

Re-Depth Anything: Test-Time Depth Refinement via Self-Supervised Re-lighting

Ananta R. Bhattarai
Bielefeld University

Helge Rhodin
Bielefeld University

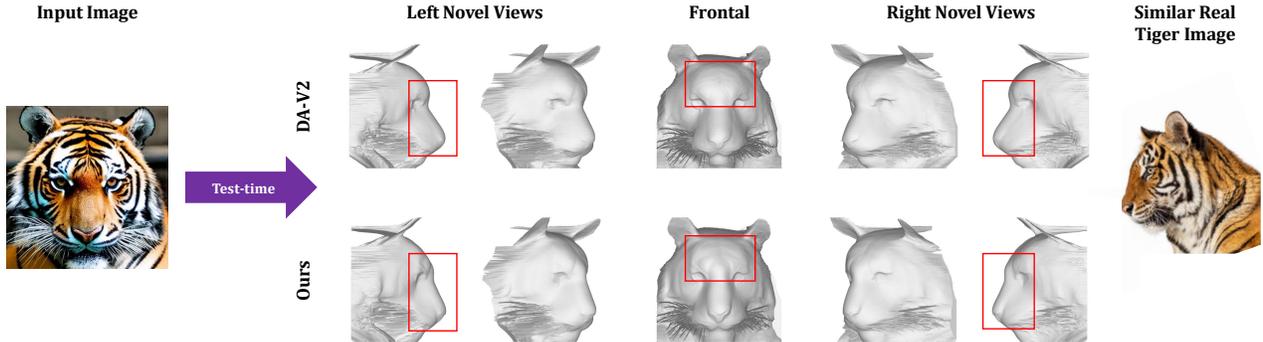


Figure 1. **Re-Depth Anything** refines the prediction of Depth Anything V2 [53] by re-lighting the reconstructed geometry and extracting knowledge from diffusion models in a self-supervised manner. In this example, the test-time optimization enhances facial detail (see frontal view) and refines the nose shape to look more like a tiger (side view), correcting the dog-like initial resemblance likely originating from a biased training distribution. The key contribution is a re-synthesis method that replaces photometric reconstruction for self-supervision.

Abstract

Monocular depth estimation remains challenging as recent foundation models, such as Depth Anything V2 (DA-V2), struggle with real-world images that are far from the training distribution. We introduce Re-Depth Anything, a test-time self-supervision framework that bridges this domain gap by fusing DA-V2 with the powerful priors of large-scale 2D diffusion models. Our method performs label-free refinement directly on the input image by re-lighting predicted depth maps and augmenting the input. This re-synthesis method replaces classical photometric reconstruction by leveraging shape from shading (SfS) cues in a new, generative context with Score Distillation Sampling (SDS). To prevent optimization collapse, our framework employs a targeted optimization strategy: rather than optimizing depth directly or fine-tuning the full model, we freeze the encoder and only update intermediate embeddings while also fine-tuning the decoder. Across diverse benchmarks, Re-Depth Anything yields substantial gains in depth accuracy and realism over the DA-V2, showcasing new avenues for self-supervision by augmenting geometric reasoning.

1. Introduction

Monocular Depth Estimation (MDE) aims to predict dense, per-pixel depth from a single RGB image, enabling numerous applications, including 3D reconstruction [13, 26, 55], autonomous driving [46], robotic navigation [50], and virtual or augmented reality [34]. Lately, high-quality depth maps for diverse images were enabled by foundation models using Vision Transformers (ViTs) [7] with dense prediction heads [33]. Crucially, MiDaS [32] pioneered the training on a myriad of labeled datasets and Depth Anything V2 (DA-V2) [53] showed that the sparse and often noisy depth labels can be enhanced with a teacher model trained on synthetic data. While these foundational models are setting new state-of-the-art performance, inaccuracies on in-the-wild reconstruction remain (see Fig. 1). Monocular depth estimation remains one of the fundamental yet challenging problems in computer vision.

In this work, we introduce Re-Depth Anything, a test-time optimization framework designed to close the domain gap for DA-V2 by self-supervision through 2D generative models. Given a single input image, our framework adapts the pre-trained DA-V2 model to the specific image content.

The core idea is to re-light the DA-V2 depth map in random illumination conditions and to superimpose these onto the input image. The depth map is then refined by using a 2D diffusion model as a prior for scoring how realistic the augmented shading is. This plausibility estimate is backpropagated to the depth map through a differentiable renderer using the Blinn-Phong illumination model [3] and the SDS loss. Crucially, instead of directly optimizing depth or fine-tuning the entire network, we propose a targeted optimization strategy to jointly optimize only the intermediate feature embeddings fed to the Dense Prediction Transformer (DPT) decoder and the decoder’s weights. This targeted approach preserves the strong geometric knowledge encoded in the embeddings while refining the final output.

3D knowledge distillation from 2D image diffusion models [30, 36, 37] using the SDS loss and shading cues was pioneered by DreamFusion [31] for text-to-3D generation. Their key ingredient is to optimize a 3D NeRF representation [26], such that its 2D renderings are perceived as realistic by the diffusion model. This and subsequent works [20, 22, 43, 47, 63] have enabled reconstruction of real images by pairing virtual views with the photometric reconstruction of the real one. However, this line of purely self-supervised learning from geometric relations suffers from cue ambiguities and lags behind supervised models. Our key advance is to apply the benefits of self-supervised learning on top of supervised methods by re-lighting the predicted depth map. This re-synthesizing and augmenting of the input image is fundamentally different to the photometric reconstruction with a full-fledged NeRF [26] or Gaussian Splatting [13] renderer in prior self-supervised work and alleviates the problem of reconstructing appearance pixel-perfect. Our main contributions are:

- We propose Re-Depth Anything, a novel test-time optimization framework that adapts pre-trained DA-V2 to real world images using a 2D diffusion prior on re-synthesized depth predictions, requiring no additional labeled data.
- We propose a single-image re-lighting model that differentially links the predicted depth map to the input image, enabling the use of an SDS loss for self-supervised geometry refinement from a single view.
- We introduce a targeted optimization scheme that jointly optimizes the decoder’s input embeddings and its weights, which we show is crucial for avoiding overfitting and preserving geometric structure.

2. Related Work

Our work, Re-Depth Anything, builds upon progress in three primary research areas: monocular depth estimation, test-time adaptation for the monocular depth estimation task, and the use of 2D diffusion models as priors for 3D reconstruction. Below, we review the most related methods from each of these domains.

Monocular Depth Estimation. Monocular depth estimation has been a long-standing challenge in computer vision. Early approaches [1, 8, 15, 18, 40, 51, 54, 56, 58] relied heavily on supervised learning, training on datasets with ground-truth depth, such as KITTI [9] for outdoor driving and NYU Depth V2 [41] for indoor scenes. More recently, the field has shifted towards building general-purpose foundation models for depth. MiDaS [32] enabled joint training on multiple datasets by predicting disparity instead of depth and by normalizing predictions to unit range, enabling zero-shot generalization to new domains. Even better performance is possible with DPT heads [33] and large diffusion models [12, 61]. DA [52] and its successor, DA-V2 [53], use the same relative prediction and gained further improvements by training on massive-scale datasets of images by aligning the predictions of a teacher model to sparse ground truth measurements. It mitigates but does not resolve the noise in the LiDAR-based ground truth on in-the-wild images.

Another line of research predicts absolute depth without the disparity normalization, from single [4, 29, 44, 57] or multiple images [5, 14, 45]. We focus on relative depth prediction as these models typically lead to higher surface detail, which we aim to improve. Moreover, the absolute depth scale is invariant to shading cues. Hence, we use the recent and popular DA-V2 model as our foundation, aiming to correct its errors rather than retraining it, such that it applies to underrepresented and out-of-distribution inputs.

Test-Time Adaptation for MDE. Test-Time Adaptation (TTA) or Test-Time Optimization (TTO) aims to adapt a pre-trained model to a specific test input at inference time.

In the context of MDE, TTA often relies on self-supervision signals available from the input itself. For video inputs, temporal and photometric consistency between frames is a powerful signal used to fine-tune a depth network [16, 19, 28].

However, these methods are not applicable to the single-image setting, which is more challenging due to the scarcity of self-supervision cues. Some single-image TTA methods adapt a relative depth model to predict metric depth using strong external priors like 3D human meshes [62] or using sparse 3D points from an external source [23]. In contrast, our method requires no such specific external data, and instead leverages a general-purpose 2D diffusion model to refine relative depth for any arbitrary scene. Crucially, instead of relying solely on the input image’s internal consistency, we introduce a powerful prior to provide a dense, geometry-aware supervisory signal for adaptation.

2D Diffusion Models as Priors for 3D 2D image diffusion models [30, 36, 37], trained on internet-scale image and text data, have learned incredibly rich priors about the

visual world. A recent line of work has focused on leveraging these 2D priors for 3D tasks. The most influential is DreamFusion [31], which introduced the SDS loss, enabling the use of a pre-trained text-to-image diffusion model as a loss function to optimize a 3D NeRF representation from scratch using only a text prompt. This concept has been extended and improved by numerous follow-ups [22, 25, 43, 43, 63], including using mesh representations [20, 47] and Gaussian Splatting [6, 42] as differentiable renderers. Reconstructing a real input image brings additional challenges. RealFusion [24] proposes a method to combine real and synthetic views through a diffusion model fine-tuning and carefully selecting compatible virtual views while others directly fine-tune a 2D diffusion model for multi-view reconstruction [21]. Our re-lighting principle is inspired by the DreamTexture [2] variant, which utilizes shape from texture cues through virtual augmentation. However, all of these fully unsupervised methods have not demonstrated advantages over the latest supervised depth estimation techniques.

Our work builds directly on this idea but applies it in a novel context. Instead of generating a 3D shape or optimizing a NeRF on multiple views, we use the SDS loss on a single image to refine the parameters of a pre-trained, feed-forward depth estimation model (DA-V2) using re-lighting as opposed to photometric reconstruction. This test-time optimization adapts the model’s prediction to the specific image content, guided by the diffusion model’s knowledge about the shading of natural objects.

Single-View Geometry and Shading. Shape-from-Shading (SfS) [10, 60] is a classical attempt to recover 3D shape from shading variations in a single 2D image. The idea is to decompose the image into (piecewise-constant) albedo and shading components, e.g. using the diffuse and specular components of a Blinn-Phong [3] model. However, classical SfS [10, 60], and other shape from X methods, such as Shape-from-Texture (SfT) [48, 49] are highly ill-posed, relying on strong assumptions about lighting, texture regularity, and material properties, which are rarely met in practice.

DreamFusion and RealFusion were the first to revisit the SfS principle in a modern context using generative models. Although they could lift some assumptions, RealFusion still follows the classical reconstruction principle of rendering a synthetic object that is optimized to reconstruct the input image. However, real illumination and material properties are complex, leading to artifacts around specular highlights when attempting to match them with simplistic shading models in the deployed differentiable renderers.

By contrast, we do not attempt to solve the full, ill-posed inverse graphics problem. Instead, one of our key contributions is to augment with additional shading cues, which is

achieved with a modification of a simple, lightweight Blinn-Phong [3] renderer. This allows the diffusion prior to critique the plausibility of the 3D shape as expressed through its shading, effectively linking the underlying geometry and image content without photometric reconstruction.

3. Preliminaries and Notation

Score Distillation Sampling Given a 2D image \mathbf{X} , rendered from a differentiable representation with parameters θ , SDS [31] utilizes a pre-trained diffusion model ϕ to optimize θ via gradient descent. Specifically, a gradient toward a more likely image is obtained from the noise ϵ_ϕ predicted by ϕ given a noisy image \mathbf{X}_t , text embedding c , and the noise level t ,

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\mathbf{X}, c) = \mathbb{E}_{\epsilon, t} \left[w(t) (\epsilon_\phi(\mathbf{X}_t; c, t) - \epsilon) \frac{\partial \mathbf{X}}{\partial \theta} \right], \quad (1)$$

where $w(t)$ is a weighting function and $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. In DreamFusion [31], the SDS loss matches the 2D renderings from random angles to the text prompt. In our work, we adapt apply the SDS loss on the re-lit image and adapt it to optimize both the intermediate feature embeddings and decoder weights in DA-V2 to enhance depth prediction.

Depth Anything V2 Architecture DA-V2 [53] follows the recent trend of using ViT encoders as in DINOv2 [27] and DPT [33] for the disparity decoder. Specifically, the encoder transforms input tokens into new feature representations using L transformer layers. Features from four selected layers are extracted and passed to the DPT head for disparity prediction, with layer selection depending on the ViT variant. For example, layers $l = \{3, 6, 9, 12\}$ are used in the ViT-Small configuration. Given an input image \mathbf{I} , we denote the extracted feature representations as embeddings \mathbf{W} . The final disparity $\hat{\mathbf{D}}_{\text{disp}} = f(\mathbf{W}; \theta)$ is predicted by the DPT head $f(\cdot)$, parameterized by pre-trained weights θ and taking embeddings \mathbf{W} as input.

4. Method

Given an input image $\mathbf{I} \in \mathbb{R}^{C \times H \times W}$, our goal is to refine an initial disparity map estimate $\hat{\mathbf{D}}_{\text{init}} \in \mathbb{R}^{H \times W}$ predicted by the pre-trained DA-V2 [53] model. As illustrated in Fig. 2, Re-Depth Anything, is a self-supervised, test-time optimization framework that adapts the DA-V2 model to the specific input image.

Key to our method is how we use the SDS [31] loss as a 2D diffusion prior for 3D refinement. To this end, we first introduce a differentiable rendering function that links the predicted disparity map $\hat{\mathbf{D}}_{\text{disp}}$ to a re-illuminated image $\hat{\mathbf{I}}$ through augmentation. We then use the SDS loss on $\hat{\mathbf{I}}$ to jointly optimize the decoder’s input embeddings \mathbf{W} and the

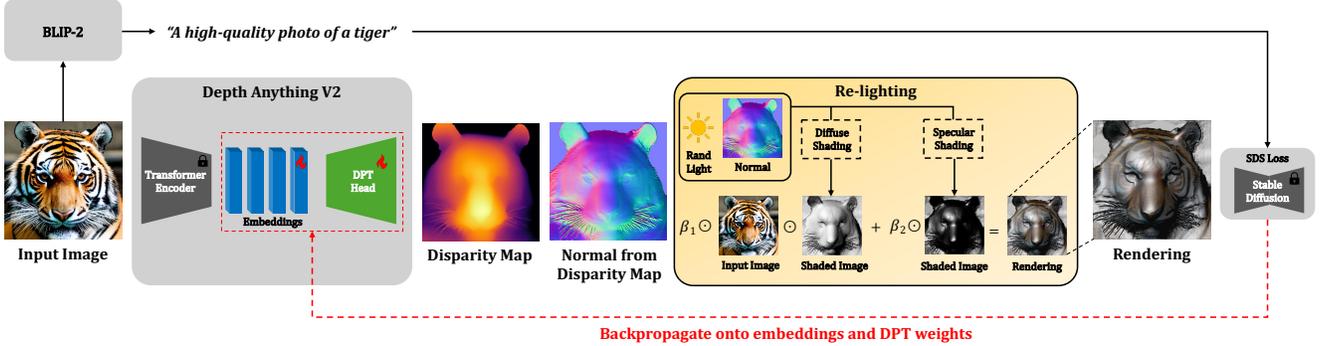


Figure 2. **Re-Depth Anything overview.** Our main contribution is the re-lighting module that randomizes light conditions and shades the estimated geometry on the input. Notably, the re-lighting does not need to look physically accurate as we are only augmenting not photometrically reconstructing the image. Key is also the SDS optimization of embeddings and decoder, while leaving the encoder frozen.

weights θ to better align re-illuminated and original image. We describe each component in detail below.

4.1. Shaded Depth Rendering

To leverage a 2D diffusion prior, we must establish a differentiable link between the depth map $\hat{\mathbf{D}}$ and a 2D image. A full inverse rendering approach, which would decompose the scene into albedo, lighting, materials, and 3D geometry, is highly ill-posed from a single image.

Instead, we propose to augment the input image with additional shading effects by re-lighting. We synthesize diffuse and specular reflectance maps with the classical Blinn-Phong shading model [3], as a function of the normals \mathbf{N} of the depth map. Computing the normals at pixel coordinates u, v requires a camera model. We test a scaled orthographic and a perspective projection with intrinsic matrices \mathbf{K}_{orth} and $\mathbf{K}_{\text{persp}}$, respectively. We then unproject the depth map element-wise into a 3D mesh with vertices

$$\mathbf{X} = \mathbf{K}_{\text{persp}}^{-1} \begin{pmatrix} \mathbf{U}\hat{\mathbf{D}} \\ \mathbf{V}\hat{\mathbf{D}} \\ \hat{\mathbf{D}} \end{pmatrix} \text{ or } \mathbf{X} = \mathbf{K}_{\text{orth}}^{-1} \begin{pmatrix} \mathbf{U} \\ \mathbf{V} \\ \hat{\mathbf{D}} \end{pmatrix}, \quad (2)$$

with \mathbf{U} and \mathbf{V} the tensor of horizontal and vertical pixel coordinate ranging from -1 to 1 . The normal \mathbf{N} is orthogonal to the spatial gradients $(\nabla \mathbf{X}_u, \nabla \mathbf{X}_v)$, computed across the entire image using the element-wise cross product,

$$\mathbf{N} = \frac{\nabla \mathbf{X}_v \times \nabla \mathbf{X}_u}{\|\nabla \mathbf{X}_v \times \nabla \mathbf{X}_u\|_2}. \quad (3)$$

Crucial for our re-lighting of the input image is to use the inverse-tonemapped input image, $\tau^{-1}(\mathbf{I})$, as a proxy for the scene’s diffuse albedo and we assume that specular highlights are colorless, not affected by albedo. While not physically accurate, it exploits that illumination effects are linear and there remains ambiguity between light and surface color, providing a sufficient, differentiable connection between geometry and re-lit appearance. Specifically,

we synthesize a re-illuminated image $\hat{\mathbf{I}} \in \mathbb{R}^{C \times H \times W}$ using Blinn-Phong shading [3],

$$\hat{\mathbf{I}} = \tau(\beta_1 \max(\mathbf{N} \cdot \mathbf{l}, 0) \odot \tau^{-1}(\mathbf{I}) + \beta_2 \max(\mathbf{N} \cdot \mathbf{h}, 0)^\alpha), \quad (4)$$

where $\mathbf{N} \in \mathbb{R}^{3 \times H \times W}$ is the per-pixel normal map derived from the depth gradients of $\hat{\mathbf{D}}$, $\mathbf{l} \in \mathbb{R}^3$ is the light direction, $\mathbf{h} \in \mathbb{R}^3$ is the halfway vector between \mathbf{l} and the view direction $\mathbf{v} = [0, 0, 1]^T$, and α, β_1 , and β_2 are material and light intensity parameters that are determined in the next section. The tone-mapping function $\tau(\mathbf{I}) = \mathbf{I}^{1/\gamma}$, with $\gamma = 2.2$, ensures that shading is performed in the linear RGB color space.

Note that using \mathbf{I} as the albedo can double the shading effects (e.g. existing shadows get darker) and the addition of specular maps may saturate the image, which, however, rarely in our experiments and is similar to specular highlights appearing white in photographs. We keep values within bounds by clamping the rendered output $\hat{\mathbf{I}}$ to the range $[10^{-3}, 1]$.

Handling normalized relative depth. DA-V2 outputs normalized relative depths in the form of the normalized disparity map $\hat{\mathbf{D}}_{\text{disp}} = (1/\hat{\mathbf{D}} - m)s$, where m is the minimum disparity and s is one over the maximum–minimum range of disparity. Converting to absolute depth requires

$$\hat{\mathbf{D}} = \frac{1}{\hat{\mathbf{D}}_{\text{disp}}/s + m} = \frac{s}{\hat{\mathbf{D}}_{\text{disp}} + ms}, \quad (5)$$

where neither the scaling s nor the offset m is known at test time. However, the normal is invariant to the global scale and hence we only have to optimize the unknown scalar $b = ms$ alongside the depth refinement. Notably, we found the optimization to be insensitive to these parameters, likely because shading is about the relative angle between the light and the normal, not absolute orientation. Specifically, it is sufficient to fix $b = 0.1$ with scaled orthographic projection.

4.2. Augmentation Objective

Our goal is to refine the depth map to yield shading effects that yield plausible re-lightings of the input image. This augmentation principle lets us choose random light and material properties instead of having to estimate parameters for the potentially absent or complex shading effects in the input image when doing photometric reconstruction. At each optimization step, we randomly sample the light direction \mathbf{l} and diffuse and specular intensities (β_1, β_2) and exponent α to ensure the refined geometry is consistent with the image across diverse shading conditions.

Loss Function. Plausibility of the augmentation is measured by the total loss combining the generative prior with a smoothness regularizer,

$$\mathcal{L}(\hat{\mathbf{I}}, c, \hat{\mathbf{D}}_{\text{disp}}) = \mathcal{L}_{\text{SDS}}(\hat{\mathbf{I}}, c) + \frac{\lambda_1}{hw} \sum_{i,j} \|\Delta \hat{\mathbf{D}}_{\text{disp}}^{i,j}\|^2, \quad (6)$$

where \mathcal{L}_{SDS} is the SDS loss defined in Eq. (1), and $\hat{\mathbf{I}}$ is the image rendered from $\hat{\mathbf{D}}$ using Eq. (4). The second term is an L1 regularizer on the disparity gradients, which encourages smoother surfaces and prevents noisy artifacts.

To obtain the conditioning text prompt c , we employ BLIP-2 [17], a state-of-the-art image-to-text model, to generate a descriptive caption for the input image \mathbf{I} .

4.3. Optimization Scheme

Our goal is to refine the output depth map of feed-forward depth estimators, specifically the DA-V2 model. However, directly optimizing the depth map on the proposed re-lighting loss remains an ambiguous problem with many plausible solutions. Instead, we propose optimizing the latent feature space and weights of DA-V2, thereby leveraging the prior on 3D shapes learned at training time.

Candidates for optimization are the entire DA-V2 weights or its components. We found that fine-tuning the entire DA-V2 tends to fall into poor local minima or cause the geometry to overfit to image textures. To address this, we jointly optimize only the intermediate embeddings \mathbf{W} (the intermediate feature embeddings of the frozen ViT encoder) and the DPT head’s weights θ .

$$\mathbf{W}^*, \theta^* = \arg \min_{\mathbf{W}, \theta} \mathcal{L}(\hat{\mathbf{I}}, c, \hat{\mathbf{D}}_{\text{disp}}) \quad (7)$$

where $\hat{\mathbf{D}}_{\text{disp}} = f(\mathbf{W}; \theta)$ and $\hat{\mathbf{I}}$ is a function of $\hat{\mathbf{D}}_{\text{disp}}$.

Depth Map Ensembling. The stochastic nature of the SDS loss, primarily due to the random sampling of noise ϵ and timestep t , can lead to high variance in the optimization results. Consequently, disparity predictions can vary noticeably across different runs.

To stabilize the final prediction, inspired by the ensembling in Marigold [12], we perform the optimization N times with a different random seed. We then aggregate the resulting disparity maps by a simple mean operation. The final disparity map \mathbf{D}_{disp} is obtained as

$$\mathbf{D}_{\text{disp}} = \frac{1}{N} \sum_{i=1}^N f(\mathbf{W}_i^*; \theta_i^*), \quad (8)$$

where $(\mathbf{W}_i^*, \theta_i^*)$ are the optimized embeddings and decoder weights from the i -th run.

5. Experiments

We study the effectiveness and accuracy of our self-supervised re-lighting approach on three benchmarks, demonstrating consistent improvements over the DA-V2 baseline across the established metrics, including a relative improvement up to 12.6%. Fig. 3 shows exemplary cases that are improved by removing noise from flat areas and adding missing details. Additional results are shown in the supplemental document.

Baselines. We utilize the small variant of the DA-V2 architecture as our base model, which is one of the most popular monocular relative depth estimation methods. Besides various baseline variants using only parts of our contributions, we also compare against a simple Shape from Shading implementation, to demonstrate the advantage of re-lighting rather than photometric reconstruction.

Implementation Details. For all experiments, input images are first resized, maintaining their original aspect ratio, such that at least one side measures 518 pixels, to match the training resolution of the DA-V2 model. Before applying Stable Diffusion, the images are further zero-padded if they have a non-square aspect ratio and resized to 512×512 .

Our entire pipeline is implemented in PyTorch. For SDS guidance, we employ v1.5 of Stable Diffusion [36]. We optimize the encoder’s embeddings and DPT weights for 1000 iterations using the AdamW optimizer. We set a learning rate of 1×10^{-3} for the embeddings and 2×10^{-6} for the DPT weights. We set the regularization weight λ_1 to 1.0. At each optimization step, we uniformly sample two coefficients $(\beta_1, \beta_2) \sim \mathcal{U}[0, 1]$ and an exponent α , where $\alpha = 2^k$ and $k \sim \mathcal{U}[2, 8]$. β_1 and β_2 are subsequently normalized to ensure their sum is 1.0 (i.e. $\beta_1 + \beta_2 = 1$). Similarly, the light direction vector $\mathbf{l} = (L_x, L_y, L_z)$ is sampled by drawing the X and Y coordinates from a uniform distribution, $L_x, L_y \sim \mathcal{U}[-1, 1]$, while fixing the Z coordinate to $L_z = 1$. The resulting vector \mathbf{l} is then L2-normalized. In the absence of known camera parameters on in-the-wild images, our default is a scaled orthographic camera and we

Table 1. Comparison with DA-V2 across datasets. Relative error reduction of **Ours** over DA-V2 is shown in the last row of each dataset.

Dataset	Method	<i>Higher is better</i> \uparrow			<i>Lower is better</i> \downarrow					
		δ_1	δ_2	δ_3	AbsRel	RMSE	log10	RMSE log	SI log	SqRel
CO3D	DA-V2	1.0	1.0	1.0	0.00227	0.0602	0.000985	0.00321	0.321	0.000244
	Ours	1.0	1.0	1.0	0.00223	0.0588	0.000968	0.00314	0.314	0.000235
	<i>Rel. Δ (%)</i>	-	-	-	1.75	2.26	1.74	2.24	2.24	3.66
KITTI	DA-V2	0.568	0.796	0.902	0.305	7.01	0.118	0.348	33.6	2.49
	Ours	0.593	0.818	0.917	0.283	6.71	0.110	0.319	30.7	2.20
	<i>Rel. Δ (%)</i>	5.73	10.9	15.3	7.10	4.29	6.55	8.51	8.51	11.4
ETH3D	DA-V2	0.884	0.956	0.978	0.113	0.955	0.0448	0.153	15.1	0.391
	Ours	0.898	0.965	0.982	0.104	0.875	0.0413	0.143	14.1	0.347
	<i>Rel. Δ (%)</i>	12.2	21.1	19.5	8.30	8.39	7.72	6.44	6.22	11.1

optimize the scaling starting from seven. For perspective camera, the focal length is initialized to two and refined alongside the disparity optimization. The final prediction is generated by aggregating the results from 10 optimization runs. Each of these 10 runs is initialized from the original pre-trained weights of the DA-V2 model. We conduct our evaluation at the resolution of the initial 518-pixel-side resized input image. One run takes approximately 80 seconds on a single NVIDIA RTX 5000.

Datasets. We evaluate Re-Depth Anything on three standard benchmarking datasets: CO3Dv2 [35] contains multi-view images of several objects across 50 categories, with camera poses, 3D point clouds, foreground masks, and depth maps that we utilize for sparse depth ground truth. From each sequence with a valid depth map, we randomly selected two images. This pre-processing step yielded a total of 80 images from 20 object categories.

KITTI [9] is a large-scale autonomous driving dataset featuring sparse metric depth captured by a LiDAR sensor. We randomly sampled 10 images from each sequence of the official validation set, resulting in a total of 130 images.

ETH3D [39] is a high-resolution benchmark with both indoor and outdoor scenes. We randomly sampled ten raw images and their corresponding depth maps from 13 scenes. This process resulted in a total of 130 images for evaluation.

Evaluation protocol. We apply the widely adopted metrics for assessing the quality of monocular depth estimation: δ_1 , δ_2 , δ_3 , AbsRel, RMSE, log10, RMSE log, SI log and SqRel. We compute these metrics on absolute depth maps, obtained by first finding the least-squares affine fit to the label in disparity, then converting to depth, and subsequently finding the affine fit in depth, as is typical on these benchmarks for relative depth prediction [53].

5.1. Quantitative Evaluation

Table 1 shows that our test-time refinement improves on DA-V2 across CO3D, KITTI, and ETH3D on all nine eval-

uation metrics. Notably, we achieve significant relative reductions in error, including 8.5% in SI log and RMSE log on KITTI, alongside an 8.4% in AbsRel on ETH3D.

On CO3D, the delta metrics are saturated and errors are smaller overall, leading to smaller yet still consistent improvements. This robustness over nine different metrics validates the efficacy of our self-supervised approach, demonstrating considerable potential for fine-tuning foundational models.

Notably, the improvements are consistent across CO3D, which covers single objects in a close-up setting, to street scenes in the KITTI dataset, and indoor scenes in the ETH3D dataset. This highlights the strong generalization capability of our test-time re-lighting approach, inherited from the robustness of generative image models. The nature of the improvement is best explained at examples.

5.2. Qualitative Evaluation

For visual assessment, we present a qualitative comparison against DA-V2 in Fig. 3. Re-Depth Anything produces visibly superior results by enhancing fine-grained details such as the threads on a ball (first image), balcony railings, and electricity wires (second-to-last image), but also removing noise from flat surfaces, as indicated by an arrow in the fourth example. These qualitative improvements are consistent with the quantitative gains reported in Table 1.

We also compare against classical shape from shading, which, however, fails drastically if its strict assumptions are violated. For instance, even on the relatively simple ball example in Fig. 4 (first row), the discoloring of the leather leads to spurious and noisy normals. This highlights the importance of our re-lighting augmentation strategy, which does not make an assumption on albedo constancy and does not suffer from the seam artifacts typically associated with shape from shading.

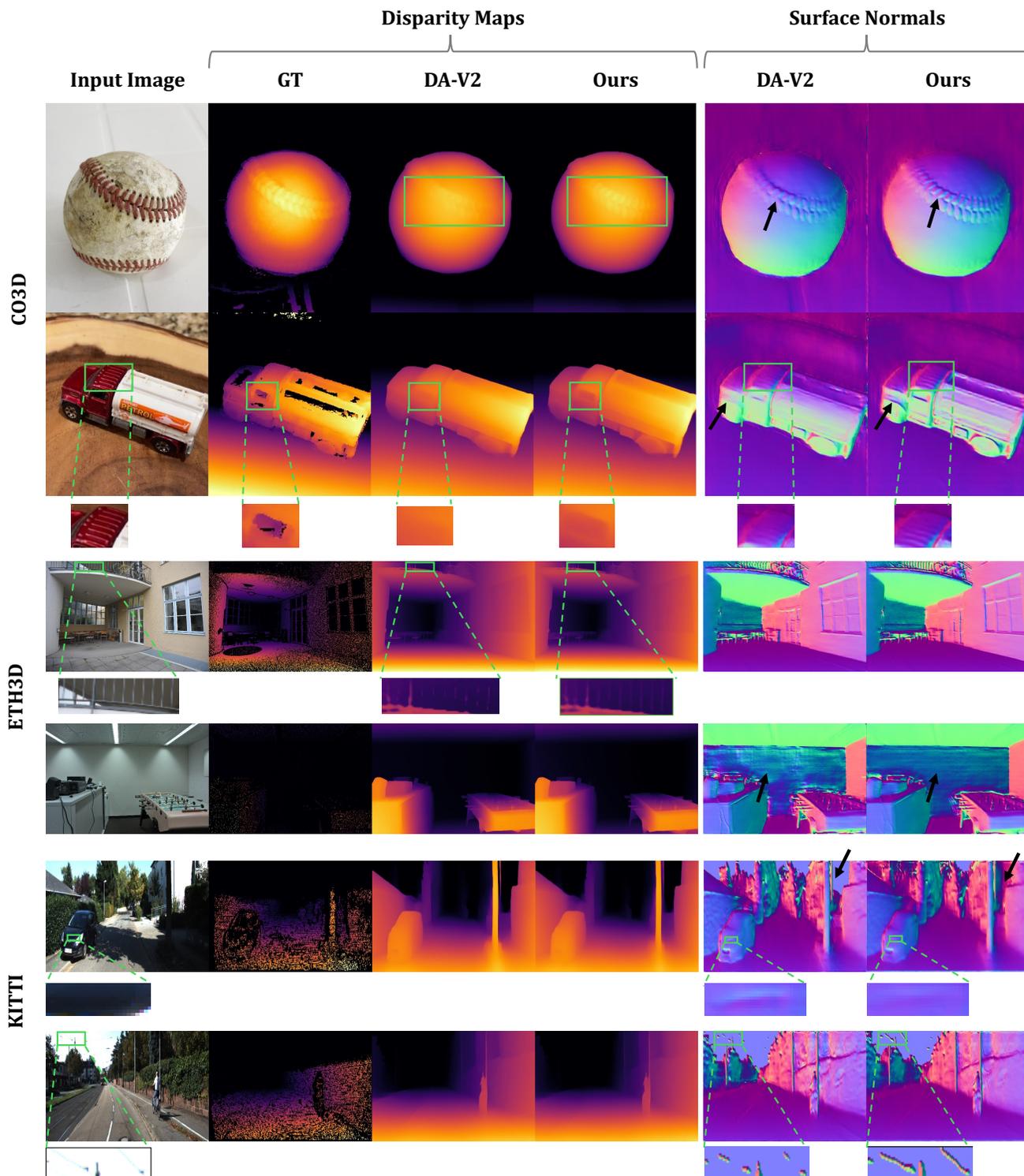


Figure 3. **Qualitative Comparison**, highlighting the added detail (rows 1,2,3,6) and noise-removal effects (rows 4,5).

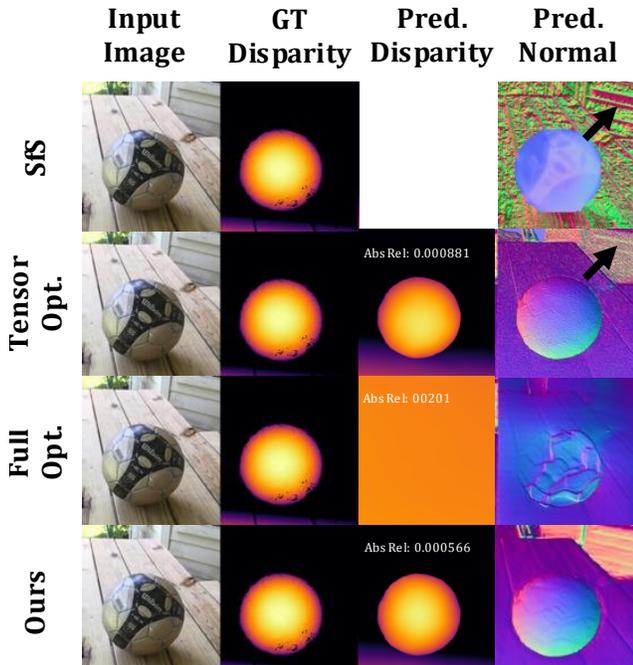


Figure 4. **Qualitative ablation** showing that optimizing depth directly or fine-tuning the whole network at once are detrimental. The listed error values relate to visual and quantitative improvements.

Limitations. We rarely observed small hallucinated edges as for the sticker on the truck in Fig. 3, which is plausible but incorrect. Sometimes the method extends geometry into the sky and over-smooths fine details, such as trees in dark regions, as seen on the KITTI example (see inset Fig. 5), which could potentially be handled by thresholding. In general, we found single objects in the CO3D dataset show stronger detail improvements, while in room and street scenes the largest gains are from removing suspicious details in the initial DA-V2 predictions, which leads to unrealistic re-lighting highlights, and are hence effectively removed by our method while preserving actual detail.

5.3. Ablation study

Optimization We present an ablation study on the CO3D evaluation set, comparing our full pipeline to a baseline lacking the SDS loss, and four optimization variants: (1) direct depth pixel optimization, (2) full DA-V2 fine-tuning, (3) fine-tuning only embeddings and DPT weights, and (4) a two-stage version of (3) that first optimizes the embeddings

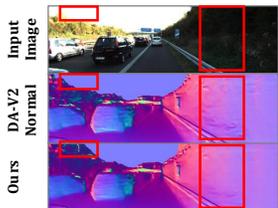


Figure 5. **Limitations.**

Table 2. **Ablation on the CO3D dataset** showing the significance of the major design choices.

Method	AbsRel ↓	RMSE ↓	log10 ↓	RMSE log ↓	SI log ↓	SqRel ↓
w/o \mathcal{L}_{SDS}	0.00427	0.0993	0.00185	0.00532	0.532	0.000661
Tensor Model	0.00226	0.0601	0.000985	0.00321	0.321	0.000244
Full Model	0.00331	0.0779	0.00143	0.00418	0.418	0.000412
Two Stage	0.00225	0.0597	0.000979	0.00319	0.319	0.000241
Ours	0.00223	0.0588	0.000968	0.00314	0.314	0.000235

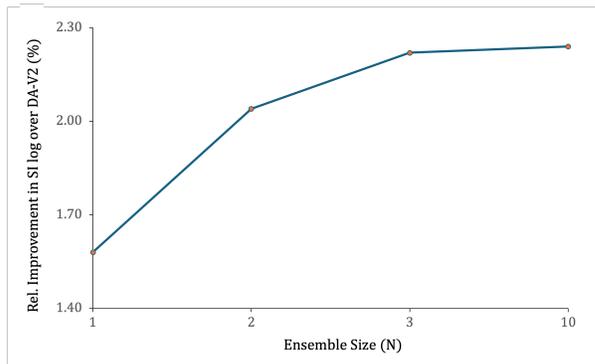


Figure 6. **Ensemble size** vs. improvement in SI log on CO3D.

and subsequently fine-tunes the decoder.

The qualitative results in Fig. 4 show that direct optimization (1) creates noise artifacts (first row), while full finetuning (2) causes the geometry to collapse (second row), even-though we reduced to a very small learning rate of 2×10^{-6} for optimizing both the ViT encoder and the DPT decoder.

Our chosen approach (3) strikes a balance, enhancing detail without degradation. These visual observations are corroborated by the quantitative metrics in Table 2. Although variants (3) and (4) are visually comparable, (3) achieves lower errors, validating our design choice.

Ensembling predictions. We investigate the impact of ensembling predictions via mean aggregation. The results, shown in Fig. 6, demonstrate clear benefits but with rapidly diminishing returns. While a prediction from a single run achieves a 1.58% improvement on SI log over DA-V2, ensembling predictions from 3 runs significantly increases this to 2.22%. This benefit quickly saturates, with 10 predictions offering only a negligible further gain (2.24%).

6. Conclusion

Re-Depth Anything presents a new method for test-time optimization by re-lighting. The key contribution is to use generative models for scoring the shading-image alignment instead of requiring a photometric reconstruction. This alleviates the need to construct a photorealistic renderer for inverse graphics and shows consistent improvements over all tested datasets and metrics. In future work, we plan to explore alternative re-synthesis approaches and explore fine-tuning foundational models at scale on in-the-wild footage.

Re-Depth Anything: Test-Time Depth Refinement via Self-Supervised Re-lighting

Supplementary Material

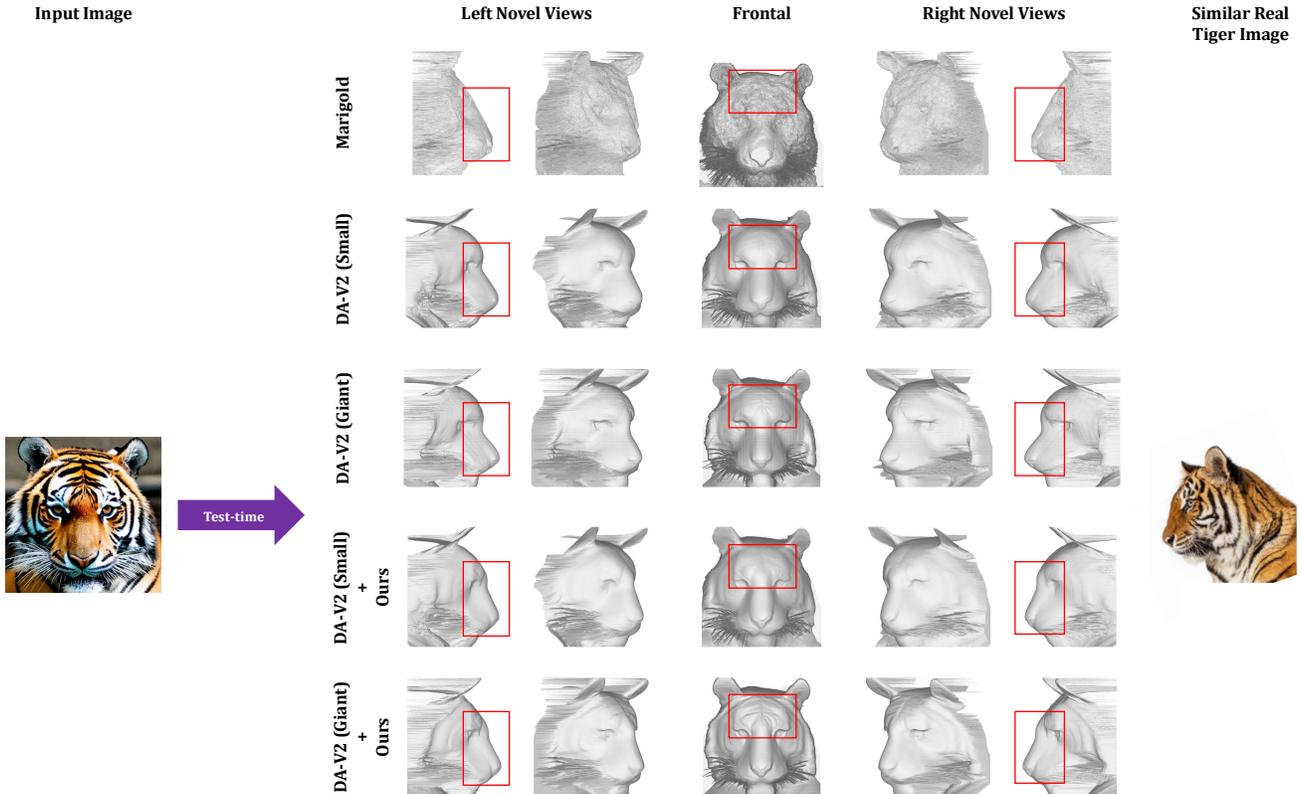


Figure 7. **Biased predictions by baselines and our correction via re-lighting.** **Top:** Marigold and the larger DA-V2 variants struggle with the tiger image in ways similar to the small variant shown in the main paper teaser. **Bottom:** Our method applies to both the large and small variants of DA-V2, correcting the overall shape in each case and adding more details when using the large variant.

This supplemental document provides additional qualitative and quantitative comparisons, justifies the choice of DA-V2 as the baseline, and explains the challenge of comparing state-of-the-art depth estimation models due to mismatching depth representations and evaluation protocols.

7. Additional Comparisons

Fig. 7 exemplifies how different depth estimation methods all struggle with the tiger image from the main paper. The reconstructed dog-like shapes indicate a bias in the training data rather than an issue with the model architecture or depth representation. Our re-lighting is agnostic to the training and corrects the bias from a dog-like to a tiger-

like shape, regardless of whether the giant or small variant of DA-V2 is used. Quantitative improvements over other methods are listed in Table 3 and discussed further below.

We present additional qualitative results in Figs. 9, 10, 11, and 12. These figures also illustrate several remaining limitations of our method, highlighted by red boxes and discussed in the main document. Fig. 8 further shows the per-sample change in SqRel relative to DA-V2, highlighting that our method delivers more frequent and larger-magnitude improvements across all three datasets.

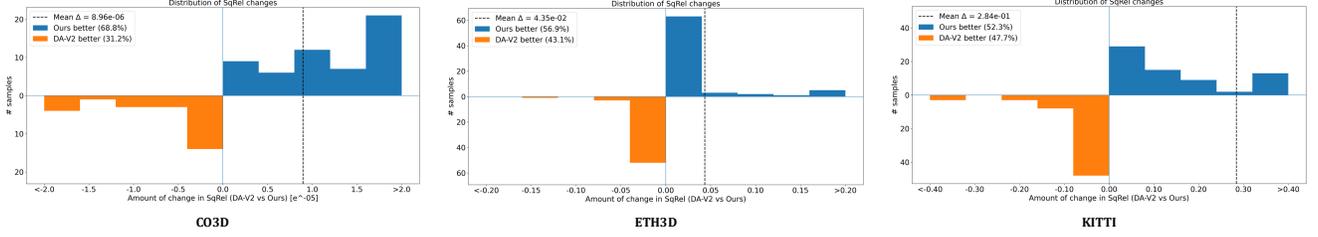


Figure 8. Histograms of per-sample change in SqRel with respect to DA-V2 on samples from CO3D, ETH3D and KITTI. The x-axis shows $\Delta \text{SqRel} = \text{SqRel}_{\text{DA-V2}} - \text{SqRel}_{\text{Ours}}$ (for CO3D scaled by 10^{-5} as indicated on the axis): blue bars on the right correspond to samples where our method achieves lower SqRel than DA-V2, and orange bars on the left to samples where DA-V2 is better. The dashed vertical line marks the mean ΔSqRel for each dataset. On CO3D, the distribution is clearly skewed towards positive changes with relatively few and smaller-magnitude failures. ETH3D and KITTI both exhibit a heavier positive tail. Overall, our method improves SqRel on more samples than it degrades, and the gains are larger in magnitude than the occasional losses, leading to a positive mean ΔSqRel on all three datasets.

8. Results with the DA-V2 Giant Backbone

We used the ViT-S backbone for the main experiments because it performs almost as well as the larger ones while being significantly smaller (in terms of embedding size) and therefore more efficient to optimize. We also found that the giant model with the ViT-G encoder suffers from similar biases, such as the dog-like prediction for the tiger image in Fig. 7. Therefore, we used the smaller model for development and for the major comparisons.

9. Evaluation Protocol Details

Since monocular depth estimation is inherently scale-ambiguous, various depth representations and corresponding evaluation protocols have been proposed to support different invariances. Unfortunately, even among relative depth estimation methods, there are discrepancies in the evaluation protocols that are used. The classical approach is to align the relative depth prediction with the ground-truth depth in a least-squares sense before applying a range of

Table 3. **Quantitative comparison with DA-V2 and Marigold.** DA-V2 outperforms the other baselines on CO3D, and our method further improves its performance. We compared different alignment metrics (see Sec. 9) and highlighted the comparable ones in gray and black, respectively. We advocate using the protocol marked in solid black.

Dataset	Method	Prediction			<i>Higher is better</i> \uparrow			<i>Lower is better</i> \downarrow					
		Rel. Disp.	Rel. Depth	Abs. Depth	δ_1	δ_2	δ_3	AbsRel	RMSE	log10	RMSE log	SI log	SqRel
CO3D	Marigold _{ls-depth}		✓		1.0	1.0	1.0	0.00276	0.0685	0.001200	0.00366	0.366	0.000320
	Marigold _{ls-depth-disp}		✓		1.0	1.0	1.0	0.00275	0.0685	0.001196	0.00366	0.366	0.000320
	DepthPro _{ls-depth}			✓	1.0	1.0	1.0	0.00242	0.0625	0.001053	0.00334	0.334	0.000286
	DepthPro _{ls-depth-disp}			✓	1.0	1.0	1.0	0.00241	0.0625	0.001046	0.00334	0.334	0.000286
	DA-V2 _{ls-disp-depth}	✓			1.0	1.0	1.0	0.00227	0.0602	0.000985	0.00321	0.321	0.000244
	Ours _{ls-disp-depth}	✓			1.0	1.0	1.0	0.00223	0.0588	0.000968	0.00314	0.314	0.000235
	DA-V2 _{ls-disp}	✓			1.0	1.0	1.0	0.00226	0.0602	0.000981	0.00321	0.321	0.000244
	Ours _{ls-disp}	✓			1.0	1.0	1.0	0.00222	0.0588	0.000964	0.00314	0.314	0.000235
KITTI	Marigold _{ls-depth}		✓		0.889	0.978	0.992	0.109	3.86	0.047	0.162	16.1	0.53
	Marigold _{ls-depth-disp}		✓		0.876	0.953	0.980	0.110	5.63	0.051	0.181	17.9	0.84
	DepthPro _{ls-depth}			✓	0.937	0.987	0.995	0.086	2.74	0.037	0.132	13.0	0.30
	DepthPro _{ls-depth-disp}			✓	0.896	0.963	0.983	0.096	6.65	0.045	0.186	18.3	3.21
	DA-V2 _{ls-disp-depth}	✓			0.568	0.796	0.902	0.305	7.01	0.118	0.348	33.6	2.49
	Ours _{ls-disp-depth}	✓			0.593	0.818	0.917	0.283	6.71	0.110	0.319	30.7	2.20
	DA-V2 _{ls-disp}	✓			0.818	0.937	0.974	0.323	130	0.062	0.256	25.4	1756
	Ours _{ls-disp}	✓			0.823	0.940	0.976	0.276	105	0.060	0.243	24.1	1335
ETH3D	Marigold _{ls-depth}		✓		0.963	0.994	0.998	0.058	0.476	0.0255	0.101	10.0	0.083
	Marigold _{ls-depth-disp}		✓		0.952	0.988	0.995	0.072	4.666	0.0318	0.178	17.7	30.25
	DepthPro _{ls-depth}			✓	0.966	0.993	0.997	0.058	0.498	0.0237	0.077	7.69	0.158
	DepthPro _{ls-depth-disp}			✓	0.962	0.990	0.994	0.065	5.489	0.0299	0.206	20.5	19.9
	DA-V2 _{ls-disp-depth}	✓			0.884	0.956	0.978	0.113	0.955	0.0448	0.153	15.1	0.391
	Ours _{ls-disp-depth}	✓			0.898	0.965	0.982	0.104	0.875	0.0413	0.143	14.1	0.347
	DA-V2 _{ls-disp}	✓			0.968	0.991	0.995	0.198	35.69	0.0253	0.0998	9.93	1339
	Ours _{ls-disp}	✓			0.968	0.992	0.996	0.148	23.27	0.0253	0.0941	9.36	850.9

metrics, which may include log transformations to reduce the impact of uncertainties. The benefit of aligning to the ground truth is that it provides interpretable results in metric space that are comparable to methods using absolute depth prediction.

DA-V2 deviated from this path by predicting disparity instead of relative depth. Hence, they evaluated directly in disparity space, using the same alignment procedure and metrics. However, this makes it incomparable to models that evaluate on depth. An established alternative is to perform the alignment in disparity space and then apply the metrics after converting disparity back to depth. However, least-squares fitting in disparity handles outliers very differently than when computed in depth. For instance, in disparity, far estimates are less pronounced, leading to an unfair comparison with methods that perform alignment directly in depth space.

To provide an as-fair-as-possible comparison consistent with widely used evaluation protocols, we followed the procedure described in the main document: we first align to obtain absolute disparity, and then perform a second alignment to minimize least-squares errors in the same space used by methods that output relative depth. We refer to this as least-squares disparity-and-depth (ls-disp-depth) alignment. To shed light on the effect of the different alignment methods, we compare them in in Tab. 9. The least-squares alignment performed directly in disparity space followed by depth conversion (ls-disp), as used in [12], performs better on CO3D than ls-disp-depth, but worse on the other two datasets. To further analyze this effect for methods predicting depth, we also mapped depth predictions into disparity space for a second alignment (ls-depth-disp), and observed the same trend. This experiment highlights the importance of the alignment procedure, and we conclude that the fairest comparison is to perform the alignment in the same (depth) space, i.e., to use ls-disp-depth for methods operating on disparity and ls-depth for methods outputting depth directly. Note that the initial alignment in disparity space (ls-disp-depth) is inevitable for methods that output disparity, as it is required to obtain absolute disparity before converting disparity to depth.

10. Baseline Selection

To determine the suitability of DA-V2 as a baseline, we evaluated several existing methods (including Marigold, DepthPro, and DA-V2) on the CO3D dataset. We selected CO3D as our benchmark because its objects exhibit high detail and are relatively close to the camera (e.g., a toy truck instead of a real truck in KITTI). As a result, the depth estimates are more reliable, which aligns well with our goal of detail refinement, as motivated in the main document.

On CO3D, DA-V2 consistently outperformed all other methods. Marigold (diffusion-based foundational model for

Table 4. **Ablation** of the camera model and b parameter on the CO3D dataset.

Method	AbsRel ↓	RMSE ↓	log10 ↓	RMSE log ↓	SI log ↓	SqRel ↓
Ours ($K_{\text{orth}}, b = 0.1$)	0.00223	0.0588	0.000968	0.00314	0.314	0.000235
Ours-Perspective-Fixed ($K_{\text{persp}}, b = 0.1$)	0.00224	0.0591	0.000973	0.00315	0.315	0.000236
Ours-Perspective-Opt ($K_{\text{persp}}, b_{\text{init}} = 0.01$)	0.00225	0.0592	0.000976	0.00316	0.316	0.000237

depth estimation) captures coarse geometric structure, its reconstructed mesh, as shown in Fig. 7, exhibits significant high-frequency noise and artifacts. DepthPro provides an estimate of absolute scale, but it is often misled by toy versions of real objects (e.g., a toy truck), and even after normalizing for scale (and shift), it does not outperform DA-V2. This justifies our selection of DA-V2 as the baseline.

Notably, Table 3 shows that DA-V2 performs better on CO3D but slightly worse on the other two datasets. This is consistent with the evaluations in DepthPro and Marigold, which report improvements over DA-V2 on this datasets. These results support the visual observation that DA-V2 excels in predicting fine details while lacking a bit in capturing overall scene composition and the relative scale of objects.

Another closely related work is BetterDepth [61], which, similar to our approach, focuses on refining the output of a pre-trained depth foundation model. However, their implementation is not publicly available.

11. Ablation - Perspective Camera

We conduct an ablation study on the camera model (Eq. 2 in the main paper) and the choice of $b = ms$ (Eq. 5 in the main paper). We compare a scaled orthographic camera, with its scale factor optimized from an initial value of 7.0, against a perspective camera, with its focal length jointly optimized from an initial value of 2.0. As shown in Table 4, the scaled orthographic model with $b = 0.1$ yields the lowest errors. We therefore adopt this configuration for all main experiments.

12. SfS Implementation Details

We implement a simple shape from shading algorithm similar to [11]. We assume a Lambertian surface, with constant albedo and a single light at infinity such that the light direction \mathbf{l} is unknown yet constant across all pixels. We consider only brightness variations and therefore convert the input image to grayscale and set the incoming light intensity to $L_{\text{in}} = \max(\mathbf{I})$, where \mathbf{I} is the input image.

Under these assumptions, the image formation model reduces to a Lambertian dot product between the surface normal and the light direction:

$$\hat{\mathbf{I}}(u, v) = \max(0, \mathbf{N}(u, v) \cdot \mathbf{l}),$$

where $\hat{\mathbf{I}}$ is the rendered image.

Motivated by the shape from intensity gradient technique [59], we initialize the normals using image gradients,

$$\mathbf{N}(u, v) = (-\mathbf{I}_u(u, v), -\mathbf{I}_v(u, v), 1),$$

where \mathbf{I}_u and \mathbf{I}_v denote the partial derivatives of \mathbf{I} with respect to the spatial coordinates u and v . We compute these derivatives using Scharr filters [38].

To optimize light direction and surface normal, we minimize a combination of a photometric loss and a smoothness loss,

$$\mathcal{L} = \mathcal{L}_{\text{smooth}} + \lambda \mathcal{L}_{\text{photo}},$$

where λ is a regularization parameter. The photometric loss minimizes the difference between the input and the rendered images,

$$\mathcal{L}_{\text{photo}} = \frac{1}{|\Omega|} \sum_{(u,v) \in \Omega} (\mathbf{I}(u, v) - \hat{\mathbf{I}}(u, v))^2,$$

where Ω denotes the valid region defined by the object mask. The smoothness loss penalizes spatial variation in the normals,

$$\mathcal{L}_{\text{smooth}} = \frac{1}{|\Omega|} \sum_{(u,v) \in \Omega} (\|\nabla \mathbf{N}_u(u, v)\|_2^2 + \|\nabla \mathbf{N}_v(u, v)\|_2^2),$$

where \mathbf{N}_u and \mathbf{N}_v denote the first two components of the normal field.

This simple shading model closely matches the re-lighting procedure used in our main method, highlighting the benefit of our shading-based augmentation. Unlike this full photometric reconstruction, which produces artifacts at texture boundaries or under complex real-world illumination, our re-lighting refinement succeeds with a simple and robust illumination model.

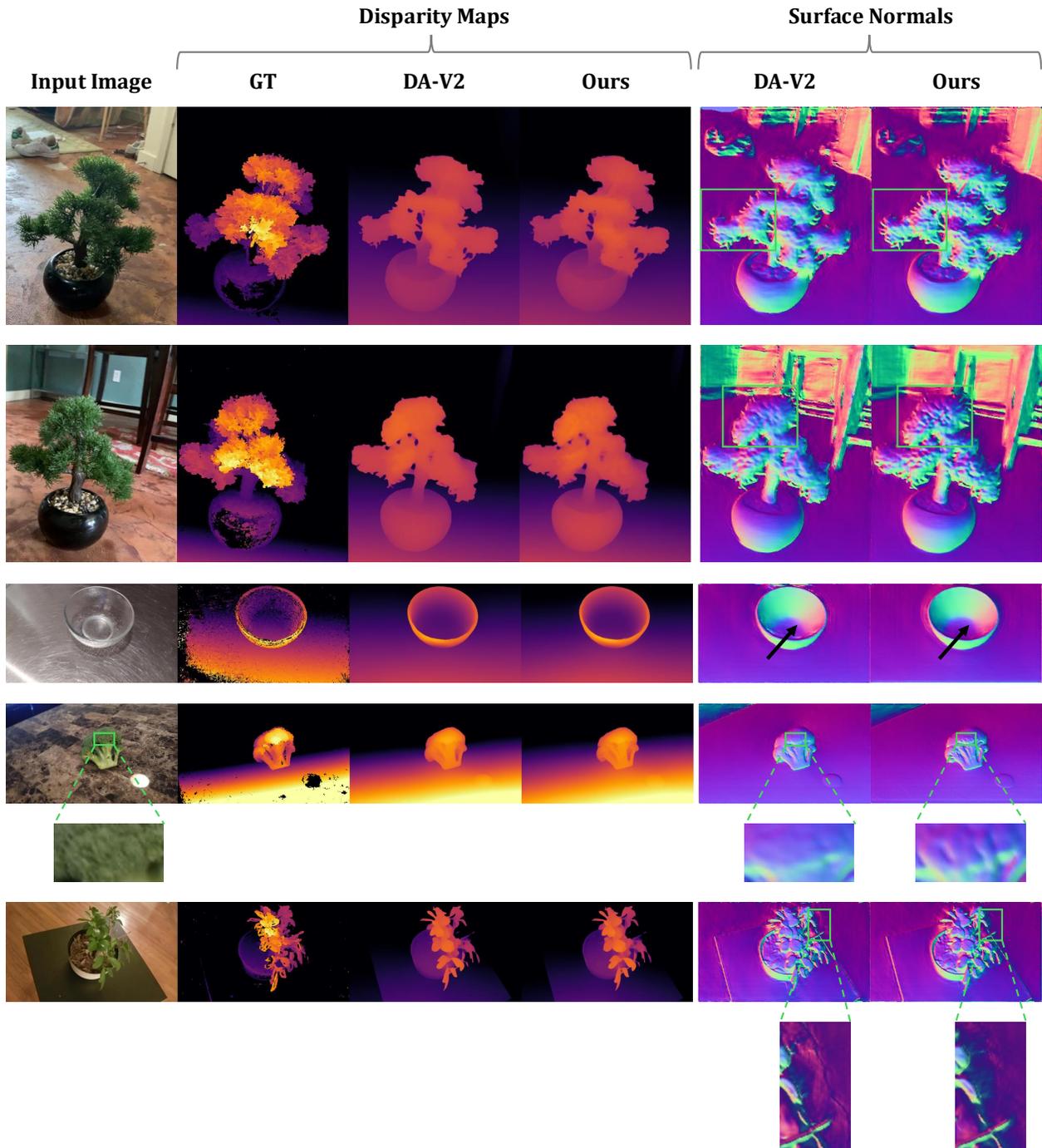


Figure 9. Additional qualitative comparison on the CO3D dataset.

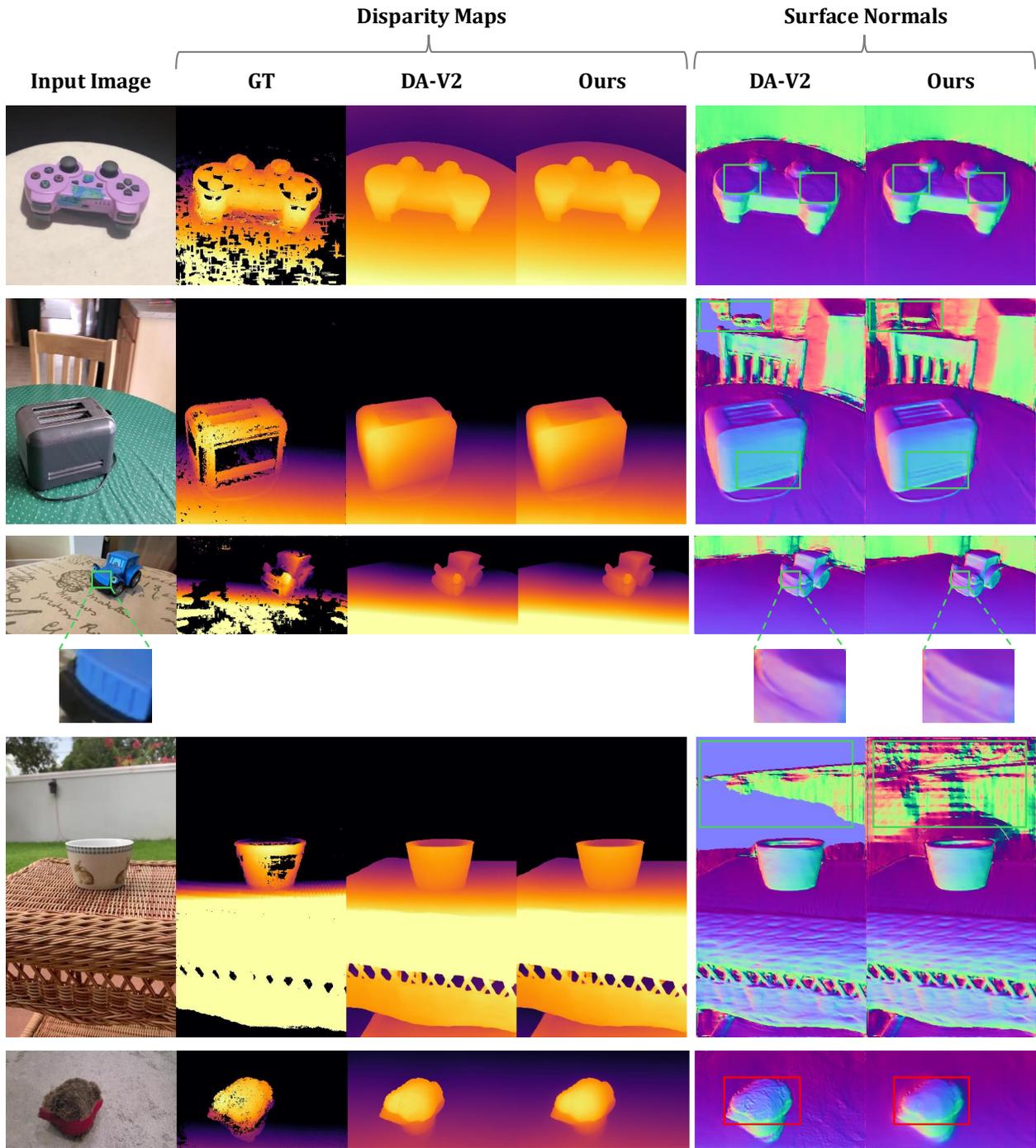


Figure 10. **Additional qualitative comparison on the CO3D dataset.** Red square highlights occasional oversmoothing, a limitation of our method.

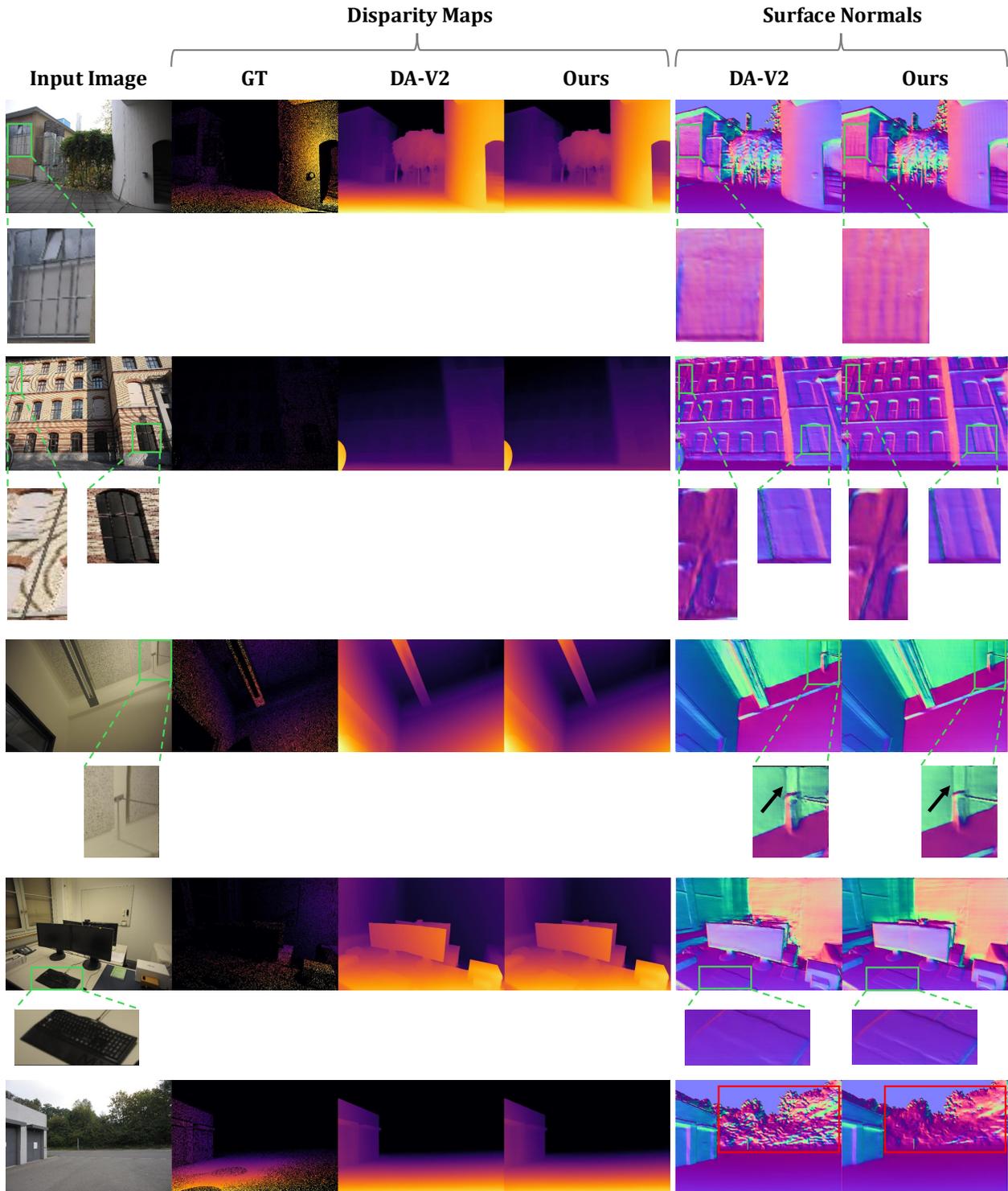


Figure 11. Additional qualitative comparison on the ETH3D dataset. Red square highlights occasional oversmoothing, a limitation of our method.

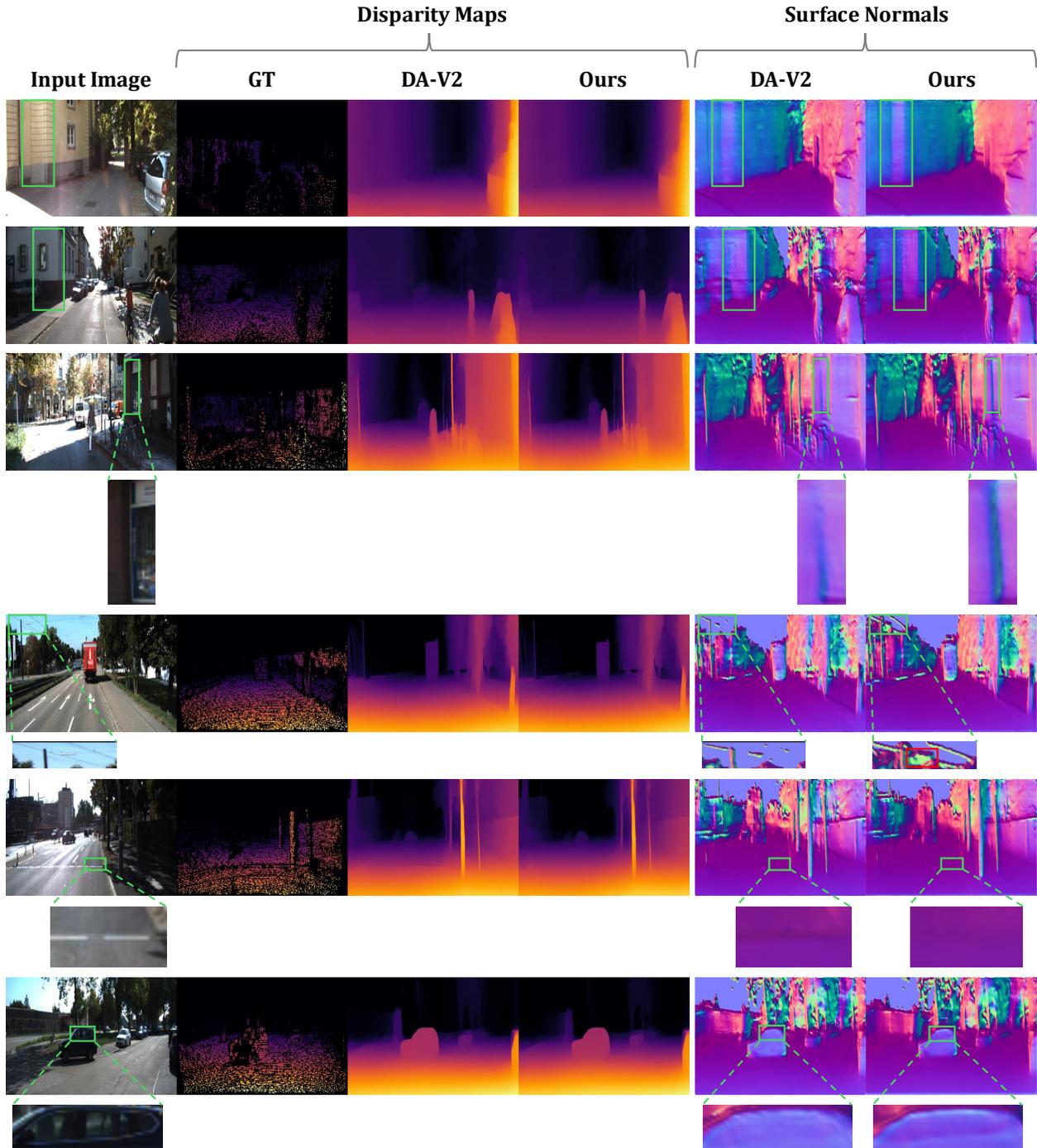


Figure 12. **Additional qualitative comparison on the KITTI dataset.** Red square highlights hallucination on sky areas, a limitation of our method.

References

- [1] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [2] Ananta R Bhattarai, Xingzhe He, Alla Sheffer, and Helge Rhodin. Dreamtexture: Shape from virtual texture with analysis by augmentation. *arXiv preprint arXiv:2503.16412*, 2025. 3
- [3] James F. Blinn. Models of light reflection for computer synthesized pictures. In *Computer Graphics and Interactive Techniques*, 1977. 2, 3, 4
- [4] Aleksei Bochkovskii, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R Richter, and Vladlen Koltun. Depth pro: Sharp monocular metric depth in less than a second. In *International Conference on Learning Representations (ICLR)*, 2025. 2
- [5] Johann Cabon, Lucas Stoffl, Leonid Antsfeld, Gabriela Csurka, Boris Chidlovskii, Jerome Revaud, and Vincent Leroy. Must3r: Multi-view network for stereo 3d reconstruction. In *Computer Vision and Pattern Recognition Conference (CVPR)*, 2025. 2
- [6] Zilong Chen, Feng Wang, and Huaping Liu. Text-to-3d using gaussian splatting. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. 1
- [8] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Neural Information Processing Systems (NeurIPS)*, 2014. 2
- [9] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32:1231 – 1237, 2013. 2, 6
- [10] Berthold K. P. Horn and Michael J. Brooks. *Shape from Shading*. MIT Press, Cambridge, Massachusetts, 1989. 3
- [11] Berthold K. P. Horn and Michael J. Brooks. *The variational approach to shape from shading*, page 173–214. MIT Press, Cambridge, MA, USA, 1989. 3
- [12] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing diffusion-based image generators for monocular depth estimation. In *Computer Vision and Pattern Recognition Conference (CVPR)*, 2024. 2, 5, 3
- [13] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 42:1 – 14, 2023. 1, 2
- [14] Vincent Leroy, Johann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. In *European Conference on Computer Vision (ECCV)*, 2024. 2
- [15] Bo Li, Chunhua Shen, Yuchao Dai, Anton van den Hengel, and Mingyi He. Depth and surface normal estimation from monocular images using regression on deep features and hierarchical crfs. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2
- [16] Huan Li, Matteo Poggi, Fabio Tosi, and Stefano Mattoccia. On-site adaptation for monocular depth estimation with a static camera. In *British Machine Vision Conference (BMVC)*, 2023. 2
- [17] Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning (ICML)*, 2023. 5
- [18] Zhenyu Li, Xuyang Wang, Xianming Liu, and Junjun Jiang. Binsformer: Revisiting adaptive bins for monocular depth estimation. *IEEE Transactions on Image Processing*, 33: 3964–3976, 2022. 2
- [19] Zhi Li, Shaoshuai Shi, Bernt Schiele, and Dengxin Dai. Test-time domain adaptation for monocular depth estimation. In *International Conference on Robotics and Automation (ICRA)*, 2023. 2
- [20] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3
- [21] 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 *International Conference on Computer Vision (ICCV)*, 2023. 3
- [22] Jonathan Lorraine, Kevin Xie, Xiaohui Zeng, Chen-Hsuan Lin, Towaki Takikawa, Nicholas Sharp, Tsung-Yi Lin, Ming-Yu Liu, Sanja Fidler, and James Lucas. Att3d: Amortized text-to-3d object synthesis. In *International Conference on Computer Vision (ICCV)*, 2023. 2, 3
- [23] R’emi Marsal, Alexandre Chapoutot, Philippe Xu, and David Filliat. A simple yet effective test-time adaptation for zero-shot monocular metric depth estimation. 2024. 2
- [24] Luke Melas-Kyriazi, Iro Laina, C. Rupprecht, and Andrea Vedaldi. Realfusion 360° reconstruction of any object from a single image. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [25] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [26] 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 *European Conference on Computer Vision (ECCV)*, 2020. 1, 2
- [27] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Q. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russ Howes, Po-Yao (Bernie) Huang, Shang-Wen Li, Ishan Misra,

- Michael G. Rabbat, Vasu Sharma, Gabriel Synnaeve, Huijiao Xu, Hervé Jégou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. *ArXiv*, abs/2304.07193, 2023. 3
- [28] Hyungseob Park, Anjali W. Gupta, and Alex Wong. Test-time adaptation for depth completion. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [29] Luigi Piccinelli, Christos Sakaridis, Yung-Hsu Yang, Mattia Segu, Siyuan Li, Wim Abbeloos, and Luc Van Gool. Unidepthv2: Universal monocular metric depth estimation made simpler. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 2
- [30] Dustin Podell, Zion English, Kyle Lacey, A. Blattmann, Tim Dockhorn, Jonas Muller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. In *International Conference on Learning Representations (ICLR)*, 2023. 2
- [31] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *International Conference on Learning Representations (ICLR)*, 2023. 2, 3
- [32] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *Transactions on Pattern Analysis and Machine Intelligence*, 44:1623–1637, 2019. 1, 2
- [33] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 3
- [34] Alex Rasla and Michael Beyeler. The relative importance of depth cues and semantic edges for indoor mobility using simulated prosthetic vision in immersive virtual reality. *ACM Symposium on Virtual Reality Software and Technology*, 2022. 1
- [35] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *International Conference on Computer Vision (ICCV)*, 2021. 6
- [36] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 5
- [37] 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. In *Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [38] Hanno Schar. *Optimal operators in digital image processing*. PhD thesis, University of Heidelberg, Germany, 2000. 4
- [39] Thomas Schops, Torsten Sattler, and Marc Pollefeys. Bad slam: Bundle adjusted direct rgb-d slam. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 6
- [40] Shuwei Shao, Zhongcai Pei, Weihai Chen, Xingming Wu, and Zhengguo Li. Ndepth: Normal-distance assisted monocular depth estimation. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [41] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *European Conference on Computer Vision (ECCV)*, 2012. 2
- [42] Jiayang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. In *International Conference on Learning Representations (ICLR)*, 2023. 3
- [43] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Gregory Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3
- [44] Ruicheng Wang, Sicheng Xu, Cassie Dai, Jianfeng Xiang, Yu Deng, Xin Tong, and Jialong Yang. Moge: Unlocking accurate monocular geometry estimation for open-domain images with optimal training supervision. In *Computer Vision and Pattern Recognition Conference (CVPR)*, 2025. 2
- [45] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [46] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark E. Campbell, and Kilian Q. Weinberger. Pseudolidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [47] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. In *Neural Information Processing Systems (NeurIPS)*, 2023. 2, 3
- [48] Ryan White and David A Forsyth. Combining cues: Shape from shading and texture. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2006. 3
- [49] Andrew P. Witkin. Recovering surface shape and orientation from texture. *Artificial Intelligence*, 17(1-3):17–45, 1981. 3
- [50] Diana Wofk, Fangchang Ma, Tien-Ju Yang, Sertac Karaman, and Vivienne Sze. Fastdepth: Fast monocular depth estimation on embedded systems. In *International Conference on Robotics and Automation (ICRA)*, 2019. 1
- [51] Ke Xian, Jianming Zhang, Oliver Wang, Long Mai, Zhe L. Lin, and ZHIGUO CAO. Structure-guided ranking loss for single image depth prediction. In *Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [52] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [53] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. In *Neural Information Processing Systems (NeurIPS)*, 2024. 1, 2, 3, 6

- [54] Xiaodong Yang, Zhuang Ma, Zhiyu Ji, and Zhe Ren. Gedept: Ground embedding for monocular depth estimation. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [55] Chongjie Ye, Yinyu Nie, Jiahao Chang, Yuantao Chen, Yihao Zhi, and Xiaoguang Han. Gaustudio: A modular framework for 3d gaussian splatting and beyond. *ArXiv*, abs/2403.19632, 2024. 1
- [56] Wei Yin, Yifan Liu, Chunhua Shen, and Youliang Yan. Enforcing geometric constraints of virtual normal for depth prediction. In *International Conference on Computer Vision (ICCV)*, 2019. 2
- [57] Wei Yin, Chi Zhang, Hao Chen, Zhipeng Cai, Gang Yu, Kaixuan Wang, Xiaozhi Chen, and Chunhua Shen. Metric3d: Towards zero-shot metric 3d prediction from a single image. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- [58] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. Neural window fully-connected crfs for monocular depth estimation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3906–3915, 2022. 2
- [59] Ruo Zhang and Mubarak Shah. Shape from intensity gradient. *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, 29(3):318–325, 1999. 4
- [60] Ruo Zhang, Ping-Sing Tsai, James Edwin Cryer, and Mubarak Shah. Shape-from-shading: a survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(8):690–706, 1999. 3
- [61] Xiang Zhang, Bingxin Ke, Hayko Riemenschneider, Nando Metzger, Anton Obukhov, Markus Gross, Konrad Schindler, and Christopher Schroers. Betterdepth: Plug-and-play diffusion refiner for zero-shot monocular depth estimation. In *Neural Information Processing Systems (NeurIPS)*, 2024. 2, 3
- [62] Yizhou Zhao, Hengwei Bian, Kaihua Chen, Pengliang Ji, Liao Qu, Shao yu Lin, Weichen Yu, Haoran Li, Hao Chen, Jun Shen, Bhiksha Raj, and Min Xu. Metric from human: Zero-shot monocular metric depth estimation via test-time adaptation. In *Neural Information Processing Systems (NeurIPS)*, 2024. 2
- [63] Junzhe Zhu, Peiye Zhuang, and Oluwasanmi Koyejo. Hifa: High-fidelity text-to-3d generation with advanced diffusion guidance. In *International Conference on Learning Representations (ICLR)*, 2023. 2, 3