.

Discriminative surface normal estimation. Surface normal estimation has been studied for decades. Early work [13, 14, 18] used hand-crafted features with learning-based classification, exemplified by [18] that discretized normals. With deep learning, convolutional neural networks (CNNs) drastically improved this task. Wang *et al.* [35] combined CNNs with vanishing point analysis. Do *et al.* [11] introduced a spatial rectifier to align tilted images with high-likelihood training distributions. Bae *et al.* [3] leveraged aleatoric uncertainty for improved robustness and performance in small structures. Eftekhar *et al.* [12] compiled over 12 million images from diverse scenes and camera intrinsics, training a U-Net on this massive dataset. Its successor Omnidata v2 [20] utilized a transformer-based model with advanced 3D augmentation and cross-task con-

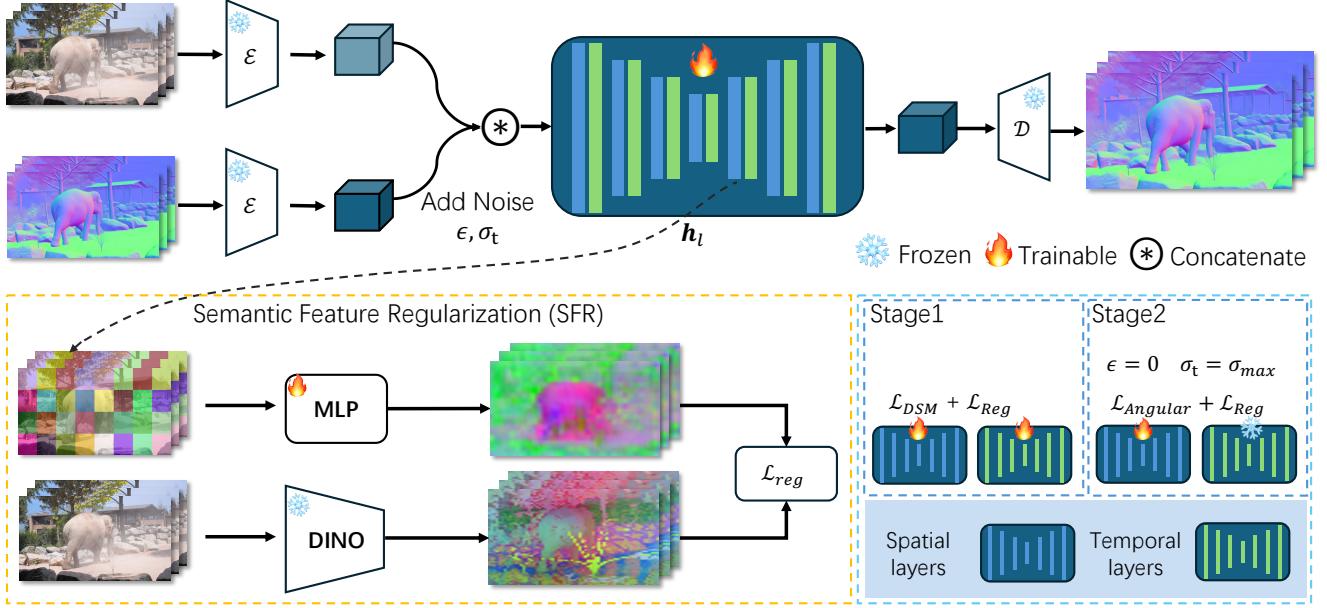


Figure 3. **Overview of our NormalCrafter.** We model the video normal estimation task with a video diffusion model conditioned on input RGB frames. We propose Semantic Feature Regularization (SFR) \mathcal{L}_{reg} to align the diffusion features with robust semantic representations from DINO encoder, encouraging the model to concentrate on the intrinsic semantics for accurate and detailed normal estimation. Our training protocol consists of two stages: 1) training the entire U-Net in the latent space with diffusion score matching \mathcal{L}_{DSM} and SFR \mathcal{L}_{reg} ; 2) fine-tuning only the spatial layers in pixel space with angular loss $\mathcal{L}_{Angular}$ and SFR \mathcal{L}_{reg} .

sistency. Recently, DSINE [2] achieved state-of-the-art results by incorporating per-pixel ray directions and modeling relationships between neighboring normals, providing a strong baseline for our approach.

Diffusion-based surface normal estimation. Recently, pre-trained diffusion models have gained strong attention. Marigold [21] fine-tuned Stable Diffusion (SD) [29] for dense prediction tasks conditioned on images. Concurrently, Geowizard [15] fine-tuned SD to output both depth and normal maps. Although effective, these models required iterative denoising, causing high computational overhead. To address this, some works [26, 36] replaced multi-step denoising with a single-step approach, sacrificing detailed geometry. Addressing this trade-off, Lotus [16] added an image reconstruction objective to enhance details, while StableNormal [37] used a coarse-to-fine scheme for sharper results. Despite their strong priors, these methods overlook temporal context and often produce flickering artifacts in videos. Concurrent with our work, BufferAnytime [23] augmented Marigold-E2E-FT [26] with temporal layers, using optical-flow-based supervision to stabilize results. However, optical flow alone cannot guarantee correct normal correspondences in consecutive frames, as it overlooks camera motion and scene dynamics. In contrast, our approach learns video normal estimation directly from large-scale labeled data and pre-trained diffusion priors, delivering a comprehensive spatio-temporal understanding of the scene. As they neither release the model nor the evalua-

tion data, we exclude it from comparisons.

Video diffusion model. Recent advances in video generation increasingly rely on diffusion models [17, 32, 33] to synthesize temporally coherent frames conditioned on text or images. Latent Diffusion Models (LDMs) [29] offer improved efficiency by operating in a compressed latent space, enabling high-resolution image generation with reduced computational cost. Building on LDMs, Blattmann *et al.* [5] added temporal convolution and attention layers to SD, training these on video data. Stable Video Diffusion (SVD) [4] further refined this approach with extensive training strategy and curated video data. SVD produces high-quality videos and serves as a model prior for diverse video-related tasks [19, 30]. In this paper, we leverage the rich spatio-temporal priors of SVD for high-fidelity, consistent video normal estimation.

3. Method

We present NormalCrafter, a reliable video normal estimator derived from video diffusion models (VDMs). The overall pipeline of NormalCrafter is illustrated in Fig. 3. Given a video $c \in \mathbb{R}^{F \times W \times H \times 3}$ with frame number F , our objective is to generate normal estimations $n \in \mathbb{R}^{F \times W \times H \times 3}$ that are spatially accurate and temporally consistent.

3.1. Normal Estimator with VDMs

To alleviate computational overhead, modern VDMs typically operate in a compressed latent space by leveraging a Variational Autoencoder (VAE) for efficient encoding and decoding of video frames. Since normal maps share the same dimensions as RGB image frames, we seamlessly utilize the same VAE for both the normal maps \mathbf{n} and the corresponding video \mathbf{c} :

$$\mathbf{z}^x = \mathcal{E}(\mathbf{x}), \quad \hat{\mathbf{x}} = \mathcal{D}(\mathbf{z}^x), \quad (1)$$

where \mathcal{E} and \mathcal{D} denote the encoder and decoder of the VAE, respectively, \mathbf{x} may represent either \mathbf{n} or \mathbf{c} , and $\hat{\mathbf{x}}$ is the reconstructed counterpart of \mathbf{x} . However, most existing VAEs are pre-trained on RGB frames, which is suboptimal for normal maps. Therefore, we specifically fine-tune the VAE decoder on normal data to bolster the reconstruction quality.

Diffusion-based normal estimation. In diffusion framework, normal estimation is formulated as a transformation between a simple noise distribution to a target data distribution $p(\mathbf{z}^n | \mathbf{z}^c)$ conditioned on the input video latents \mathbf{z}^c . On the one hand, to map $p(\mathbf{z}^n | \mathbf{z}^c)$ into the noise distribution, a forward diffusion sequence is applied by injecting Gaussian noise with variance σ_t^2 into the latent normal sequence \mathbf{z}_0^n at each time step t :

$$\mathbf{z}_t^n = \mathbf{z}_0^n + \sigma_t^2 \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (2)$$

where $\mathbf{z}_t^n \sim p(\mathbf{z}^n; \sigma_t)$ denotes the noisy latent normal sequence. When σ_t becomes sufficiently large, the noisy latent distribution $p(\mathbf{z}^n; \sigma_t)$ becomes statistically indistinguishable from a pure Gaussian prior. On the other hand, to transform the noise distribution to $p(\mathbf{z}^n | \mathbf{z}^c)$, a reverse denoising process begins by drawing a noise sample $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_{\max}^2 \mathbf{I})$ and iteratively transforms it into \mathbf{z}_0^n through a learned denoiser D_θ . This denoiser is trained via denoising score matching (DSM) [4]:

$$\mathcal{L}_{\text{DSM}} := \mathbb{E}_{\mathbf{z}^n \sim p(\mathbf{z}^n; \sigma_t), \sigma_t \sim p(\sigma)} \lambda(\sigma_t) \| D_\theta(\mathbf{z}_t^n; \sigma_t; \mathbf{z}^c) - \mathbf{z}_0^n \|_2^2, \quad (3)$$

where $p(\sigma)$ is the noise level distribution during training, and $\lambda(\sigma_t) = (1 + \sigma_t^2) \sigma_t^{-2}$ is the weight function. The denoiser function D_θ specifies the noise-level distribution during training. We build our NormalCrafter on top of the SVD model, which is originally designed for generating videos from an input image. We adapt this SVD framework into our NormalCrafter model by substituting the image input with a frame-wise concatenation of the noisy normal latent \mathbf{z}_t^n and the conditional video latent \mathbf{z}^c , as shown in Fig. 3.

3.2. Semantic Feature Regularization

SVD was originally designed for conditioning on a single input image, and may therefore struggle to effectively accumulate contextual information when extended to sequences

of multiple frames. As illustrated in Fig. 2, the initial SVD intermediate features exhibit semantic ambiguity; for instance, the stone region in the background is over-blurred, contradicting the detailed geometry evident in the original frames. As a result, leveraging SVD directly leads to over-smooth normal maps, lacking the intricate structural details in the corresponding areas. To verify whether high-level semantic representations can preserve such geometric details, we visualized the features of the DINO encoder [8] by applying PCA [1]. As shown in Fig. 2, the DINO features exhibit a strong correlation with the geometric structures of the input frames, exemplified by their refined representations of both stone and plant regions.

This motivates us to incorporate semantic features into the diffusion model to further elevate the quality of normal estimation. To this end, the most straightforward approach is to augment the diffusion model with DINO features as an additional conditioning. However, such a design leads to substantial computational and memory overheads during training and inference. Therefore, we propose a Semantic Feature Regularization (SFR) method to rectify the semantic ambiguities in SVD features by aligning them with robust semantic representations throughout training, inspired by REPA [38]. This alignment encourages the diffusion model to concentrate on the intrinsic semantics of the input frames, yielding more accurate and finely detailed normal maps. Moreover, SFR introduces overhead solely during training, leaving inference unaltered with no extra costs.

Specifically, as shown in Fig. 3, we initially derive the DINO features $\mathbf{h}_{\text{dino}} = f(\mathbf{c}) \in \mathbb{R}^{N \times D}$ from the input video frames \mathbf{c} , where N and D indicate the number of patches and the embedding dimension, respectively. Then, we extract the intermediate features \mathbf{h}_l from the l -th layer of the diffusion model, and project them into the DINO feature space using a learnable multilayer perceptron h_ϕ . Finally, we regularize the projected features to align with the DINO features by maximizing the patch-wise cosine similarities:

$$\mathcal{L}_{\text{reg}}(\theta, \phi) := -\mathbb{E}_{\mathbf{c}} \left[\frac{1}{N} \sum_{n=1}^N \text{cossim}(\mathbf{h}_{\text{dino}}^{[n]}, h_\phi(\mathbf{h}_l^{[n]})) \right], \quad (4)$$

where n is the patch index, and cossim is the cosine similarity function between two vectors.

3.3. Two-Stage Training Protocol

Although training NormalCrafter in the latent space with the loss $\mathcal{L}_{\text{DSM}} + \mathcal{L}_{\text{reg}}$ is feasible, it may not yield optimal results in terms of accuracy or efficiency, as highlighted in [26]. Instead, it proposes to fine-tune the image diffusion model in a single end-to-end step for depth and normal estimation, directly optimizing the pixel-wise loss in the image space, thereby achieving superior spatial fidelity alongside improved efficiency. However, extending such an approach

to video normal estimation heavily restricts the length of training clips, since it requires employing VAE to decode the latent normal sequence into pixel space to compute the loss, which drastically elevates memory requirements, especially for long sequences.

To this end, we propose a two-stage training protocol that artfully balances the need for long temporal context modeling with high-precision spatial fidelity. As shown in Fig. 3, we first train NormalCrafter in the latent space under the combined objectives of $\mathcal{L}_{\text{DSM}} + \mathcal{L}_{\text{reg}}$. The sequence length in this stage is randomly sampled from [1, 14], enabling NormalCrafter to flexibly adapt to diverse video durations. Moreover, this setup facilitates training on both single-frame and multi-frame video datasets. In the second stage, we fine-tune only the spatial layers by decoding the latent normal sequence into pixel space and employing the loss $\mathcal{L}_{\text{angular}} + \mathcal{L}_{\text{reg}}$. Here, $\mathcal{L}_{\text{angular}}$ is defined as:

$$\mathcal{L}_{\text{angular}} = \frac{1}{HW} \sum_{i,j} \arccos \left(\frac{\mathbf{n}_{i,j}^* \cdot \hat{\mathbf{n}}_{i,j}}{\|\mathbf{n}_{i,j}^*\| \|\hat{\mathbf{n}}_{i,j}\|} \right), \quad (5)$$

where $\mathbf{n}_{i,j}^*$ is the ground-truth normal at pixel (i, j) , and $\hat{\mathbf{n}}_{i,j}$ is the predicted normal. During this second stage, the sequence length is randomly sampled from [1, 4] frames, thereby easing GPU memory constraints. Since the model has already absorbed long-range temporal cues in the first stage, and only the spatial layers are refined in the second, this two-stage protocol allows the model to enjoy the benefits of end-to-end fine-tuning while preserving its capacity to process extensive sequences.

4. Experiment

4.1. Experimental Setup

Implementation details. We build our NormalCrafter upon the SVD [4] model. For the SFR, we resize the input images to make the DINO feature match the size of the U-Net intermediate features. h_ϕ is a three-layer perceptron while h_l is the output features of the second up blocks of U-Net’s decoder. We fine-tune the VAE for 20,000 iterations employing a base learning rate of 1×10^{-5} . For the U-Net, we train the first stage for 20,000 iterations using a learning rate of 3×10^{-5} and subsequently conduct end-to-end fine-tuning for 10,000 iterations with 1×10^{-5} learning rate at the second stage. In the first stage, we use a hybrid approach: with probability 0.5, we set the noise level σ_t to a fixed value of 700; otherwise, we sample σ_t from a noise level distribution $p(\sigma) = \mathcal{N}(0.7, 1.6)$ following SVD. In both stages, we resize the short edge of input clips to 576 without changing the aspect ratio. All training processes utilize the AdamW optimizer with an exponential learning rate decay strategies following a 100-step warm-up. We conduct training on eight GPUs with a total batch size of

eight. The U-Net training spans approximately 1.5 days, while VAE fine-tuning requires about one day.

Training datasets. Following [37], we train our model using five meticulously selected datasets, encompassing both single-frame and video types, each with high-resolution frames and ground-truth normal maps from synthetic environments. For single-frame datasets, we utilize 49,494 images from Replica [34] for indoor scenes and 45,620 frames from 3D Ken Burns [27] for outdoor scenes. For video datasets, we employ Hypersim [28], MatrixCity [24], and Objaverse [10] to cover indoor scenes, outdoor scenes, and object sequences, respectively. For Hypersim, we utilize the training subset and chain frames from each scene in sequence, yielding 613 videos. We further segment them into 1,780 short clips, each containing between 30 and 60 frames for balanced sampling during training. For MatrixCity, we draw on the training subset of the Big City scene, restructuring frames based on camera extrinsic to produce 2,316 videos, which will then be further divided into 7601 short clips. The normal maps are generated from ground-truth depth maps using cross-product-based methods [2]. For Objaverse, we render 45,081 objects under randomly sampled continuous camera trajectories with diverse lighting to form 45,081 videos. During training, these datasets are sampled proportionally to the number of frames in order to balance the overall training process.

4.2. Evaluations

Evaluation protocols. We thoroughly evaluate NormalCrafter on four widely recognized benchmarks: NYUv2 [31], iBims-1 [22], ScanNet [9], and Sintel [6]. Among these benchmarks, NYUv2 and iBims-1 cater to single-image normal estimation, whereas ScanNet and Sintel contain video sequences. For Sintel, we adopt the consecutive-frame split from DSINE [2] to assess temporal consistency across 1064 frames from 23 scenes. For ScanNet, we sample 20 different scenes, each providing 50 continuous frames for thorough evaluation. We adhere to the DSINE [2] evaluation protocols, computing angular deviations (measured in degrees) between the estimated normal maps and their ground-truth counterparts. We compare mean and median angular errors, where lower values indicate superior performance, and the proportion of pixels with angular errors below certain thresholds (i.e., 11.25° , 22.5° and 30°), where higher values reflect greater precision.

Baselines. We comprehensively evaluate the performance of NormalCrafter against six representative baselines: DSINE [2], GeoWizard [15], GenPercept [36], StableNormal [37], Marigold-E2E-FT [26], and Lotus-D [16]. DSINE stands as the leading method among all discriminative approaches, whereas StableNormal, Marigold-E2E-FT, and Lotus-D establish the frontier among diffusion-based solutions. All of these baselines are devised primarily for

Table 1. **Quantitative evaluations.** The top section shows the results on single-image benchmarks, while the bottom section shows the results on video benchmarks. “mean” and “med” denote the mean and median angular error, respectively. The last column shows the average ranking across all metrics. The best, 2nd-best, and 3rd-best results are highlighted.

Method	NYUv2 [31] (Single-image Benchmark)						iBims [22] (Single-image Benchmark)					
	mean ↓	med ↓	11.25° ↑	22.5° ↑	30° ↑	Rank ↓	mean ↓	med ↓	11.25° ↑	22.5° ↑	30° ↑	Rank ↓
DSINE [2]	16.4	8.4	59.6	77.7	83.5	4.8	17.1	6.1	67.4	79.0	82.3	5.0
GeoWizard [15]	18.6	12.0	46.4	76.1	83.0	7.0	20.5	10.9	51.5	75.2	80.1	7.0
GenPercept [36]	16.4	8.0	60.9	78.3	83.7	3.2	16.3	6.3	69.5	81.1	84.1	2.4
StableNormal [37]	17.7	10.3	54.2	78.1	84.1	4.6	17.0	7.0	68.0	80.9	84.2	3.4
Lotus-D [16]	16.2	8.4	59.8	78.0	83.9	3.4	17.1	6.8	66.4	79.4	83.0	5.2
Marigold-E2E-FT [26]	16.2	7.6	61.4	77.9	83.5	2.8	15.8	5.5	69.9	80.6	83.9	1.8
Ours	15.4	7.9	61.4	79.4	85.1	1.2	16.1	5.9	68.9	80.4	83.7	3.0

Method	ScanNet [9] (Video Benchmark)						Sintel [6] (Video Benchmark)					
	mean ↓	med ↓	11.25° ↑	22.5° ↑	30° ↑	Rank ↓	mean ↓	med ↓	11.25° ↑	22.5° ↑	30° ↑	Rank ↓
DSINE [2]	15.5	8.0	62.4	79.5	84.9	5.4	34.9	28.1	21.5	41.5	52.7	4.6
GeoWizard [15]	18.9	13.1	41.7	75.1	83.0	7.0	37.6	32.0	11.7	32.8	46.8	6.4
GenPercept [36]	14.5	7.2	66.0	81.8	86.7	3.4	34.6	26.2	18.4	43.8	55.8	3.6
StableNormal [37]	15.9	10.0	57.0	81.9	87.0	4.4	38.8	32.7	17.9	36.1	46.6	6.6
Lotus-D [16]	14.3	7.1	65.6	81.4	86.5	3.8	32.3	25.5	22.4	44.9	57.0	2.0
Marigold-E2E-FT [26]	14.1	6.3	67.6	81.7	86.4	2.6	33.5	27.0	21.5	43.0	54.3	3.6
Ours	13.3	6.8	67.4	82.9	87.9	1.4	30.7	23.9	23.5	47.5	60.1	1.0

single-image normal estimation.

Quantitative comparison. We first quantitatively compare our model with baseline normal estimators on both single-image benchmarks (NYUv2 and iBims) and video benchmarks (ScanNet and Sintel) in Tab. 1. We can observe that our NormalCrafter achieves state-of-the-art performance on all video datasets, surpassing existing approaches by a considerable margin. Particularly on the Sintel dataset, characterized by its substantial camera motion and fast-moving objects, NormalCrafter outperforms the second-best method across all metrics, most notably improving mean angular error (1.6°), median angular error (1.6°), and the proportion of pixels with angular errors below 22.5° (2.6) and 30° (3.1). Moreover, on the ScanNet dataset, despite its limited camera movement and static scenes, NormalCrafter still attains the highest performance. Compared with the second-best method, Marigold-E2E-FT [26], NormalCrafter yields a 0.8° improvement in mean angular error, alongside enhancements of 1.2 and 1.5 in the proportions of pixels with angular errors below 22.5° and 30° , respectively, while delivering comparable performance on median angular error and angular errors under 11.25° . The superior performance of NormalCrafter can be attributed to our model’s ability to effectively capture temporal context and SFR to extract intrinsic semantics.

Although our model is primarily designed for video normal estimation, it can also perform single-image normal estimation by setting the frame length to one. As shown in Tab. 1, NormalCrafter demonstrates either state-of-the-art or competitive performance on image-based datasets,

outperforming the second-best method on the NYUv2 dataset in terms of mean angular error (0.8°), as well as the proportions of pixels with angular errors below 22.5° (1.5) and 30° (1.6). On the iBims dataset, our method remains on par with other single-image normal estimation approaches. These results demonstrate the adaptability and robust performance of NormalCrafter, as it can effectively address both video and single-image normal estimation tasks.

Qualitative results. To qualitatively evaluate the performance of NormalCrafter, we compare it with StableNormal and Marigold-E2E-FT on the DAVIS dataset [7] and Sora-generated videos [25], as illustrated in Fig. 4. StableNormal is designed for robust single-image normal estimation, while Marigold-E2E-FT represents a cutting-edge normal estimator. To more vividly illustrate the temporal consistency of the results, we profile the y-t slices for each output within red boxes, obtained by extracting normal values along the temporal axis at designated red line positions, following [19]. We can observe that NormalCrafter consistently yields temporally coherent normal sequences, as evidenced by the smooth y-t slices in all examined examples, whereas both StableNormal and Marigold-E2E-FT exhibit zigzag patterns, indicating flickering artifacts in their estimations. Moreover, NormalCrafter’s predictions exhibit finer-grained details compared to those of StableNormal and Marigold-E2E-FT, thanks to the SFR, which accentuates fine-grained details. More qualitative results are provided in the supplementary material.

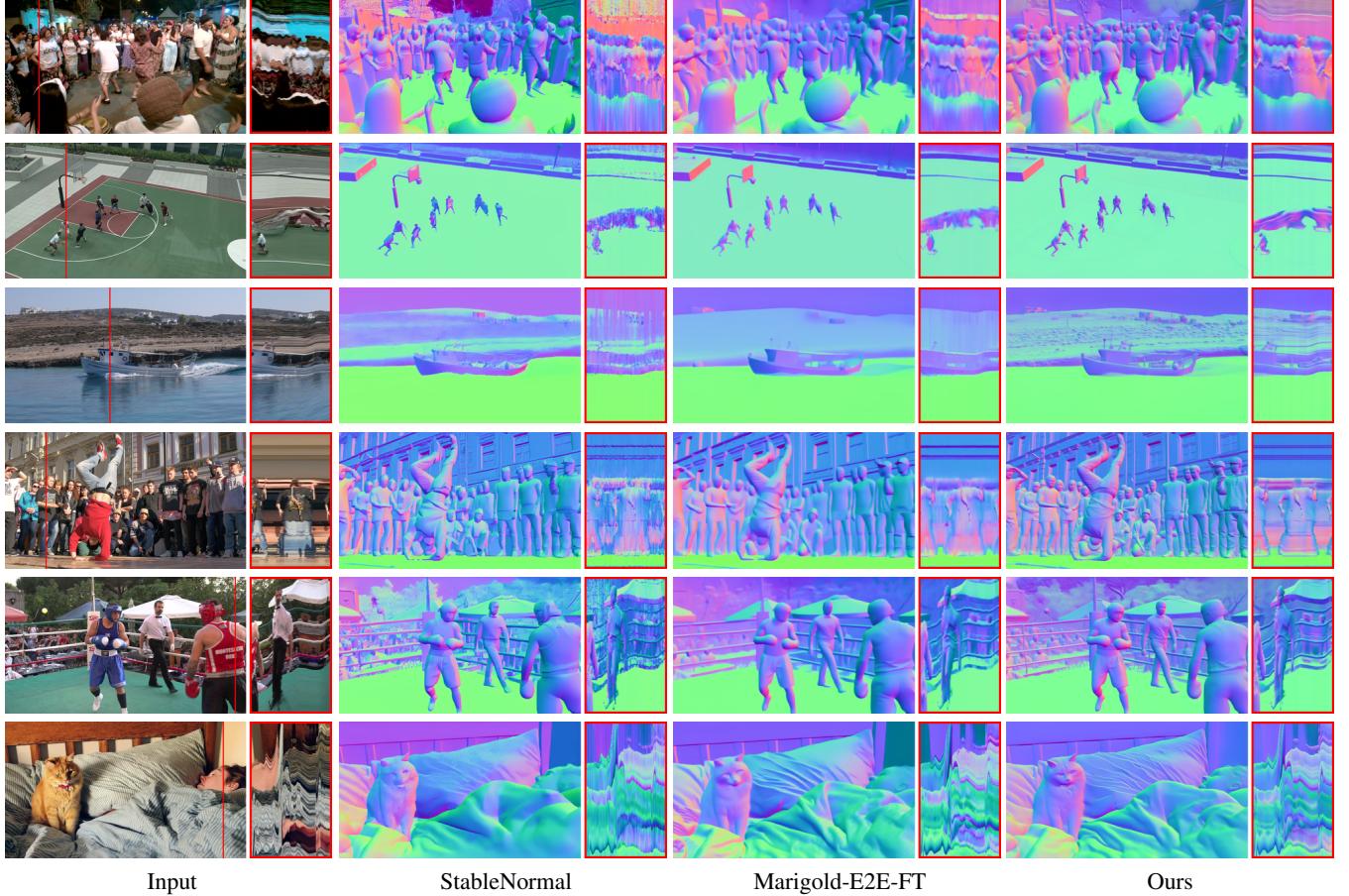


Figure 4. **Qualitative comparisons.** The input videos are sampled from the DAVIS dataset [7] and Sora-generated videos. To highlight the temporal consistency, the y-t slices at the designated red line positions are displayed in red boxes.

Table 2. **Ablation study.** We ablate the effectiveness of Semantic Feature Regularization (SFR), Two-Stage Training strategy (w/o Stage1 and w/o Stage2), and fine-tuning VAE decoder (VAE-FT).

Method	Scannet [9] (Video Benchmark)						Sintel [6] (Video Benchmark)					
	mean ↓	med ↓	11.25° ↑	22.5° ↑	30° ↑	Rank ↓	mean ↓	med ↓	11.25° ↑	22.5° ↑	30° ↑	Rank ↓
Ours w/o VAE-FT	13.4	6.8	67.3	82.8	87.8	1.8	30.8	23.9	23.4	47.4	60.1	1.8
Ours w/o Stage1	13.4	6.8	67.3	82.8	87.6	2.0	30.7	24.2	21.5	46.9	60.0	2.6
Ours w/o Stage2	14.2	8.1	63.7	82.0	87.4	4.8	31.6	25.3	19.7	44.6	57.9	5.0
Ours w/o SFR	13.7	7.0	67.1	82.5	87.4	4.0	31.1	24.7	21.4	45.8	58.9	4.0
Ours	13.3	6.8	67.4	82.9	87.9	1.0	30.7	23.9	23.5	47.5	60.1	1.0

4.3. Ablation study

Effectiveness of Semantic Feature Regularization (SFR).

We compare the performance of NormalCrafter with and without SFR. As shown in Tab. 2, NormalCrafter consistently outperforms the variant without SFR across all metrics on both ScanNet and Sintel datasets. The qualitative comparison in Fig. 5 further illustrates the benefits of SFR, demonstrating SFR’s capability to direct the diffusion model concentrate on intrinsic semantics, thereby enabling accurate and detailed normal predictions.

Influence of SFR location. The U-Net consists of four encoder blocks (“Down0-3”), one middle block (“Mid”), and four decoder blocks (“Up0-3”). We investigate the impact of SFR location by applying SFR at different layers, from “Down1” to “Up2”. As shown in Tab. 3, the performance improvement peaks at “Up1”, indicating that the optimal location for SFR is in the middle of the network. We suspect this is because shallow layers primarily capture low-level information, while deeper layers have too few subsequent layers to effectively map semantics to normal maps.

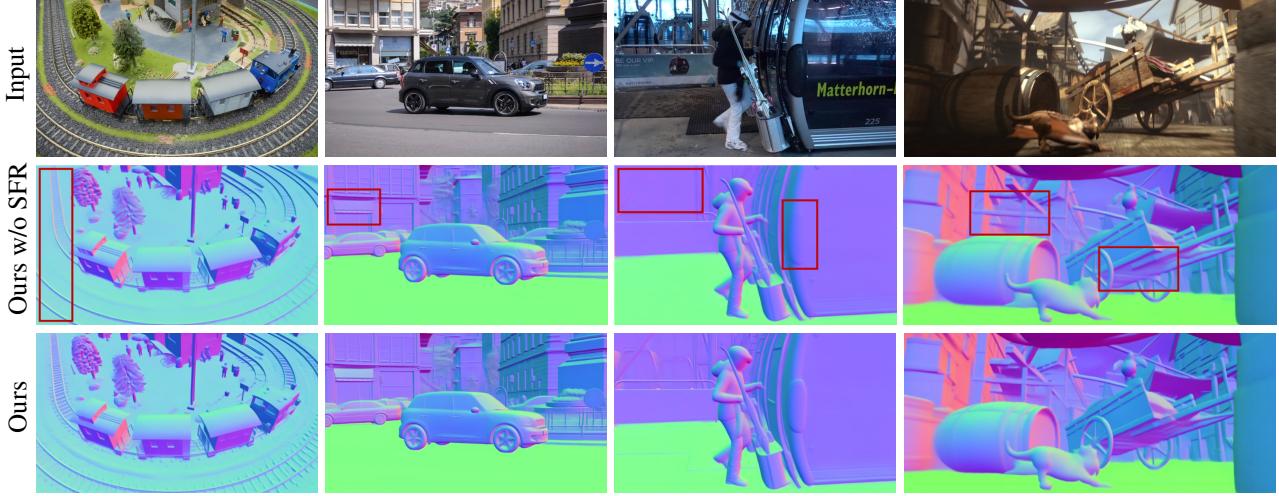


Figure 5. **Ablation results** with Semantic Feature Regularization (SFR). Red boxes highlight the significant differences.

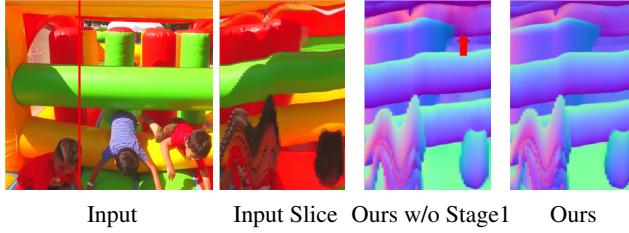


Figure 6. **Qualitative Ablation Results** of two-stage fine-tuning strategy. Without Stage1, the model suffers from temporal consistency due to the limited number of frames in training.

Effectiveness of two-stage training. We ablate the effectiveness of the two-stage training strategy by training the model using stage 1 (w/o stage 2) or stage 2 only (w/o stage 1). From Tab. 2, the model without stage2 (w/o Stage2) performs significantly worse. On the other hand, although the model without stage1 (w/o Stage1) performs comparably with ours in spatial accuracy, it falls short in temporal consistency as shown in Fig. 6. The above observation demonstrates that the two-stage training strategy significantly improves spatial accuracy without compromising temporal consistency. More qualitative comparisons are provided in the supplementary material.

Effectiveness of fine-tuning VAE. We evaluate the effectiveness of fine-tuning the VAE decoder. The reconstruction error of VAE decreases after fine-tuning, with mean angular error reducing from 5.75 to **4.07** and PSNR improving from 25.58 to **28.00**. This superior decoder further positively affects the training of the normal estimator. As shown in Tab. 2, the improved performance of **Ours VAE-FT** demonstrates the effectiveness of fine-tuning VAE.

Table 3. **Influence of SFR location.** We apply SFR at different locations in the U-Net architecture, from “Down1” to “Up2”, and analyze its impact on performance.

Method	Scannet [9]					
	mean ↓	med ↓	11.25° ↑	22.5° ↑	30° ↑	Rank↓
w/o SFR	13.7	7.0	67.1	82.5	87.4	7.0
Down1	13.5	6.8	67.6	82.5	87.4	3.8
Down2	13.6	6.8	67.2	82.4	87.3	6.4
Down3	13.5	6.8	67.4	82.5	87.5	4.2
Mid	13.5	6.8	67.5	82.7	87.7	3.4
Up0	13.4	6.8	67.5	82.8	87.8	2.4
Up1(Ours)	13.3	6.8	67.4	82.9	87.9	2.0
Up2	13.4	6.7	67.8	82.9	87.8	1.4

Method	Sintel [6]					
	mean ↓	med ↓	11.25° ↑	22.5° ↑	30° ↑	Rank↓
w/o SFR	31.1	24.7	21.4	45.8	58.9	6.2
Down1	31.0	24.3	21.4	46.6	59.6	3.8
Down2	31.2	24.9	21.2	45.7	58.7	7.8
Down3	31.1	24.5	21.2	46.2	59.3	6.0
Mid	31.0	24.3	21.5	46.6	59.7	3.2
Up0	30.7	23.8	21.5	47.5	60.5	1.4
Up1(Ours)	30.7	23.9	23.5	47.5	60.1	1.4
Up2	31.1	24.3	21.7	46.7	59.5	3.6

4.4. Limitations

Although our method achieves the state-of-the-art performance in terms of spatial accuracy and temporal consistency in video normal estimation, its large parameter size poses challenges for deployment on mobile devices. Therefore, optimizing the model’s efficiency through model pruning, model quantization and distillation techniques could be a potential direction for future work.

5. Conclusion

We present NormalCrafter, a video normal estimator that can generate temporally consistent normal sequences with fine-grained details for open-world videos. The temporal consistency is achieved by leveraging video diffusion priors, while the spatial accuracy with details is enhanced by semantic feature regularization. Additionally, a two-stage training strategy further improved spatial accuracy while maintaining long temporal context by leveraging both latent and pixel space learning. Extensive evaluations have demonstrated that NormalCrafter achieves state-of-the-art performance in open-world video normal estimation under zero-shot settings. We hope our work can provide inspiration for future investigations in this domain.

References

- [1] Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, 2(4):433–459, 2010. [4](#)
- [2] Gwangbin Bae and Andrew J Davison. Rethinking inductive biases for surface normal estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9535–9545, 2024. [1](#), [3](#), [5](#), [6](#)
- [3] Gwangbin Bae, Ignas Budvytis, and Roberto Cipolla. Estimating and exploiting the aleatoric uncertainty in surface normal estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13137–13146, 2021. [1](#), [2](#)
- [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](#), [4](#), [5](#)
- [5] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22563–22575, 2023. [3](#)
- [6] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In *European Conf. on Computer Vision (ECCV)*, pages 611–625. Springer-Verlag, 2012. [5](#), [6](#), [7](#), [8](#)
- [7] Sergi Caelles, Jordi Pont-Tuset, Federico Perazzi, Alberto Montes, Kevis-Kokitsi Maninis, and Luc Van Gool. The 2019 davis challenge on vos: Unsupervised multi-object segmentation. *arXiv:1905.00737*, 2019. [6](#), [7](#)
- [8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. [2](#), [4](#)
- [9] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017. [5](#), [6](#), [7](#), [8](#)
- [10] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-xl: A universe of 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023. [5](#)
- [11] Tien Do, Khiem Vuong, Stergios I Roumeliotis, and Hyun Soo Park. Surface normal estimation of tilted images via spatial rectifier. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 265–280. Springer, 2020. [1](#), [2](#)
- [12] 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. [1](#), [2](#)
- [13] David F Fouhey, Abhinav Gupta, and Martial Hebert. Data-driven 3d primitives for single image understanding. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3392–3399, 2013. [2](#)
- [14] David Ford Fouhey, Abhinav Gupta, and Martial Hebert. Unfolding an indoor origami world. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VI 13*, pages 687–702. Springer, 2014. [2](#)
- [15] 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, 2024. [1](#), [3](#), [5](#), [6](#)
- [16] Jing He, Haodong Li, Wei Yin, Yixun Liang, Leheng Li, Kaiqiang Zhou, Hongbo Zhang, Bingbing Liu, and Ying-Cong Chen. Lotus: Diffusion-based visual foundation model for high-quality dense prediction. *arXiv preprint arXiv:2409.18124*, 2024. [1](#), [3](#), [5](#), [6](#)
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. [3](#)
- [18] Derek Hoiem, Alexei A Efros, and Martial Hebert. Automatic photo pop-up. In *ACM SIGGRAPH 2005 Papers*, pages 577–584. 2005. [2](#)
- [19] 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. [3](#), [6](#)
- [20] Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18963–18974, 2022. [1](#), [2](#)
- [21] 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. 3
- [22] Tobias Koch, Lukas Liebel, Friedrich Fraundorfer, and Marco Körner. Evaluation of cnn-based single-image depth estimation methods. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018. 5, 6
- [23] Zhengfei Kuang, Tianyuan Zhang, Kai Zhang, Hao Tan, Sai Bi, Yiwei Hu, Zexiang Xu, Milos Hasan, Gordon Wetzstein, and Fujun Luan. Buffer anytime: Zero-shot video depth and normal from image priors. *arXiv preprint arXiv:2411.17249*, 2024. 3
- [24] Yixuan Li, Lihan Jiang, Lining Xu, Yuanbo Xiangli, Zhenzhi Wang, Dahua Lin, and Bo Dai. Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3205–3215, 2023. 5
- [25] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024. 6
- [26] 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. In *WACV*, 2025. 1, 3, 4, 5, 6
- [27] Simon Niklaus, Long Mai, Jimei Yang, and Feng Liu. 3d ken burns effect from a single image. *ACM Transactions on Graphics*, 38(6):184:1–184:15, 2019. 5
- [28] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M. Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *International Conference on Computer Vision (ICCV) 2021*, 2021. 5
- [29] 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. 3
- [30] 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. *arXiv preprint arXiv:2406.01493*, 2024. 3
- [31] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7–13, 2012, Proceedings, Part V 12*, pages 746–760. Springer, 2012. 5, 6
- [32] 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
- [33] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 3
- [34] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 5
- [35] Rui Wang, David Geraghty, Kevin Matzen, Richard Szeliski, and Jan-Michael Frahm. Vplnet: Deep single view normal estimation with vanishing points and lines. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 689–698, 2020. 1, 2
- [36] 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. 1, 3, 5, 6
- [37] Chongjie Ye, Lingteng Qiu, Xiaodong Gu, Qi Zuo, Yushuang Wu, Zilong Dong, Liefeng Bo, Yuliang Xiu, and Xiaoguang Han. Stablenormal: Reducing diffusion variance for stable and sharp normal. *ACM Transactions on Graphics (TOG)*, 43(6):1–18, 2024. 1, 3, 5, 6
- [38] Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining Xie. Representation alignment for generation: Training diffusion transformers is easier than you think. *arXiv preprint arXiv:2410.06940*, 2024. 4