

UniK3D: Universal Camera Monocular 3D Estimation

Luigi Piccinelli¹ Christos Sakaridis¹ Mattia Segu¹
Yung-Hsu Yang¹ Siyuan Li¹ Wim Abbeloos² Luc Van Gool^{1,3}

¹ETH Zürich ²Toyota Motor Europe ³INSAIT, Sofia University St. Kliment Ohridski

Abstract

Monocular 3D estimation is crucial for visual perception. However, current methods fall short by relying on oversimplified assumptions, such as pinhole camera models or rectified images. These limitations severely restrict their general applicability, causing poor performance in real-world scenarios with fisheye or panoramic images and resulting in substantial context loss. To address this, we present UniK3D¹, the first generalizable method for monocular 3D estimation able to model any camera. Our method introduces a spherical 3D representation which allows for better disentanglement of camera and scene geometry and enables accurate metric 3D reconstruction for unconstrained camera models. Our camera component features a novel, model-independent representation of the pencil of rays, achieved through a learned superposition of spherical harmonics. We also introduce an angular loss, which, together with the camera module design, prevents the contraction of the 3D outputs for wide-view cameras. A comprehensive zero-shot evaluation on 13 diverse datasets demonstrates the state-of-the-art performance of UniK3D across 3D, depth, and camera metrics, with substantial gains in challenging large-field-of-view and panoramic settings, while maintaining top accuracy in conventional pinhole small-field-of-view domains. Code and models are available at github.com/lpiccinelli-eth/unik3d.

1. Introduction

Estimating 3D scene geometry is a fundamental task in computer vision since such 3D information serves as a crucial cue for action planning and execution [14, 89]. The scene’s geometry 3D estimation task is vital for a wide range of applications, including autonomous navigation [56, 76] and 3D modeling [13], where accurate spatial understanding is essential. Recent advances in generalizable monocular depth estimation (MDE) [32, 63, 81] deliver impressive performance and visual quality across various domains, but these mod-

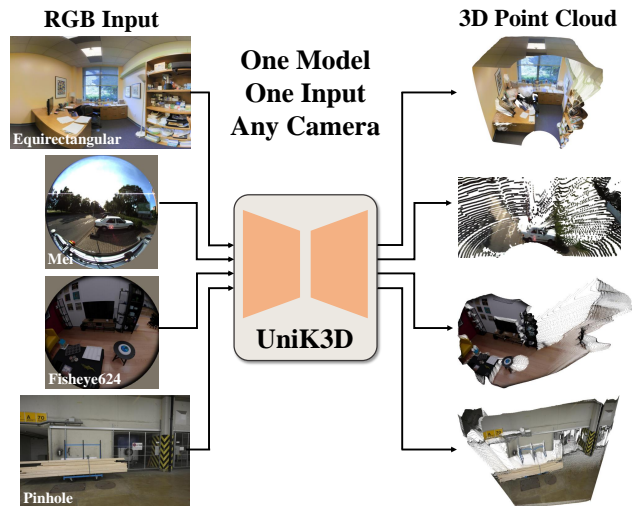


Figure 1. **UniK3D** introduces a novel and versatile approach that delivers precise metric 3D geometry estimation from a single image and for any camera type, ranging from pinhole to panoramic, without requiring any camera information. By leveraging (i) a flexible and general spherical formulation both for the radial dimension of 3D space and for the two camera-model-dependent orientation dimensions and (ii) advanced conditioning strategies. UniK3D outperforms traditional models without needing camera calibration or domain-specific tuning.

els are constrained to a relative output scale. Nonetheless, for practical applications, a consistent and reliable *metric-scaled* monocular depth estimate (MMDE) is crucial, as it enables accurate 3D reconstruction and geometric scene understanding necessary for embodied agents.

Existing methods have made considerable strides in the above direction of metric estimation. Earlier approaches assumed known camera intrinsics at test time [24, 85], while more recent works have relaxed this assumption [9, 60, 61]. However, these approaches still impose restrictive assumptions about input cameras, such as relying on a basic pinhole camera model [9, 60] or requiring access to ground-truth rectification parameters [85]. These simplifications substantially hinder the applicability and degrade the performance

¹Pronounced “Unique-3D”, with **K** denoting the intrinsics matrix.

of the above methods in real-world settings, where a wide range of camera projection models with strong non-linear deformations are common, such as fisheye or panoramic lenses. This limitation is more pronounced when estimating complete metric 3D geometry instead of only depth maps, as the former depends more heavily on the quality of camera estimation. Due to the restrictive assumptions in existing models, general camera estimation can not be effectively learned, even when models are exposed to images from varied camera types. Furthermore, the output space of previous state-of-the-art MMDE methods has inherent limitations, *e.g.* both disparity and log-depth prediction become mathematically ill-posed when the field of view (FoV) exceeds 180 degrees.

To address these challenges, we introduce *UniK3D*, the first framework for monocular metric 3D scene’s geometry estimation that generalizes across a wide variety of camera models, from pinhole to fisheye and panoramic configurations, as shown in Fig. 1. Our method proposes a novel formulation for monocular 3D estimation which is spherical in two senses. First, UniK3D leverages a *fully spherical* output 3D space, modeling the range dimension through *radial distance* instead of perpendicular depth. This approach is especially beneficial at large angles from the optical axis, effectively resolving the ill-posed nature of traditional methods at extreme fields of view. Second, while building on the recently proposed decomposition [60] of camera prediction from depth estimation, UniK3D newly presents a general *spherical harmonics basis* as the *direct* output space of the camera module that represents the pencil of rays. Unlike previous works [9, 60] which predict explicit pinhole camera parameters and then encode [60] induced rays using a spherical basis, we remove the camera assumption and directly model the rays. As a result, UniK3D spans an unrestricted space of possible camera models, allowing for flexible and accurate depth prediction regardless of camera intrinsics. Our assumption-free spherical camera representation, with its flexibility, ensures that our model is well-suited for real-world deployment, where capturing scenes with non-standard cameras is common.

Our key contribution is the first camera-universal model for monocular 3D estimation that can accommodate any camera projective geometry. We achieve this through our unified spherical output representation that supports all inverse projection problems. By employing a fully spherical framework, our method ensures a complete disentanglement of projective *vs.* 3D scene geometry, as the dimension of an object projection on the image is a univocal function only w.r.t. radial distance and not w.r.t. depth. This disentanglement allows more consistent 3D reconstructions and enhances the stability of 3D outputs near the *xy*-plane, where depth approaches zero. Moreover, UniK3D models the camera rays as a decomposition across a finite spherical harmonics basis. This choice ensures representation generality and versatility,

and at the same time maintains an accurate and compact representation for the resulting pencils of rays, also introducing inductive biases such as continuity and differentiability. In addition, we propose multiple novel strategies to ensure robust camera conditioning of our *radial module* such as an asymmetric angular loss based on quantile regression, static encoding, and curriculum learning.

We validate our approach through extensive zero-shot experiments on 13 widely used metric depth datasets, where UniK3D not only achieves state-of-the-art performance in monocular metric depth and 3D estimation, but also generalizes very well across various camera models, without either preprocessing or specific camera domains during training.

2. Related Work

Monocular Depth Estimation. The introduction of end-to-end neural networks for MDE, first demonstrated by [16], revolutionized the field by enabling depth prediction through direct optimization, utilizing the Scale-Invariant log loss (SI_{log}). Since then, the field has evolved with increasingly sophisticated models, ranging from convolutional architectures [19, 36, 45, 58] to recent advancements using transformers [6, 59, 80, 86]. While these approaches have pushed the boundaries of MDE performance in controlled benchmarks, they often fail when faced with zero-shot scenarios, highlighting a persistent challenge: ensuring robust generalization across varying camera and scene domains and diverse geometric and visual conditions.

Generalizable Monocular Depth Estimation. To address the limitations of domain-specific models, recent research has focused on developing generalizable and zero-shot MDE techniques. These methods can be categorized into scale-agnostic approaches [32, 63, 73, 81, 82], which aim to mitigate scale ambiguity and emphasize perceptual depth quality, and metric depth models [7, 9, 24, 28, 60, 61, 85], which prioritize accurate geometric reconstruction. However, most existing MDE methods fall short of achieving truly zero-shot monocular metric 3D scene estimation. In particular, scale-agnostic approaches often require additional information to resolve scale ambiguities, while most of the metric-based models depend on a known camera or assume a simplistic pinhole camera configuration. Even the few models which are designed for zero-shot 3D scene estimation [9, 60, 85] remain constrained: they either explicitly assume a pinhole camera model [9, 60] or necessitate image rectification [85], effectively requiring test-time camera information and limiting their zero-shot generalizability to pinhole cameras.

On the contrary, UniK3D addresses these limitations by offering a unified solution that can handle any inverse projection problem. Our model can recover a coherent 3D point cloud from any single image, regardless of camera intrinsics, without any rectification or camera information at test time. This generality sets UniK3D apart, enabling robust and uni-

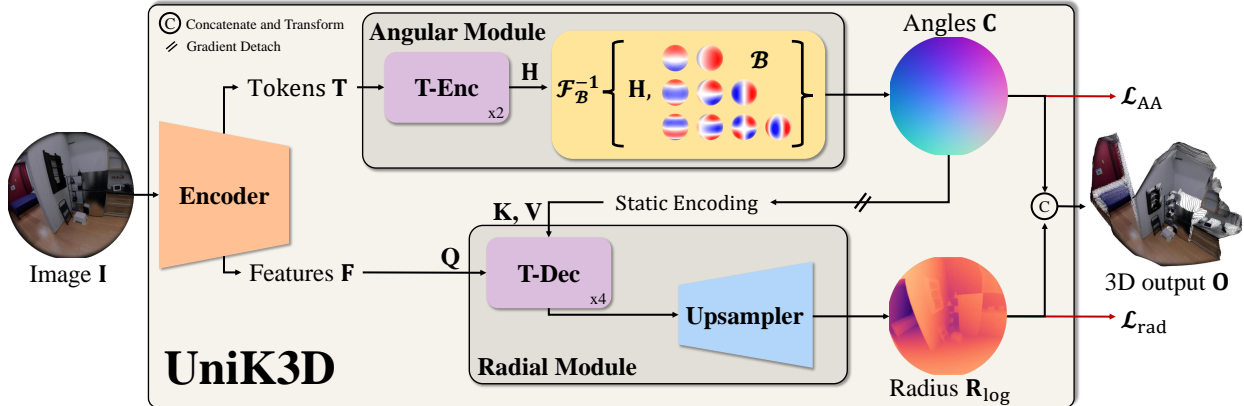


Figure 2. **Model architecture.** UniK3D utilizes solely the single input image to generate the 3D output point cloud (\mathbf{O}) for any camera. The projective geometry of the camera is predicted by the Angular Module. The camera representation corresponds to azimuth and polar angles (\mathbf{C}) of the backprojected pencil of rays on the unit sphere \mathbb{S}^3 . The class tokens from the Encoder are processed by 2 Transformer Encoder (T-Enc) layers to obtain the 15 coefficients (\mathbf{H}) of the inverse Spherical transform $\mathcal{F}_{\mathcal{B}}^{-1}\{\mathbf{H}\}$ defined by a finite basis (\mathcal{B}) of spherical harmonics up to degree 3 with no constant component. Stop-gradient is applied to the angular information which conditions the Radial Module, simulating external information flow. The “static encoding” refers to sinusoidal encoding which matches the internal feature dimensionality. The Radial Module is composed of Transformer Decoder (T-Dec) blocks, one for each input resolution, which is utilized to condition the Encoder features on the bootstrapped camera representation. This conditioning injects prior knowledge on scene scale and projective geometry. The radial output (\mathbf{R}_{\log}) is obtained by processing the camera-aware features via a learnable upsampling module. The final output is the concatenation of the camera and radial tensors ($\mathbf{C} \parallel \mathbf{R}_{\log}$). A closed-form coordinate transform is applied to obtain the Cartesian 3D output, but supervision is applied directly on angular coordinates, with our asymmetric angular loss \mathcal{L}_{AA} , and radial coordinates.

versal monocular metric 3D estimation that is required in diverse and challenging real-world applications.

Camera Calibration. Camera calibration is essential for estimating intrinsic parameters like focal length, principal point, and distortion coefficients to model the mapping from 3D world points to 2D image coordinates. Traditional parametric models, such as the pinhole model, Kannala-Brandt [30], Mei [49], Omnidirection [64], Unified Camera Model (UCM) [22], Enhanced UCM [33], and Double Sphere [69] models are effective for narrow- and wide-angle lenses but require controlled environments for accurate calibration. As models grow more complex, the risk of errors or divergence increases, especially under varying lighting or sensor noise. Additionally, each model has inherent limitations, *e.g.* UCM cannot represent tangential distortion, and Kannala-Brandt struggles beyond a 210° FoV.

By contrast, we take a different approach and model the camera backprojection as a linear combination of *spherical basis* functions, *i.e.* via an inverse spherical harmonics transformation, where the model simply infers the scalar expansion coefficients and the spherical domain boundaries.

3. UniK3D

Generalizable depth or 3D scene estimation models often face significant challenges when adapting to diverse camera configurations. Existing methods typically rely on rigid and camera-specific assumptions, such as the pinhole model or equirectangular models, or require extensive preprocessing

steps like rectification. These constraints limit their applicability to real-world scenarios with non-standard camera projective geometries. By contrast, our model, UniK3D, introduces a novel framework that enables monocular 3D geometry estimation for any scene and any camera setup.

We begin by introducing the design of our 3D output space and the internal representation of the camera in Sec. 3.1. Our representation is intentionally formulated to be as general as possible, allowing to handle all inverse projection problems. Through our preliminary studies, we observed a consistent issue: the network predictions contracted to a reduced FoV, even when trained on a diverse set of camera types including large FoVs. Simple data re-balancing strategies proved insufficient to address this phenomenon. To overcome this, we have developed a series of architectural and design interventions, detailed in Sec. 3.2, aimed at preventing the backprojection contraction. In Sec. 3.3, we describe the architecture of our model, our optimization strategy, and the specific design and loss functions underpinning our approach. Fig. 2 displays an overview of our method.

3.1. Representations

Output Space. The output representation of UniK3D is designed to be universally compatible with any scene and camera configuration, providing a direct metric 3D scene estimate for each input image. Drawing from the disentanglement strategy presented in [60], our approach separates camera parameters from scene geometry. Specifically, we represent the camera using a dense tensor $\mathbf{C} = \theta \parallel \phi$, where θ

is the polar angle and ϕ is the azimuthal angle, consistent with standard spherical coordinates. However, we use the Euclidean *radius* (distance from the camera center) as the scene range component within a *fully spherical* framework, instead of relying on traditional perpendicular-depth-based representations. This design choice ensures that dimensions of projected objects in the image vary *univocally* with radius, a property that does not characterize the depth representation and renders the latter much harder to learn. Furthermore, the spherical framework enhances numerical stability when handling points near the xy -plane, a region where previous methods typically face challenges due to large gradients. We convert the spherical representation to Cartesian coordinates using a bijective transformation, accurately capturing the 3D geometry of the scene as the output 3D point cloud \mathbf{O} .

Camera Internal Space. In UniK3D, the dense pencil of rays which represents the viewing directions for the various pixels is expressed through a basis decomposition, providing a flexible and comprehensive angular representation. As shown in Fig. 2, our Angular Module predicts a tensor of coefficients \mathbf{H} , which is derived from the encoder’s class tokens, denoted as \mathbf{T} . These coefficients correspond to a predefined basis: the Spherical Harmonics (SH) basis. We reconstruct the pencil of rays from \mathbf{H} as follows:

$$\mathbf{C} = \mathcal{F}_{\mathcal{B}}^{-1}\{\mathbf{H}\} = \sum_{l=0}^L \sum_{m=-l}^l \mathbf{H}_{lm} \mathcal{B}_{lm}(\theta, \phi), \quad (1)$$

where \mathbf{C} represents the reconstructed angular field and $\mathcal{F}_{\mathcal{B}}^{-1}$ denotes the inverse transform from the coefficient space to the angular space, using the SH basis \mathcal{B} . $\mathcal{B}_{lm}(\theta, \phi)$ are the SH basis functions, *i.e.* Legendre polynomials, and \mathbf{H}_{lm} are the predicted coefficients. Here, l and m index the degree and order of the harmonics, respectively. This inverse transform is implemented as an inner product that maps from $\mathbb{R}^n \times \mathbb{S}^3$ to \mathbb{S}^3 . The SH basis domain is defined by 4 parameters: the generalized “principal point” of the reference frame, *i.e.* the pole, and the horizontal and vertical FoVs. This formulation allows us to describe complex ray distributions *compactly* and implicitly, while ensuring important properties of the output, such as continuity and differentiability.

Additionally, the SH basis offers high sparsity, requiring only 15 harmonics for a 3rd degree basis without constant component and an equal number of coefficients to accurately represent intrinsics for most camera types. By leveraging this SH-based representation and defining the domain through the pole and FoV parameters, UniK3D achieves a robust and flexible framework that can handle virtually any camera geometry with only 19 parameters.

3.2. Preventing Distribution Contraction

Asymmetric Angular Loss. Neural networks tend to regress towards the most frequent modes in the training data, often

neglecting the distribution tails. In our case, this bias would cause UniK3D to underrepresent wide-FoV angles in its outputs, since most visual datasets are heavily skewed towards small-FoV pinhole cameras. This leads to poor performance in scenarios requiring accurate wide-angle predictions. To overcome this issue, we introduce an asymmetric angular loss based on quantile regression, inspired by robust statistical estimators and decision theory principles, *i.e.* type-I and type-II errors [51]. Our loss function is defined as:

$$\mathcal{L}_{\text{AA}}^{\alpha}(\hat{\theta}, \theta^*) = \alpha \sum_{\hat{\theta} > \theta^*} |\hat{\theta} - \theta^*| + (1 - \alpha) \sum_{\hat{\theta} \leq \theta^*} |\hat{\theta} - \theta^*|, \quad (2)$$

where $0 \leq \alpha \leq 1$ is the target quantile, $\hat{\theta}$ is the predicted angle, and θ^* is the ground-truth angle. This formulation adjusts the weighting of over- and underestimations of the polar angle θ . When $\alpha = 0.5$, the loss degenerates to the standard Mean Absolute Error (MAE), but by tuning α , we can emphasize underrepresented angles and balance the regression more effectively. Unlike naive dataset rebalancing—which would alter the underlying 3D scene diversity and introduce significant complexity, especially across multiple datasets—our loss addresses the angular imbalance directly and efficiently. By using quantile regression, we minimize the complexity to a simple search over the interval $[0, 1]$ for α , making our method well-suited for large-scale and diverse training scenarios. This quantile-based strategy allows us to address the angular distribution bias without sacrificing simplicity and diversity, making it a robust and scalable solution. **Enhancing Camera Conditioning.** In our initial experiments, we observed that our model struggled to effectively utilize camera conditioning following previous works [60], even when explicitly supplied with ground-truth camera rays during both training and testing. This issue was subtle for small-FoV pinhole cameras, but it became significant for large-FoV configurations. The root of the problem lies in weak conditioning: the model fails to disentangle camera parameters from geometric features, causing it to route local aberrations back to the encoder features’ space, without integrating essential FoV information. As a result, even when prompted with accurate camera parameters at test time, the model might ignore, or be misled by, this information.

To address this, we hypothesize that camera data must be clear and explicitly structured from the beginning of training. To this end, we implement in UniK3D a static (non-learnable) encoding of camera rays and adopt a curriculum learning strategy, transitioning gradually from feeding GT camera parameters to predicted ones to the Radial Module. In particular, the GT camera is fed to the Radial Module with probability $1 - \tanh(\frac{s}{10^5})$, where s is the current optimization step. To reinforce external conditioning, we detach gradients from the camera output that is fed to the Radial Module, hence preventing the model from relying on feedback mechanisms that could undermine the conditioning on

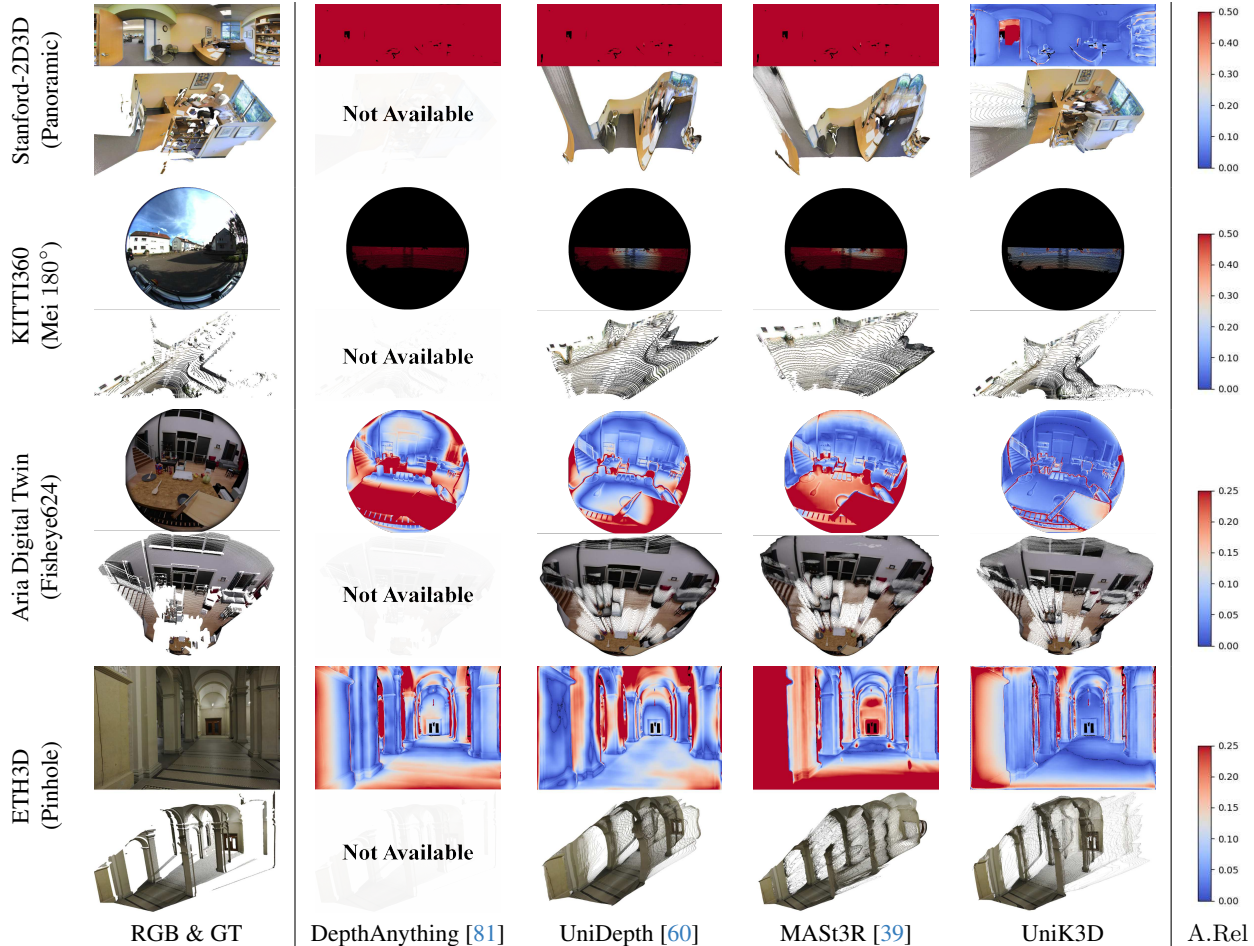


Figure 3. **Qualitative comparisons.** Each pair of consecutive rows represents one test sample. Each odd row displays the input RGB image and the 2D error map, color-coded with the *coolwarm* colormap based on absolute relative error (for panoramic images, the error is computed on distance rather than depth). To ensure a fair comparison, errors are calculated on GT-based shifted and scaled outputs for all models. Each even row shows the ground truth and predictions of the 3D point cloud. The last column displays the specific colormap ranges for absolute relative error. Key observations for each rows pair: (1) competing methods are limited to only positive depth and heavily distort the scenes for larger FoV; (2) in the case of representable but large FoV (180°), UniK3D output is the only one that does not suffer from pronounced FoV contraction; (3) for moderate-FoV images but with strong boundary distortion, *e.g.* fisheye, UniK3D can maintain planarity and overall scene structure; (4) our approach also delivers accurate 3D estimates for standard pinhole images.

the camera. Additionally, we disable learnable gains, such as LayerScale [68], in the cross-attention layers of the Radial Module’s transformer decoder, to avoid shortcuts of the conditioning. These strategies ensure that the model effectively leverages camera information to adjust its encoder features, enhancing the robustness of 3D predictions.

3.3. Network Design

Architecture. Our network consists of an Encoder Backbone, an Angular Module, and a Radial Module, as illustrated in Fig. 2. Our encoder is ViT-based [15] and we extract dense features $\mathbf{F} \in \mathbb{R}^{h \times w \times C \times 4}$ —where $(h, w) = (\frac{H}{14}, \frac{W}{14})$ —along with class tokens \mathbf{T} . The Angular Module processes these class tokens, projecting them onto 512-channel representations that are split into 3 domain parameters and 15

spherical coefficient prototypes. These tokens pass through two layers of a Transformer Encoder (T-Enc) with 8 heads and are then projected onto scalar values. The values for the 3 domain parameters define the principal point (2) and the horizontal FoV (1), determining the intervals for the harmonics. We assume square pixels and thus do not learn an extra, fourth parameter for the vertical FoV, but rather compute this fourth parameter directly from the horizontal FoV. The 15 spherical coefficients undergo an inverse SH transformation according to (1), using a 3-degree SH basis. The gradient flowing from the Angular Module to the class tokens is multiplied by 0.1, as the magnitude of the camera-induced gradient for the encoder weights was empirically found to be ca. 10x higher than the radial-induced gradient.

Table 1. **Comparison on zero-shot evaluation for diverse camera domains.** Validation sets: *S.FoV* includes NYU, KITTI, IBims-1, ETH-3D, nuScenes, and Diode Indoor; *S.FoV_{Dist}* includes IBims-1, ETH-3D, and Diode Indoor with synthetic distortion; *L.FoV* includes ADT, ScanNet++ (DSLr), and KITTI360; *Pano* uses Stanford-2D3D. All models use a ViT-L backbone. Missing values (-) indicate the model’s inability to produce the respective output. Metric3D and Metric3Dv2 cannot be evaluated on panoramic images as focal lengths are undefined. †: Requires ground-truth (GT) camera for 3D reconstruction. ‡: Requires GT camera for 2D depth map inference.

Method	S.FoV			S.FoV _{Dist}			L.FoV			Pano		
	$\delta_1^{\text{SSI}} \uparrow$	F _A \uparrow	$\rho_A \uparrow$	$\delta_1^{\text{SSI}} \uparrow$	F _A \uparrow	$\rho_A \uparrow$	$\delta_1^{\text{SSI}} \uparrow$	F _A \uparrow	$\rho_A \uparrow$	$\delta_1^{\text{SSI}} \uparrow$	F _A \uparrow	$\rho_A \uparrow$
DepthAnything [81]	92.2	-	-	94.3	-	-	47.5	-	-	10.4	-	-
DepthAnythingv2 [82]	92.4	-	-	88.9	-	-	48.7	-	-	11.3	-	-
Metric3D ^{†‡} [85]	86.4	43.1	-	88.0	36.7	-	58.7	26.0	-	-	-	-
Metric3Dv2 ^{†‡} [28]	91.1	59.7	-	89.4	47.1	-	69.2	24.7	-	-	-	-
ZoeDepth [†] [7]	88.9	53.3	-	90.3	40.1	-	65.3	6.4	-	32.7	9.9	-
UniDepth [60]	94.9	59.0	85.0	94.0	43.0	70.5	68.6	16.9	19.8	33.0	2.0	1.7
MASt3R [39]	88.0	37.8	80.8	89.9	35.2	<u>77.1</u>	67.1	29.7	25.1	32.3	3.7	2.1
DepthPro [9]	87.4	56.0	79.6	80.6	29.4	71.7	64.5	26.1	32.1	31.8	1.9	1.9
UniK3D-Small	94.3	61.3	85.7	95.1	48.4	73.8	84.5	55.5	70.1	81.3	72.5	53.7
UniK3D-Base	95.5	64.9	86.1	96.5	50.2	75.1	87.4	67.7	79.9	83.6	73.7	53.7
UniK3D-Large	96.1	68.1	89.4	97.3	54.5	78.8	91.2	71.6	81.9	<u>81.4</u>	80.2	57.1

Table 2. **Zero-shot comparison with equirectangular-specialized methods.** All methods are zero-shot tested on Stanford-2D3D [2]. Competing methods are all trained on equirectangular images. Our training set includes Matterport3D [11] with 2% sampling.

Method	Train	$\delta_1 \uparrow$	A.Rel \downarrow
BiFuse [†] [71]	Matterport3D	86.2	12.0
BiFuse++ [†] [72]	Matterport3D	<u>91.4</u>	10.7
UniFuse [†] [29]	Matterport3D	91.3	<u>9.42</u>
UniK3D	Ours	96.8	8.01

The Radial Module first processes the dense encoder features \mathbf{F} through a Transformer Decoder (T-Dec) with 4 parallel layers, one for each resolution, and 1 head. These layers condition \mathbf{F} on the sine-encoded angular representation \mathbf{C} (cf. supplement for details). The conditioned features are then projected onto a 512-channel tensor, forming radial features $\mathbf{D} \in \mathbb{R}^{h \times w \times 512}$. These radial features are afterwards upsampled to the input resolution using residual convolutional blocks and learnable upsampling techniques, *i.e.* bilinear upsampling followed by a single 1×1 convolution. The radial log-scale output $\mathbf{R}_{\log} \in \mathbb{R}^{H \times W}$ is computed from the upsampled features and transformed to \mathbf{R} via element-wise exponentiation. The final 3D spherical output $\mathbf{O} = \mathbf{C} \parallel \mathbf{R}$ is converted to a Cartesian point cloud $\mathbf{O} \in \mathbb{R}^{H \times W \times 3}$ using a spherical-to-Cartesian coordinate transformation. Also, we predict a confidence map (Σ) for the radial outputs by including a second projection head fed with upsampled \mathbf{D} features, besides the first head of the Radial Module which computes \mathbf{R}_{\log} .

Optimization. The optimization process is defined by three different losses. The angular loss \mathcal{L}_{AA} is applied on θ and ϕ separately, with $\mathcal{L}_{AA}^{0.7}$ and $\mathcal{L}_{AA}^{0.5}$ for θ and ϕ , respectively. The final angular loss can be expressed as

$$\mathcal{L}_A(\hat{\mathbf{C}}, \mathbf{C}^*) = \beta \mathcal{L}_{AA}^{0.7}(\hat{\theta}, \theta^*) + (1 - \beta) \mathcal{L}_{AA}^{0.5}(\hat{\phi}, \phi^*), \quad (3)$$

with $(\hat{\cdot})$ and $(\cdot)^*$ serving as prediction and GT identifiers, respectively, and $\beta = 0.75$. It is worth noting that $\mathcal{L}_{AA}^{0.5}$ corresponds to the standard, symmetric L1-loss, as the azimuthal dimension ϕ w.r.t. the principal point is *not* affected by prediction contraction. Our radial loss is the L1-loss between the predicted and GT log-radius obtained by the GT camera and depth: $\mathcal{L}_{\text{rad}} = \left\| \hat{\mathbf{R}}_{\log} - \mathbf{R}_{\log}^* \right\|_1$. The confidence loss is the L1-loss between the detached radial loss and the inverse predicted confidence, Σ : $\mathcal{L}_{\text{conf}} = \left\| \hat{\mathbf{R}}_{\log} - \mathbf{R}_{\log}^* \right\|_1$. The loss is a linear combination of the three losses: $\mathcal{L}_A + \eta \mathcal{L}_{\text{rad}} + \gamma \mathcal{L}_{\text{conf}}$, with $\eta = 2$ and $\gamma = 0.1$.

4. Experiments

Training Datasets. The training dataset accounts for 26 different sources: A2D2 [23], aiMotive [48], Argoverse2 [77], ARKit-Scenes [5], ASE [17], BEDLAM [8], Blended-MVS [83], DL3DV [44], DrivingStereo [79], DynamicReplica [31], EDEN [37], FutureHouse [42], HOI4D [46], HM3D [62], Matterport3D [11], Mapillary-PSD [1], MatrixCity [40], MegaDepth [41], NianticMapFree [3], PointOdyssey [88], ScanNet [12], ScanNet++ (iPhone) [84], TartanAir [75], Taskonomy [87], Waymo [67], and WildRGBD [78]. More details are given in the supplement.

Zero-shot Testing Datasets. We evaluate the generalizability of models by testing them on 13 datasets not seen during training, grouped in 4 different domains which are defined based on their camera type: 1) small FoV (S.FoV), *i.e.* FoV $< 90^\circ$, 2) small FoV with radial and tangential distortions (S.FoV_{Dist}), 3) large FoV (L.FoV), *i.e.* FoV $> 120^\circ$, and 4) Panoramic (Pano) with 360° viewing angle. More specifically, the S.FoV group corresponds to the validation splits of NYU-Depth V2 [50], KITTI Eigen-split [21] and nuScenes [10], and the full IBims-1 [34], ETH-3D [65], and

Table 3. **Ablation on data.** *Data* indicates whether training images include strongly distorted cameras, either from real data or synthesized from pinhole cameras. Output representation: depth.

	Model	Data	S.FoV		S.FoV _{Dist}		L.FoV		Pano	
			F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑
1	Pinhole	✗	55.1	79.2	31.7	60.0	41.2	35.1	8.4	4.2
2	Pinhole	✓	56.1	81.1	40.4	58.2	44.9	43.1	5.9	3.0
3	SH	✗	56.1	79.1	34.5	60.2	47.1	56.7	11.3	16.1
4	SH	✓	56.2	79.4	42.1	62.7	48.5	60.8	10.9	14.8

Diode Indoor [70]; the S.FoV_{Dist} is composed by images artificially distorted from IBims-1, ETH-3D, and Diode Indoor (more details in the supplement); L.FoV is the mix of ADT [55], ScanNet++ (DSLRL) [84], and KITTI360 [43]; and Panoramic (Pano) is to the full Stanford-2D3D [2] dataset.

Evaluation Details. All methods have been re-evaluated with a fair and consistent pipeline. In particular, we do not exploit any test-time augmentations and utilize the same set of weights for all zero-shot evaluations. We use the checkpoint corresponding to the zero-shot model for each method, *i.e.* not fine-tuned on KITTI or NYU. The metrics utilized in the main experiments are δ_1^{SSI} , F_A, and ρ_A. Further metrics are reported in the supplement. δ_1^{SSI} measures scale- and shift-invariant depth estimation performance. F_A is the area under the curve (AUC) of F1-score [54] up to 1/20 of the datasets’ maximum depth and evaluates monocular 3D estimation. ρ_A evaluates the camera performance and is the AUC of the average angular error of camera rays up to 15°, 20°, 30° for S.FoV, L.FoV, and Pano, respectively. We avoid parametric evaluations, such as those based on focal length or FoV, because they lack generality across diverse camera models. Instead, our chosen metrics ensure applicability to any camera type, preserving fairness and consistency in evaluation. The supplement shows the fine-tuning ability of UniK3D by training the final checkpoint on KITTI and NYU-Depth V2 and evaluating in-domain, as per standard practice.

Implementation Details. UniK3D is implemented in PyTorch [57] and CUDA [52]. For training, we use the AdamW [47] optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.999$) with an initial learning rate of 5×10^{-5} . The learning rate is divided by a factor of 10 for the backbone weights for every experiment and weight decay is set to 0.1. We exploit Cosine Annealing as learning rate scheduler to one-tenth starting from 30% of the whole training. We run 250k optimization iterations with a batch size of 128. The training time amounts to 6 days on 16 NVIDIA 4090. The dataset sampling procedure follows a weighted sampler, where the weight of each dataset is its number of scenes. Our augmentations are both geometric and photometric, *i.e.* random resizing and cropping for the former type, and brightness, gamma, saturation, and hue shift for the latter. We randomly sample the image ratio per batch between 2:1 and 9:16. Our ViT [15] backbone is initialized with weights from DINO-pre-trained [53] models. For the ablations, we run 100k training steps with a ViT-S backbone, with training pipeline as for the main experiments.

Table 4. **Ablation on camera model.** *Model* corresponds to the type of camera model for output rays and internal conditioning: pinhole, Zernike-polynomial coefficients, SH coefficients, or non-parametric, *i.e.* predicting one ray per pixel. All experiments are with full data, augmentation, model components, and radial output.

	Model	S.FoV		S.FoV _{Dist}		L.FoV		Pano	
		F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑
1	Pinhole	55.5	79.9	52.5	73.8	45.2	47.9	24.6	16.4
2	Zernike	56.6	80.9	39.9	51.3	49.9	54.6	31.8	17.9
3	Non-Parametric	56.4	81.0	45.2	62.8	42.0	42.8	51.7	20.1
4	SH	57.3	79.8	44.6	59.3	53.5	64.8	58.6	26.3

Table 5. **Ablation on output representation.** *Output* refers to the type of the 3rd dimension of the predicted output space: either Cartesian z-axis depth or spherical radius, *i.e.* distance. All experiments are with full data and augmentation.

	Model	Output	S.FoV		S.FoV _{Dist}		L.FoV		Pano	
			F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑
1	Pinhole	depth	56.1	81.1	40.4	58.2	44.9	43.1	5.9	3.0
2	Pinhole	radius	56.0	81.1	39.5	57.6	44.4	48.9	10.1	4.9
3	SH	depth	56.2	79.4	42.1	62.7	48.5	60.8	10.9	14.8
4	SH	radius	56.8	76.7	35.0	43.7	51.8	61.1	53.8	22.0

4.1. Comparison with The State of The Art

Table 1 presents a comprehensive comparison of UniK3D against existing SotA methods across various FoV and image types. Our model consistently outperforms prior models, especially in challenging large-FoV and panoramic scenarios. For instance, in the L.FoV domain, UniK3D achieves a remarkable δ_1^{SSI} of 91.2% and F_A of 71.6%, outperforming the second-best method by more than 20% and 40%, respectively. This substantial improvement underscores the robustness of our unified spherical framework in handling wide FoVs. In the Pano category, our model’s δ_1^{SSI} and F_A scores of 71.2% and 66.1% also set the new SotA, demonstrating its ability to effectively reconstruct 3D geometry even under extreme camera setups. These results validate that our design choices, including the SH-based camera model and radial output representation, are crucial for maintaining high performance in complex and diverse camera settings.

In addition, Fig. 3 clearly shows how UniK3D can estimate the 3D geometry of scenes from various and distorted cameras. This is in contrast to other methods that fail when facing unconventional or non-pinhole camera images, as depicted by the 2nd, 3rd, and 4th columns. It is important to highlight that Metric3D, Metric3Dv2, and ZoeDepth are evaluated using GT camera parameters for the F_A score, while UniK3D, UniDepth, MAST3R, and DepthPro rely on their predicted cameras. Despite this added difficulty, UniK3D still demonstrates superior 3D reconstruction performance, showcasing its strength in handling real-world conditions where precise camera information is unavailable. Interestingly, our method does not sacrifice performance in more conventional, small-FoV scenarios. UniK3D keeps its top rank, with a δ_1^{SSI} of 94.3 in the S.FoV setting, outperforming previously leading methods. This balance highlights that our advancements in L.FoV representation do not undermine

Table 6. **Ablation on network components.** \mathcal{L}_{AA} indicates if our asymmetric angular loss is used, L1-loss otherwise. *Cond* indicates if our design for enhanced camera conditioning from Sec. 3.2 is utilized. All experiments are with full data and augmentations, radial output representation, and an SH-based camera model.

	\mathcal{L}_{AA}	Cond	S.FoV		S.FoV _{Dist}		L.FoV		Pano	
			$F_A \uparrow$	$\rho_A \uparrow$	$F_A \uparrow$	$\rho_A \uparrow$	$F_A \uparrow$	$\rho_A \uparrow$	$F_A \uparrow$	$\rho_A \uparrow$
1	✗	✗	56.8	76.7	35.0	43.7	51.8	61.1	53.8	22.0
2	✓	✗	57.7	80.9	39.5	52.1	52.9	64.2	56.1	24.4
3	✓	✓	57.3	79.8	44.6	59.3	53.5	64.8	58.6	26.3

the model’s effectiveness for S.FoV tasks. F_A scores remain high in S.FoV and the ρ_A metric shows that our model consistently provides accurate camera parameter estimation.

Moreover, UniK3D is competitive with specialized methods for equirectangular images, as demonstrated in Table 2. This shows how our model can incorporate different scene and camera domains at training time without compromising any domain-specific performance.

4.2. Ablation Studies

Data. Table 3 demonstrates the effect of training on datasets with and without large FoV and camera distortions. Incorporating images with strong camera distortions generally enhances performance across all domains, particularly in challenging cases such as S.FoV with distortion and L.FoV. This underscores the importance of diverse camera geometries in the training set to achieve better generalization. However, the improvement on Pano is limited due to the difficulty of representing panoramic images using a log-depth representation. **Camera Model.** As shown in Table 4, employing SH as the basis for camera rays yields the best overall performance, particularly on L.FoV and Pano. This highlights the effectiveness of our basis function selection in capturing diverse camera models. By contrast, the non-parametric model underperforms in F_A and ρ_A . Since the latter formulation is purely data-driven, we presume that it requires significantly more data to generalize well. It tends to underrepresent the tails of the data distribution, *i.e.* L.FoV and Pano, while performing adequately on more common domains, *i.e.* S.FoV with or without distortion. The Zernike-polynomial basis [18], typically used for modeling lens aberrations, struggles to represent spherical or equirectangular camera geometries due to its inherent planar structure.

Output Space. Table 5 compares different output representations for the third dimension of the predicted space: either the Cartesian z-axis (rows 1 and 3) or the spherical radius (rows 2 and 4). The results show that using the radius representation improves reconstruction metrics in Pano and L.FoV settings, as depth is less effective when dealing with FoVs near or exceeding 180 degrees. This improvement is realized only when the radial component is paired with a camera model capable of representing a wide range of geometries, *e.g.* our SH-based model (row 4 *vs.* row 2). However, the radius-based output space leads to poorer reconstruction for S.FoV with distortion (row 3 *vs.* row 4). This

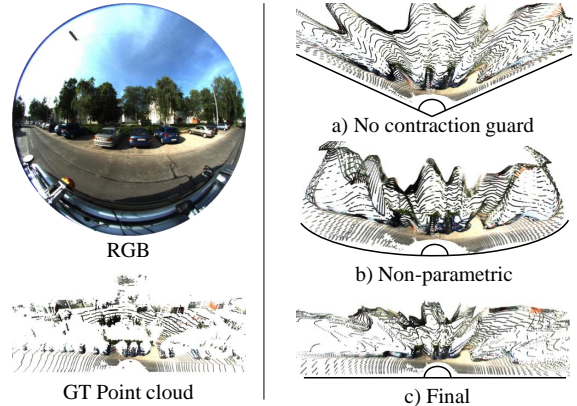


Figure 4. **FoV effects.** The image on the left showcases the challenge of representing the full 180° FoV, alongside the GT point cloud. The effect of FoV contraction occurs when no “guarding”, *i.e.* asymmetric loss (\mathcal{L}_{AA}) and camera conditioning, is put in force, as shown in a). The total absence of any prior may lead to impossible and inconsistent backprojection, as shown in b). The final UniK3D is depicted in c), clearly showing the ability to recover large FoVs with a sensible camera backprojection model.

degradation occurs because the radius representation is more sensitive to minor angular variations, which disproportionately impacts accuracy in small but highly distorted views.

Components. Table 6 examines the impact of our asymmetric angular loss (\mathcal{L}_{AA}) and our strategies designed to enhance camera conditioning. Our full model, which leverages both the asymmetric loss and the improved conditioning (row 3), significantly outperforms those that do not, especially in distorted and L.FoV domains. This demonstrates the efficacy of our combined strategies in preventing contraction in backprojection and improving angular prediction accuracy. The overall gains are rather due to the synergy of combining these contributions. Moreover, these strategies aim at mitigating extreme cases, which may not be easily represented in aggregate quantitative results, but are clearly visible in qualitative samples as in Fig. 4.

5. Conclusion

We have presented UniK3D, the first universal framework for monocular 3D estimation that generalizes seamlessly across diverse camera models, from pinhole to fisheye and panoramic. Our approach introduces strategies to prevent FOV contraction and supports accurate metric 3D estimation through a flexible and robust design for backprojection with any generic camera model. While expanding the diversity and coverage of training data could even further enhance the robustness and applicability of UniK3D, the latter already achieves compelling generalization to unseen cameras and 3D scene domains far beyond the capabilities of the previous state of the art, with only a fair quantity of data.

Acknowledgment. This work is funded by Toyota Motor Europe via the research project TRACE-Zürich.

References

- [1] Manuel Lopez Antequera, Pau Gargallo, Markus Hofinger, Samuel Rota Bulò, Yubin Kuang, and Peter Kotschieder. Mapillary planet-scale depth dataset. In *The European Conference on Computer Vision (ECCV)*, pages 589–604. Springer International Publishing, 2020. [6](#), [15](#)
- [2] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017. [6](#), [7](#), [15](#)
- [3] Eduardo Arnold, Jamie Wynn, Sara Vicente, Guillermo Garcia-Hernando, Áron Monszpart, Victor Adrian Prisacariu, Daniyar Turmukhambetov, and Eric Brachmann. Map-free visual relocalization: Metric pose relative to a single image. In *European Conference on Computer Vision (ECCV)*, 2022. [6](#), [15](#)
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. [13](#)
- [5] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, and Elad Shulman. ARKitscenes - a diverse real-world dataset for 3d indoor scene understanding using mobile RGB-d data. In *Advances in Neural Information Processing Systems (NIPS)*, 2021. [6](#), [15](#)
- [6] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4008–4017, 2020. [2](#), [16](#)
- [7] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. [2](#), [6](#), [14](#), [16](#), [17](#), [18](#), [19](#)
- [8] Michael J. Black, Priyanka Patel, Joachim Tesch, and Jinlong Yang. BEDLAM: A synthetic dataset of bodies exhibiting detailed lifelike animated motion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8726–8737, 2023. [6](#), [15](#)
- [9] Aleksei Bochkovskii, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R Richter, and Vladlen Koltun. Depth pro: Sharp monocular metric depth in less than a second. *arXiv preprint arXiv:2410.02073*, 2024. [1](#), [2](#), [6](#), [14](#), [17](#), [18](#), [19](#)
- [10] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [6](#), [15](#)
- [11] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2017. [6](#), [15](#)
- [12] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. [6](#), [15](#)
- [13] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022. [1](#)
- [14] Xingshuai Dong, Matthew A Garratt, Sreenatha G Anavatti, and Hussein A Abbass. Towards real-time monocular depth estimation for robotics: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(10):16940–16961, 2022. [1](#)
- [15] 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)*. OpenReview.net, 2021. [5](#), [7](#), [13](#)
- [16] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2366–2374. Neural information processing systems foundation, 2014. [2](#)
- [17] Jakob Engel, Kiran Somasundaram, Michael Goesele, Albert Sun, Alexander Gamino, Andrew Turner, Arjang Talattof, Arnie Yuan, Bilal Souti, Brighid Meredith, et al. Project aria: A new tool for egocentric multi-modal ai research. *arXiv preprint arXiv:2308.13561*, 2023. [6](#), [15](#)
- [18] Zernike Frits. Beugungstheorie des schneidenver-fahrens und seiner verbesserten form, der phasenkontrastmethode. *Physica: Nonlinear Phenomena*, 1:689–704, 1934. [8](#)
- [19] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2002–2011, 2018. [2](#)
- [20] Ravi Garg, B. G. Vijay Kumar, Gustavo Carneiro, and Ian Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. *Lecture Notes in Computer Science*, 9912 LNCS:740–756, 2016. [17](#)
- [21] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. [6](#), [15](#)
- [22] Christopher Geyer and Kostas Daniilidis. A unifying theory for central panoramic systems and practical applications. In *The European Conference on Computer Vision (ECCV)*, 2000. [3](#)
- [23] Jakob Geyer, Yohannes Kassahun, Mentar Mahmudi, Xavier Ricou, Rupesh Durgesh, Andrew S. Chung, Lorenz Hauswald, Viet Hoang Pham, Maximilian Mühlegg, Sebastian Dorn, Tiffany Fernandez, Martin Jänicke, Sudesh Mirashi, Chiragkumar Savani, Martin Sturm, Oleksandr Vorobiov, Martin Oelker, Sebastian Garreis, and Peter Schuberth. A2D2: Audi Autonomous Driving Dataset. *arXiv preprint arXiv:2004.06320*, 2020. [6](#), [15](#)
- [24] Vitor Guizilini, Igor Vasiljevic, Dian Chen, Rareş Ambruş, and Adrien Gaidon. Towards zero-shot scale-aware monocular depth estimation. In *Proceedings of the IEEE/CVF In-*

- ternational Conference on Computer Vision (ICCV), pages 9233–9243, 2023. 1, 2
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016-December:770–778, 2015. 13
- [26] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. *CoRR*, abs/1606.08415, 2016. 13
- [27] Noriaki Hirose and Kosuke Tahara. Depth360: Self-supervised learning for monocular depth estimation using learnable camera distortion model. *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 317–324, 2021. 20
- [28] Mu Hu, Wei Yin, Chi Zhang, Zhipeng Cai, Xiaoxiao Long, Hao Chen, Kaixuan Wang, Gang Yu, Chunhua Shen, and Shaojie Shen. Metric3d v2: A versatile monocular geometric foundation model for zero-shot metric depth and surface normal estimation. *arXiv preprint arXiv:2404.15506*, 2024. 2, 6, 14, 16, 17, 18, 19, 21
- [29] Hualie Jiang, Zhe Sheng, Siyu Zhu, Zilong Dong, and Rui Huang. Unifuse: Unidirectional fusion for 360 panorama depth estimation. *IEEE Robotics and Automation Letters (RA-L)*, 6(2):1519–1526, 2021. 6
- [30] Juho Kannala and Sami Sebastian Brandt. A generic camera model and calibration method for conventional, wide-angle, and fish-eye lenses. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 28:1335–1340, 2006. 3, 15
- [31] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Dynamicstereo: Consistent dynamic depth from stereo videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 6, 15
- [32] 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 (CVPR)*, pages 9492–9502, 2024. 1, 2
- [33] Bogdan Khomutenko, Gaëtan Garcia, and Philippe Martinet. An enhanced unified camera model. *IEEE Robotics and Automation Letters (RA-L)*, 1:137–144, 2016. 3, 19
- [34] Tobias Koch, Lukas Liebel, Marco Körner, and Friedrich Fraundorfer. Comparison of monocular depth estimation methods using geometrically relevant metrics on the IBims-1 dataset. *Computer Vision and Image Understanding (CVIU)*, 191:102877, 2020. 6, 15, 16
- [35] Varun Ravi Kumar, Senthil Kumar Yogamani, Markus Bach, Christian Witt, Stefan Milz, and Patrick Mäder. Unrect-depthnet: Self-supervised monocular depth estimation using a generic framework for handling common camera distortion models. *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8177–8183, 2020. 20
- [36] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. *Proceedings of the International Conference on 3D Vision (3DV)*, pages 239–248, 2016. 2
- [37] Hoang-An Le, Thomas Mensink, Partha Das, Sezer Karaoglu, and Theo Gevers. Eden: Multimodal synthetic dataset of enclosed garden scenes. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1579–1589, 2021. 6, 15
- [38] Jin Han Lee, Myung-Kyu Han, Dong Wook Ko, and Il Hong Suh. From big to small: Multi-scale local planar guidance for monocular depth estimation. *CoRR*, abs/1907.10326, 2019. 16
- [39] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. *arXiv preprint arXiv:2406.09756*, 2024. 5, 6, 17, 18, 19, 21
- [40] Yixuan Li, Lihan Jiang, Linning 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 (ICCV)*, pages 3205–3215, 2023. 6, 15
- [41] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2041–2050, 2018. 6, 15
- [42] Zhen Li, Lingli Wang, Xiang Huang, Cihui Pan, and Jiaqi Yang. Phyr: Physics-based inverse rendering for panoramic indoor images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 6, 15
- [43] Yiyi Liao, Jun Xie, and Andreas Geiger. KITTI-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 2022. 7, 15
- [44] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin, Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu, et al. DL3DV-10k: A large-scale scene dataset for deep learning-based 3d vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22160–22169, 2024. 6, 15
- [45] Fayao Liu, Chunhua Shen, Guosheng Lin, and Ian Reid. Learning depth from single monocular images using deep convolutional neural fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 38:2024–2039, 2015. 2
- [46] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan, Hao Shen, Boqiang Liang, Zhoujie Fu, He Wang, and Li Yi. Hoi4d: A 4d egocentric dataset for category-level human-object interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21013–21022, 2022. 6, 15
- [47] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *7th International Conference on Learning Representations, ICLR 2019*, 2017. 7, 15
- [48] Tamas Matuszka, Ivan Barton, Ádám Butykai, Péter Hajas, Dávid Kiss, Domonkos Kovács, Sándor Kunsági-Máté, Péter Lengyel, Gábor Németh, Levente Pető, Dezső Ribli, Dávid Szeghy, Szabolcs Vajna, and Balint Viktor Varga. aimotive dataset: A multimodal dataset for robust autonomous driving with long-range perception. In *International Conference on*

- Learning Representations (ICLR) Workshop on Scene Representations for Autonomous Driving*, 2023. 6, 15
- [49] Christopher Mei and Patrick Rives. Single view point omnidirectional camera calibration from planar grids. *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 3945–3950, 2007. 3, 15
- [50] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *The European Conference on Computer Vision (ECCV)*, 2012. 6, 15
- [51] Jerzy Neyman and Egon Sharpe Pearson. On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 231(694-706):289–337, 1933. 4
- [52] John Nickolls, Ian Buck, Michael Garland, and Kevin Skadron. Scalable parallel programming with cuda: Is cuda the parallel programming model that application developers have been waiting for? *Queue*, 6(2):40–53, 2008. 7
- [53] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 7
- [54] Evin Pinar Örnek, Shristi Mudgal, Johanna Wald, Yida Wang, Nassir Navab, and Federico Tombari. From 2d to 3d: Rethinking benchmarking of monocular depth prediction. *arXiv preprint arXiv:2203.08122*, 2022. 7
- [55] Xiaqing Pan, Nicholas Charron, Yongqian Yang, Scott Peters, Thomas Whelan, Chen Kong, Omkar Parkhi, Richard Newcombe, and Yuheng Carl Ren. Aria digital twin: A new benchmark dataset for egocentric 3d machine perception. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 20133–20143, 2023. 7, 15
- [56] Dennis Park, Rares Ambrus, Vitor Guizilini, Jie Li, and Adrien Gaidon. Is pseudo-lidar needed for monocular 3d object detection? In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 1
- [57] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8024–8035. Curran Associates, Inc., 2019. 7
- [58] Vaishakh Patil, Christos Sakaridis, Alexander Liniger, and Luc Van Gool. P3Depth: Monocular depth estimation with a piecewise planarity prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [59] Luigi Piccinelli, Christos Sakaridis, and Fisher Yu. iDisc: Internal discretization for monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 16
- [60] Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. Unidepth: Universal monocular metric depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10106–10116, 2024. 1, 2, 3, 4, 5, 6, 14, 17, 18, 19, 20, 21
- [61] Luigi Piccinelli, Christos Sakaridis, Yung-Hsu Yang, Mattia Segu, Siyuan Li, Wim Abbeels, and Luc Van Gool. UniDepthV2: Universal monocular metric depth estimation made simpler. *arXiv:2502.20110*, 2025. 1, 2
- [62] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexander Clegg, John M Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI. In *Advances in Neural Information Processing Systems (NIPS)*, 2021. 6, 15
- [63] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 44(3):1623–1637, 2020. 1, 2
- [64] Davide Scaramuzza. Omnidirectional camera. In *Computer Vision, A Reference Guide*, 2014. 3
- [65] Thomas Schöps, Johannes L. Schönberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 6, 15, 16
- [66] Niklaus Simon and Feng Liu. Softmax splatting for video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 16
- [67] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2446–2454, 2020. 6, 15
- [68] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 32–42, 2021. 5
- [69] Vladyslav C. Usenko, Nikolaus Demmel, and Daniel Cremers. The double sphere camera model. *International Conference on 3D Vision (3DV)*, pages 552–560, 2018. 3, 19
- [70] Igor Vasiljevic, Nicholas I. Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z. Dai, Andrea F. Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R. Walter, and Gregory Shakhnarovich. DIODE: A dense indoor and outdoor depth dataset. *CoRR*, abs/1908.00463, 2019. 7, 15, 16
- [71] Fu-En Wang, Yu-Hsuan Yeh, Min Sun, Wei-Chen Chiu, and Yi-Hsuan Tsai. Bifuse: Monocular 360 depth estimation via bi-projection fusion. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 6
- [72] Fu-En Wang, Yu-Hsuan Yeh, Yi-Hsuan Tsai, Wei-Chen Chiu, and Min Sun. Bifuse++: Self-supervised and efficient bi-projection fusion for 360 depth estimation. *IEEE Transactions*

- on *Pattern Analysis and Machine Intelligence (T-PAMI)*, 45 (5):5448–5460, 2022. 6
- [73] Ruicheng Wang, Sicheng Xu, Cassie Dai, Jianfeng Xiang, Yu Deng, Xin Tong, and Jiaolong Yang. Moge: Unlocking accurate monocular geometry estimation for open-domain images with optimal training supervision. *arXiv preprint arXiv:2410.19115*, 2024. 2
- [74] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jérôme Revaud. Dust3r: Geometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20697–20709, 2024. 19
- [75] Wenshan Wang, DeLong Zhu, Xiangwei Wang, Yaoyu Hu, Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Sebastian Scherer. Tartanair: A dataset to push the limits of visual slam. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4909–4916. IEEE, 2020. 6, 15
- [76] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudolidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8445–8453, 2019. 1
- [77] Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh Khandelwal, Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemodel Pontes, Deva Ramanan, Peter Carr, and James Hays. Argoverse 2: Next generation datasets for self-driving perception and forecasting. In *Advances in Neural Information Processing Systems*, 2021. 6, 15
- [78] Hongchi Xia, Yang Fu, Sifei Liu, and Xiaolong Wang. Rgb-d objects in the wild: Scaling real-world 3d object learning from rgb-d videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22378–22389, 2024. 6, 15
- [79] Guorun Yang, Xiao Song, Chaoqin Huang, Zhidong Deng, Jianping Shi, and Bolei Zhou. Drivingstereo: A large-scale dataset for stereo matching in autonomous driving scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 6, 15
- [80] Guanglei Yang, Hao Tang, Mingli Ding, Nicu Sebe, and Elisa Ricci. Transformer-based attention networks for continuous pixel-wise prediction. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16249–16259, 2021. 2
- [81] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10371–10381, 2024. 1, 2, 5, 6, 17, 18, 19
- [82] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv preprint arXiv:2406.09414*, 2024. 2, 6, 14, 16, 17, 18, 19
- [83] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1790–1799, 2020. 6, 15
- [84] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 6, 7, 15
- [85] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9043–9053, 2023. 1, 2, 6, 16, 17, 18, 19
- [86] Weihao Yuan, Xiaodong Gu, ZuoZhuo Dai, Siyu Zhu, and Ping Tan. Neural window fully-connected crfs for monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3906–3915. IEEE, 2022. 2, 16
- [87] Amir R Zamir, Alexander Sax, William B Shen, Leonidas Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2018. 6, 15
- [88] Yang Zheng, Adam W Harley, Bokui Shen, Gordon Wetstein, and Leonidas J Guibas. Pointodyssey: A large-scale synthetic dataset for long-term point tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 19855–19865, 2023. 6, 15
- [89] Brady Zhou, Philipp Krähenbühl, and Vladlen Koltun. Does computer vision matter for action? *Science Robotics*, 4, 2019. 1

Supplementary Material

This supplementary material offers further insights into our work. In Appendix A we describe the network architecture in more detail, necessarily Appendix A overlaps with Sec. 3. Moreover, we analyze the complexity of UniK3D and compare it with other methods in Appendix A.1. Also, we provide further alternatives to our design choices and ablate them in Appendix A.2. Appendix B outlines the training pipeline and hyperparameters chosen in Appendix B.1, altogether with training and validation data in Appendix B.2, and the camera augmentations in Appendix B.3 for completeness and reproducibility. Furthermore, Appendix C provides a more detailed quantitative evaluation with per-dataset evaluation in Appendix C.2. The results corresponding to UniK3D finetuned on KITTI and NYUv2 are reported in Appendix C.1. In Appendix D, we provide answers to possible questions that may arise. Eventually, additional visualizations are provided in Appendix E.

A. Architecture

Encoder. Our model architecture employs a Vision Transformer (ViT) [15] as the encoder, demonstrating its effectiveness across different scales, from Small to Large. The ViT backbones were originally developed for classification tasks, and as such, we modify them by removing the final three layers: the pooling layer, the fully connected layer, and the softmax layer. We extract feature maps and class tokens from the last four layers of the modified ViT backbone. These outputs are flattened and processed using LayerNorm [4] followed by a linear projection layer. The linear layer maps the features and class tokens to a common channel dimension, which is set to 512, 384, and 256 for Large, Base, and Small ViT variants, respectively. Importantly, the normalization and linear layer weights are distinct and are not shared between the different feature resolutions and the class tokens. The dense feature maps are subsequently passed to the Radial Module, while the class tokens are directed to the Angular Module.

Angular Module. The four class tokens extracted from the encoder are first projected to dimensions of $3D$, $3D$, $5D$, and $7D$, respectively. These are then divided into chunks based on the channel dimension d , yielding token groups of size 3, 3, 5, and 7. These token groups serve as the initialization for domain tokens, representing the spherical harmonics (SH) coefficients: 1st-degree, 2nd-degree, and 3rd-degree, respectively. In total, there are 18 tokens (\mathbf{T}), which are processed through two layers of a Transformer Encoder. Each Transformer Encoder layer consists of self-attention with eight heads and a Multi-Layer Perceptron (MLP) that has a single hidden layer of dimension $4C$ and uses the Gaussian Error Linear Unit (GELU) activation function [26]. Both

self-attention and MLP layers include residual connections to improve learning stability. Each of the 18 tokens is then projected to a scalar dimension. The first three tokens specifically define the domain for the spherical harmonics. The first token determines the horizontal field of view (HFov), calculated as $2\pi \cdot \sigma(\mathbf{T}_0)$, where σ denotes the sigmoid function. The second and third tokens represent the poles of the spherical harmonics, *i.e.* the center of projection relative to the image shape, computed as $c_x = \frac{\sigma(\mathbf{T}_1)W}{2}$ and $c_y = \frac{\sigma(\mathbf{T}_2)H}{2}$, respectively, where H and W are the image height and width. The vertical FoV is derived under the assumption of square pixels: $\text{HFov} \times \frac{H}{W}$. Using this domain definition, we compute the spherical harmonics up to the 3rd degree, excluding the constant component, yielding 15 harmonic tensors of size $\mathbb{R}^{H \times W \times 3}$. The pencil of rays \mathbf{C} is then constructed as a linear combination of these harmonics and the corresponding 15 processed tokens ($\mathbf{T}_{3:18}$).

Radial Module. The sine-encoded camera rays \mathbf{C} are used to condition each resolution level of the dense feature maps \mathbf{F} via a Transformer Decoder layer. In this setup, the dense features \mathbf{F} serve as the *query*, while the sine-encoded camera rays provide the *keys* and *values*. The cross-attention mechanism includes a residual connection without any learnable gain factors, such as LayerScale. The conditioned features are then refined in a Feature Pyramid Network (FPN) manner: the deepest features are processed through two Residual Convolution blocks [25], followed by bilinear upsampling and a projection step that halves the channel dimension. These upsampled features are then combined with the features from the next layer, which are similarly projected to match channel dimension and upsampled using a single 2×2 transposed convolution. This process continues until all remaining three feature maps are consumed. The final output features are upsampled to the input image resolution and projected to a single-channel dimension, yielding the log-radius \mathbf{R}_{\log} . The same projection, architectural-wise but with separate weights, is used to generate the log-confidence Σ_{\log} . The final radius and confidence values are obtained by exponentiating these tensors element-wise, transforming them from log-space to the original space.

A.1. Complexity

We perform a detailed analysis of the computational cost of UniK3D, presented in Table 7, and compare it to other state-of-the-art methods. To ensure a fair and consistent comparison, we use input sizes that are as similar as possible across all models. However, this approach introduces certain challenges. DepthPro, for instance, has an entangled and multi-resolution architecture, which complicates tuning the input size consistently across methods. Its architectural design does not easily allow for adjustments, making it difficult to

Table 7. **Parameters and efficiency comparison.** Comparison of performance of methods based on input size, latency, and number of trainable parameters. Tested on RTX3090 GPU, 16-bit precision float, and synchronized timers. The last two rows correspond to the Angular and Radial Modules evaluated independently. All models are based on ViT-L backbone.

Method	Input Size	Latency (ms)	Parameters (M)
ZoeDepth [7]	512 × 512	144.8	345.9
DepthAnything v2 [82]	518 × 518	78.1	334.7
UniDepth [60]	518 × 518	146.4	347.0
Metric3Dv2 [28]	518 × 518	135.6	441.9
MASt3R [28]	512 × 512	154.7	668.6
DepthPro [9]	1536 × 1536	808.1	952.0
UniK3D	518 × 518	88.4	358.8
Radial Module	-	21.9	38.2
Angular Module	-	3.1	12.1

align with a standardized input size. Additionally, the performance of models like DepthPro and Metric3D, as evaluated in our main experiments in Sec. 4, shows a significant drop when tested with image shapes that differ from those used during training. This sensitivity highlights a fundamental limitation: these methods are heavily optimized for specific image resolutions, and deviations from these resolutions can lead to substantial performance degradation. Consequently, while we strive to measure computation under the most equitable conditions, it is essential to acknowledge that these models are not well-suited for resolutions that differ from their training setup. In contrast, UniK3D is designed to be flexible w.r.t. image shape, maintaining robust performance across different resolutions. For our experiments, we chose the same input shape as DepthAnything v2, as it provides a balanced trade-off between computational efficiency and performance. Furthermore, to account for the asynchronous nature of CUDA kernel threading, we ensure precise inference time measurements by enabling proper synchronization and utilizing CUDA event recording. This approach guarantees an accurate reflection of computational cost, avoiding any misrepresentation caused by asynchronous operations. As shown in Table 7, UniK3D is among the most efficient models. The primary differences in computational cost, especially when compared to DepthAnything v2, stem from the inclusion of our Angular Module and Scale components. These components are essential for our model to handle absolute metric depth and camera-specific adjustments, features that relative depth estimation networks do not require. Despite this additional complexity, our model’s efficiency remains competitive, underscoring its design’s effectiveness in addressing diverse camera geometries while maintaining high performance.

A.2. Architectural Alternatives

Despite the camera conditioning has been proven superior in UniDepth [60], we ablate alternative architectural choices for both the Transformer Encoder and Decoder components. In particular, we have chosen the most typical alternatives

Table 8. **Ablation on camera conditioning design.** *Camera Cond.* corresponds to the type of camera conditioning employed to condition the depth features with camera ones. *Add* refers to a simple addition in the feature space. *Cat* represents a simple concatenation and projection from $2C$ to C channel dimension. *Prompt* is our attention-based conditioning.

Camera Cond.	S.FoV		S.FoV _{Dist}		L.FoV		Pano	
	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑
1 Add	53.0	78.9	26.3	41.6	45.0	58.5	42.5	18.2
2 Cat	54.7	79.0	28.7	44.6	46.6	58.1	42.3	18.1
3 Prompt	57.3	79.8	44.6	59.3	53.5	64.8	58.6	26.3

Table 9. **Ablation on camera tokens processing.** *T-Enc.* indicates if the camera tokens are processed in the Angular Module either via the transformer encoder layer or not, in the latter case the tokens are fed directly to the final projections.

T-Enc	S.FoV		S.FoV _{Dist}		L.FoV		Pano	
	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑	F _A ↑	ρ _A ↑
1 ✗	55.7	77.3	43.2	56.6	50.9	63.2	54.9	20.7
2 ✓	57.3	79.8	44.6	59.3	53.5	64.8	58.6	26.3

for conditioning: a simple addition or concatenation in place. While the camera tokens processing “alternative” involves an identity that shortcuts the camera tokens to the final projection layers. Table 9 shows how the camera tokens processing, via the encoder layer, does not present large changes, showing how the class tokens from different layers are already informative. However, Table 8 clearly shows how the simpler conditioning alternatives, such as addition or concatenation, underperform our attention-based conditioning. This highlights how conditioning plays an important role in final performance and how strongly designed conditioning is paramount to achieving proper generalization.

B. Training Details

B.1. Hyperparameters.

The training parameters, *i.e.* those for optimization, scheduling, and augmentations, are described in Table 10. The losses utilized, with the input and corresponding weights, are outlined in Table 11.

B.2. Data

Details of training and validation datasets are presented in Table 12 and Table 13.

Training Datasets. The datasets utilized for training are a mixture of different cameras and domains as shown in Table 12. The sequence-based datasets are sub-sampled during collection in a way that the interval between two consecutive frames is not smaller than half a second. No post-processing is applied. The total amount of training samples accounts for more than 8M samples. The datasets are sampled in each batch with a probability corresponding to the values in *Sampling* column in Table 12. This probability is related to the number of scenes present in each dataset. However, probabilities are changed based on a simple qualitative data in-

Table 10. **Training Hyperparameters.** All training hyperparameters with corresponding values are presented.

Hyperparameter	Value
Steps	250k
Batch Size	128
LR	$5 \cdot 10^{-5}$
LR Encoder	$5 \cdot 10^{-6}$
Optimizer	AdamW [47]
(β_1, β_2)	(0.9, 0.999)
Weight Decay	0.1
Gradient Clip Norm	1.0
Precision	16-bit Float
LR Scheduler	Cosine to 0.1 start after 75k iters
EMA	0.9995 start after 75k iters
Color jitter prob	80%
Color jitter intensity	[0.0, 0.5]
Gamma prob	80%
Gamma intensity	[0.5, 1.5]
Horizontal flip prob	50%
Greyscale prob	20%
Gaussian blur prob	20%
Gaussian blur sigma	[0.1, 2.0]
Random zoom	[0.5, 2.0]
Random translation	[-0.05, 0.05]
Image ratio	[1 : 2, 2 : 1]
Resolution	0.28MP [0.2MP, 0.6MP] last 50k iters

Table 11. **Training Losses.** Training losses with corresponding weight and input.

Loss	Inputs	Weight	Parameters
L1	Radius (log)	2.0 (η)	-
L1-asymmetric	Polar	0.75	$\alpha = 0.7$
L1	Azimuth	0.25	-
L1	Confidence (log), Radius error (detached)	0.1 (γ)	-

spection, such that the most diverse datasets are sampled more. Most of the datasets involve pinhole images or rectified cameras, *e.g.* MegaDepth [41] or NianticMapFree [3], other datasets provide only the pinhole calibration despite being clearly distorted, *i.e.* Mapillary [1], there the entire samples are masked out in the camera loss computation.

Validation Datasets. Table 13 presents all the validation datasets and splits them into 3 groups: small FoV, large FoV, and Panoramic. As per standard practice, KITTI Eigen-split corresponds to the corrected and accumulated GT depth maps with 45 images with inaccurate GT discarded from the original 697 images. The small FoV with distortion

Table 12. **Training Datasets.** List of the validation datasets: number of images, scene type, acquisition method, and sampling frequency are reported. SfM: Structure-from-Motion. MVS: Multi-View Stereo. Syn: Synthetic. Rec: Mesh reconstruction. KB: Kannala-Brandt [30]. Equi: Equirectangular

Dataset	Images	Scene	Acquisition	Camera	Sampling
A2D2 [23]	78k	Outdoor	LiDAR	Pinhole	2.5%
aiMotive [48]	178k	Outdoor	LiDAR	Mei [49]	0.3%
Argoverse2 [77]	403k	Outdoor	LiDAR	Pinhole	7.6%
ARKit-Scenes [5]	1.75M	Indoor	LiDAR	Pinhole	1.3%
ASE [17]	2.72M	Indoor	Syn	Fisheye624	10.1%
BEDLAM [8]	24k	Various	Syn	Pinhole	2.0%
BlendedMVS [83]	114k	Outdoor	MVS	Pinhole	2.5%
DL3DV [44]	306k	Outdoor	SfM	KB [30]	4.7%
DrivingStereo [79]	63k	Outdoor	MVS	Pinhole	2.5%
DynamicReplica [31]	120k	Indoor	Syn	Pinhole	1.3%
EDEN [37]	368k	Outdoor	Syn	Pinhole	2.5%
FutureHouse [42]	28.3	Indoor	Syn	Equi	2.5%
HOI4D [46]	59k	Egocentric	RGB-D	KB [30]	1.7%
HM3D [62]	540k	Indoor	Rec	Pinhole	5.2%
Matterport3D [11]	10.8k	Indoor	Rec	Equi	2.0%
Mapillary PSD [1]	742k	Outdoor	SfM	Pinhole	2.0%
MatrixCity [40]	190k	Outdoor	Syn	Pinhole	5.0%
MegaDepth [41]	273k	Outdoor	SfM	Pinhole	8.0%
NianticMapFree [3]	25k	Outdoor	SfM	Pinhole	2.0%
PointOdyssey [88]	33k	Various	Syn	Pinhole	1.7%
ScanNet [12]	83k	Indoor	RGB-D	Pinhole	5.0%
ScanNet++ [84]	39k	Indoor	Rec	Pinhole	3.0%
TartanAir [75]	306k	Various	Syn	Pinhole	5.5%
Taskonomy [87]	1.94M	Indoor	RGB-D	Pinhole	6.0%
Waymo [67]	223k	Outdoor	LiDAR	Pinhole	7.5%
WildRGBD [78]	1.35M	Indoor	RGB-D	Pinhole	7.5%

Table 13. **Validation Datasets.** List of the validation datasets: number of images, scene type, acquisition method, and max evaluation distance are reported. 1st group: small FoV, 2nd group: large FoV, 3rd: Panoramic. Rec: Mesh reconstruction.

Dataset	Images	Scene	Acquisition	Max Distance
KITTI [21]	652	Outdoor	LiDAR	80.0
NYU [50]	654	Indoor	RGB-D	10.0
IBims-1 [34]	100	Indoor	RGB-D	25.0
Diode [70]	325	Indoor	LiDAR	25.0
ETH3D [65]	454	Outdoor	RGB-D	50.0
NuScenes [10]	3.6k	Outdoor	LiDAR	80.0
ScanNet++ [84]	779	Indoor	Rec	10.0
ADT [55]	469	Indoor	Rec	20.0
KITTI360 [43]	527	Outdoor	LiDAR	80.0
Stanford-2D3D [2]	1413	Indoor	Rec	10.0

presented in Sec. 3 and used for evaluation is obtained based on synthesized cameras from ETH3D, Diode (Indoor), and IBims-1, all distorted images and cameras are manually checked for realism, after being generated with the pipeline presented in Appendix B.3.

B.3. Camera Augmentations

To address the limited diversity of distorted camera data, we augment images captured with pinhole cameras by artificially deforming them, thereby simulating images from distorted camera models, *e.g.* Fisheye624 or radial Kannala-Brandt [30]. The augmentation process involves two main steps. First, we compute a deformation field. This starts with unprojecting the 2D depth map obtained from a pinhole camera into a 3D point cloud. We then project these 3D points onto the image plane of a randomly sampled distorted

Table 14. **Camera Sampling for S.FoV_{Dist} generation.** The parameters to generate S.FoV_{Dist} images are listed. We employed different camera models with different parameter ranges. The sampling is uniform sampling within the ranges. The seed is 13.

Model	Probability	Parameter	Range
EUCM	0.1	α β	[0, 1] [0.25, 4]
Fisheye624	0.35	$\{k_i\}_{i=1}^6$ $\{t_i\}_{i=1}^2$ $\{s_i\}_{i=1}^4$	[0.6, 0.8] [-0.01, 0.01] [-0.01, 0.01]
Fisheye624	0.35	$\{k_i\}_{i=1}^6$ $\{t_i\}_{i=1}^2$ $\{s_i\}_{i=1}^4$	[-0.6, -0.4] [-0.01, 0.01] [-0.01, 0.01]
Fisheye624	0.2	$\{k_i\}_{i=1}^6$ $\{t_i\}_{i=1}^2$ $\{s_i\}_{i=1}^4$	[-0.2, 0.2] [-0.05, 0.05] [-0.05, 0.05]

Table 15. **Camera Sampling for Camera Augmentation.** The parameters to generate an augmented camera during training images are listed. We employed different camera models with different parameter ranges. The sampling is uniform sampling within the ranges. When some parameters are not listed, *e.g.* $\{k_i\}_{i=4}^6$ for Kannala-Brandt model, they are set to 0.

Model	Probability	Parameter	Range
EUCM	0.1	α β	[0, 1] [0.25, 4]
Fisheye624	0.15	$\{k_i\}_{i=1}^6$ $\{t_i\}_{i=1}^2$ $\{s_i\}_{i=1}^4$	[0.1, 0.5] [-0.005, 0.005] [-0.01, 0.01]
Fisheye624	0.15	$\{k_i\}_{i=1}^6$ $\{t_i\}_{i=1}^2$ $\{s_i\}_{i=1}^4$	[-0.5, -0.1] [-0.005, 0.005] [-0.01, 0.01]
Kannala-Brandt	0.2	$\{k_i\}_{i=1}^3$ $\{t_i\}_{i=1}^2$	[-0.05, 0.05] [-0.02, 0.02]
Kannala-Brandt	0.4	$\{k_i\}_{i=1}^3$ $\{t_i\}_{i=1}^2$	[-0.5, 0.5] [-0.001, 0.001]

camera model to obtain the new 2D coordinates. The deformation field is defined as the distance between the original 2D image coordinates and the newly projected 2D coordinates. This flow indicates how the original image should be warped to mimic the appearance of a distorted camera view. Next, we warp the image using softmax-based splatting [66], a technique that projects pixels based on the computed deformation field while preserving image details. To ensure the warping process does not create artifacts like holes, we use an ‘‘importance’’ metric, which is the inverse of the depth value for each pixel. This metric prioritizes closer points, ensuring that details and correct parallax are maintained during the warping. For non-synthetic images, where ground-truth depth maps are unavailable, we generate depth predictions

Table 16. **Comparison on NYU validation set.** All models are trained on NYU. The first four are trained only on NYU. The last four are fine-tuned on NYU.

Method	δ_1	δ_2	δ_3	A.Rel RMS Log ₁₀		
	<i>Higher is better</i>			<i>Lower is better</i>		
BTS [38]	88.5	97.8	99.4	10.9	0.391	0.046
AdaBins [6]	90.1	98.3	99.6	10.3	0.365	0.044
NeWCRF [86]	92.1	99.1	<u>99.8</u>	9.56	0.333	0.040
iDisc [59]	93.8	99.2	<u>99.8</u>	8.61	0.313	0.037
ZoeDepth [7]	95.2	99.5	<u>99.8</u>	7.70	0.278	0.033
Metric3Dv2 [28]	98.9	99.8	100	<u>4.70</u>	<u>0.183</u>	<u>0.020</u>
DepthAnythingv2 [82]	<u>98.4</u>	99.8	100	5.60	0.206	0.024
UniK3D	98.9	99.8	100	4.43	0.173	0.019

Table 17. **Comparison on KITTI Eigen-split validation set.** All models are trained on KITTI E-ign-split training and tested on the corresponding validator split. The first are trained only on KITTI. The last 4 are fine-tuned on KITTI.

Method	δ_1	δ_2	δ_3	A.Rel RMS RMS _{log}		
	<i>Higher is better</i>			<i>Lower is better</i>		
BTS [38]	96.2	99.4	99.8	5.63	2.43	0.089
AdaBins [6]	96.3	99.5	99.8	5.85	2.38	0.089
NeWCRF [86]	97.5	<u>99.7</u>	<u>99.9</u>	5.20	2.07	0.078
iDisc [59]	97.5	<u>99.7</u>	<u>99.9</u>	5.09	2.07	0.077
ZoeDepth [7]	96.5	99.1	99.4	5.76	2.39	0.089
Metric3Dv2 [85]	<u>98.5</u>	99.8	100	<u>4.40</u>	1.99	0.064
DepthAnythingv2 [82]	98.3	99.8	100	4.50	<u>1.86</u>	<u>0.067</u>
UniK3D	99.0	99.8	<u>99.9</u>	3.69	1.68	0.060

in an inference-only mode to compute the deformation. To ensure these predictions are accurate enough to create realistic deformations, we apply this augmentation only after the model has been trained for 10,000 steps. By this point, the model has learned a decently reliable (scale-invariant) depth representation. The specific camera parameters used to sample the new random camera are listed in Table 15.

Validation datasets generation. Generating validation datasets for testing models on distorted images with reduced fields of view presents an additional challenge, as most distortions are typically associated with large fields of view. To simulate this, we use synthetic camera parameters to deform RGB images from datasets such as ETH3D [65], IBims-1 [34], and Diode (Indoor) [70]. These datasets are chosen because they provide nearly complete ground-truth depth maps, making the deformation process well-posed and realistic. Any small gaps or holes in the depth maps are filled using inpainting. Importantly, the 3D ground-truth data remains unchanged, as it is invariant to the camera model used. To ensure realism, we manually validate each deformed image and will release both the code for data generation and the resulting validation data.

C. Additional Quantitative Results

C.1. Fine-tuning

We evaluate the fine-tuning capability of UniK3D by resuming training with either KITTI or NYU as the sole training dataset. The fine-tuning process starts from the weights and

optimizer states obtained after the large-scale pretraining phase, ensuring a fair and consistent initialization. The standard SILog loss is used as the training objective, with a batch size of 16, and the model is trained for an additional 40,000 steps. To focus the evaluation on the impact of in-domain data, we disable all augmentations except for horizontal flipping and omit the asymmetric component of the angular loss during fine-tuning. For evaluation, we adhere to the standard practices for both datasets to ensure comparability with prior work. KITTI results are reported using the Garg [20] evaluation crop, and the maximum evaluation depths for KITTI and NYU are set to 80 and 10 meters, respectively. Importantly, we do not apply any test-time augmentations or tuning, such as varying the input size, to maintain consistency and avoid introducing additional confounding factors. Our results demonstrate that UniK3D benefits significantly from in-domain fine-tuning. Table 17 highlights the model’s ability to perform exceptionally well on highly structured and calibrated datasets like KITTI, even though UniK3D is inherently designed for flexibility and cross-domain generalization. This suggests that the model can effectively adapt to well-structured data when fine-tuned. This fine-tuning analysis highlights the adaptability of UniK3D to diverse settings while maintaining its primary design focus on flexibility. Similarly, Table 16 shows that UniK3D remains competitive when fine-tuned on less structured domains like NYU, which represent typical indoor environments. These results reinforce the importance of in-domain data for achieving optimal performance, particularly on datasets with distinct properties or domain-specific challenges. In addition, the results underline the robustness of our model, as it achieves strong performance across significantly different dataset characteristics.

Table 18. **Comparison on zero-shot evaluation for NYUv2.** Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel \downarrow	RMSE \downarrow	$\delta_1^{\text{SSI}} \uparrow$	F _A \uparrow	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	97.8	-	-
DepthAnythingv2 [82]	-	-	-	97.7	-	-
Metric3D [‡] [85]	68.1	44.2	1.23	89.0	-	-
Metric3Dv2 [‡] [28]	93.4	9.1	0.399	98.1	-	-
ZoeDepth [†] [7]	94.2	8.2	0.305	98.0	-	-
UniDepth [60]	98.0	7.3	0.230	99.0	83.1	99.2
MASt3R [39]	83.9	13.0	0.435	94.8	69.6	90.7
DepthPro [9]	92.2	10.1	0.357	97.2	73.0	<u>93.1</u>
UniK3D-Small	90.4	11.2	0.351	97.4	69.1	83.0
UniK3D-Base	93.1	10.3	0.325	97.9	75.4	89.1
UniK3D-Large	<u>96.5</u>	<u>7.4</u>	<u>0.259</u>	<u>98.2</u>	<u>82.5</u>	91.2

C.2. Per-dataset Evaluation

We present results for each of the validation datasets independently in Table 18 (NYUv2), Table 19 (KITTI), Table 20 (IBims-1), Table 21 (ETH3D), Table 22 (Diode Indoor), Table 23 (nuScenes), Table 24 (IBims-1_{Dist}), Table 25

Table 19. **Comparison on zero-shot evaluation for KITTI.** Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel \downarrow	RMSE \downarrow	$\delta_1^{\text{SSI}} \uparrow$	F _A \uparrow	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	88.5	-	-
DepthAnythingv2 [82]	-	-	-	88.4	-	-
Metric3D [‡] [85]	3.3	49.8	10.35	97.0	-	-
Metric3Dv2 [‡] [28]	2.3	56.3	12.81	96.7	-	-
ZoeDepth [†] [7]	<u>93.6</u>	<u>8.2</u>	<u>3.24</u>	96.7	-	-
UniDepth [60]	98.0	4.8	2.14	98.3	85.8	97.5
MASt3R [39]	2.8	58.2	11.88	90.9	10.9	77.7
DepthPro [9]	78.2	17.2	5.27	94.8	62.4	80.9
UniK3D-Small	92.1	11.6	3.76	96.4	<u>77.7</u>	<u>85.6</u>
UniK3D-Base	93.1	12.6	3.84	<u>97.3</u>	76.6	82.7
UniK3D-Large	81.2	17.4	4.77	96.8	71.4	79.3

Table 20. **Comparison on zero-shot evaluation for IBims-1.** Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel \downarrow	RMSE \downarrow	$\delta_1^{\text{SSI}} \uparrow$	F _A \uparrow	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	97.0	-	-
DepthAnythingv2 [82]	-	-	-	98.0	-	-
Metric3D [‡] [85]	75.1	19.3	0.633	96.2	-	-
Metric3Dv2 [‡] [28]	68.4	20.7	0.700	98.8	-	-
ZoeDepth [†] [7]	49.8	21.5	0.989	95.8	-	-
UniDepth [60]	15.7	41.0	1.25	98.1	30.3	76.6
MASt3R [39]	61.0	19.7	0.883	95.1	55.7	<u>76.0</u>
DepthPro [9]	82.3	17.0	0.573	98.0	62.8	75.9
UniK3D-Small	<u>87.7</u>	13.0	0.484	97.7	67.3	74.6
UniK3D-Base	87.6	<u>12.5</u>	<u>0.452</u>	98.0	<u>67.5</u>	73.4
UniK3D-Large	91.9	10.4	0.406	<u>98.5</u>	69.8	75.4

Table 21. **Comparison on zero-shot evaluation for ETH3D.** Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel \downarrow	RMSE \downarrow	$\delta_1^{\text{SSI}} \uparrow$	F _A \uparrow	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	93.2	-	-
DepthAnythingv2 [82]	-	-	-	93.3	-	-
Metric3D [‡] [85]	19.7	136.8	10.45	81.1	-	-
Metric3Dv2 [‡] [28]	90.0	12.7	1.85	89.7	-	-
ZoeDepth [†] [7]	33.8	54.7	3.45	86.1	-	-
UniDepth [60]	18.5	53.3	3.50	93.9	27.6	42.6
MASt3R [39]	21.4	45.3	4.43	91.3	28.4	92.2
DepthPro [9]	39.7	65.2	36.31	81.1	41.2	77.4
UniK3D-Small	53.6	60.0	4.89	94.2	44.3	80.7
UniK3D-Base	68.4	28.5	3.77	<u>95.8</u>	53.8	<u>82.0</u>
UniK3D-Large	<u>68.7</u>	<u>23.6</u>	<u>2.63</u>	95.9	<u>53.6</u>	81.3

(ETH3D_{Dist}), Table 26 (Diode Indoor_{Dist}), Table 27 (ScanNet++ DSLR), Table 28 (ADT), and Table 29 (KITTI360). Note that we do not report results for the “Pano” group, as it only consists of a single dataset, Stanford-2D3D. Our results show that performance on pinhole camera models has reached a saturation point, yet UniK3D achieves the highest average metric overall, even though it does not always rank first on every individual dataset. This demonstrates the strong generalization ability of UniK3D, attributed to its flex-

Table 22. **Comparison on zero-shot evaluation for Diode (Indoor)**. Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel ↓	RMSE ↓	$\delta_1^{SSI} \uparrow$	F _A ↑	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	97.5	-	-
DepthAnythingv2 [82]	-	-	-	97.6	-	-
Metric3D ^{†‡} [85]	40.4	61.1	2.34	91.3	-	-
Metric3Dv2 ^{†‡} [28]	94.0	9.3	0.399	98.5	-	-
ZoeDepth [†] [7]	34.9	33.6	2.07	91.8	-	-
UniDepth [60]	<u>76.2</u>	17.2	0.954	97.2	63.0	96.1
MASt3R [39]	52.6	27.9	1.68	92.3	48.8	70.2
DepthPro [9]	67.1	19.9	0.900	93.9	50.3	71.5
UniK3D-Small	57.2	21.4	0.968	96.1	49.3	<u>92.5</u>
UniK3D-Base	55.1	19.6	0.859	97.4	50.1	91.2
UniK3D-Large	71.3	<u>16.1</u>	<u>0.767</u>	<u>97.9</u>	<u>53.8</u>	79.5

Table 23. **Comparison on zero-shot evaluation for nuScenes**. Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel ↓	RMSE ↓	$\delta_1^{SSI} \uparrow$	F _A ↑	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	79.0	-	-
DepthAnythingv2 [82]	-	-	-	79.4	-	-
Metric3D ^{†‡} [85]	75.4	23.7	8.94	64.0	-	-
Metric3Dv2 ^{†‡} [28]	84.1	23.6	9.40	64.8	-	-
ZoeDepth [†] [7]	33.8	42.0	<u>7.41</u>	64.8	-	-
UniDepth [60]	<u>84.6</u>	12.7	4.56	83.1	<u>64.4</u>	<u>97.7</u>
MASt3R [39]	2.7	65.6	13.76	63.5	13.6	78.3
DepthPro [9]	56.6	28.7	11.29	59.1	46.5	79.1
UniK3D-Small	80.9	18.9	8.43	83.8	59.4	95.8
UniK3D-Base	84.9	<u>16.7</u>	9.15	<u>86.7</u>	65.5	97.8
UniK3D-Large	84.0	18.9	10.83	87.0	60.3	86.9

Table 24. **Comparison on zero-shot evaluation for IBims-1Dist**. Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel ↓	RMSE ↓	$\delta_1^{SSI} \uparrow$	F _A ↑	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	97.1	-	-
DepthAnythingv2 [82]	-	-	-	93.4	-	-
Metric3D ^{†‡} [85]	56.8	26.5	0.947	93.3	-	-
Metric3Dv2 ^{†‡} [28]	61.3	22.1	0.940	93.3	-	-
ZoeDepth [†] [7]	30.0	28.0	1.28	94.5	-	-
UniDepth [60]	48.7	23.0	0.966	97.2	53.3	69.3
MASt3R [39]	31.8	31.9	1.30	92.8	44.1	69.7
DepthPro [9]	27.2	47.6	1.86	83.0	32.4	69.5
UniK3D-Small	<u>67.2</u>	<u>17.1</u>	0.726	97.6	<u>62.6</u>	71.5
UniK3D-Base	66.0	17.9	<u>0.695</u>	<u>98.3</u>	59.8	<u>72.7</u>
UniK3D-Large	70.9	15.0	0.615	98.6	67.9	77.3

ible design and large-scale training, which enables robust performance across diverse domains without overfitting to any specific one. We report additional and more typical metrics such as absolute relative error as A.Rel as a percentage and root-means-squared error RSME using meter as unit.

Table 25. **Comparison on zero-shot evaluation for ETH3D_{Dist}**. Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel ↓	RMSE ↓	$\delta_1^{SSI} \uparrow$	F _A ↑	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	91.8	-	-
DepthAnythingv2 [82]	-	-	-	83.9	-	-
Metric3D ^{†‡} [85]	19.6	123.6	11.05	80.9	-	-
Metric3Dv2 ^{†‡} [28]	42.8	104.3	9.87	83.5	-	-
ZoeDepth [†] [7]	25.4	45.9	4.12	86.1	-	-
UniDepth [60]	27.6	43.8	4.69	90.1	38.5	67.5
MASt3R [39]	14.6	51.8	5.37	87.7	32.0	<u>78.5</u>
DepthPro [9]	16.1	72.8	18.77	72.7	29.1	69.9
UniK3D-Small	42.1	125.3	12.14	92.9	49.9	68.4
UniK3D-Base	<u>47.9</u>	<u>36.5</u>	<u>3.54</u>	<u>95.1</u>	<u>53.5</u>	67.1
UniK3D-Large	67.0	22.1	2.75	95.5	63.6	83.1

Table 26. **Comparison on zero-shot evaluation for Diode_{Dist} (Indoor)**. Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel ↓	RMSE ↓	$\delta_1^{SSI} \uparrow$	F _A ↑	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	94.2	-	-
DepthAnythingv2 [82]	-	-	-	89.3	-	-
Metric3D ^{†‡} [85]	26.4	124.0	4.08	89.7	-	-
Metric3Dv2 ^{†‡} [28]	34.1	35.2	1.61	91.6	-	-
ZoeDepth [†] [7]	24.0	39.8	2.32	90.1	-	-
UniDepth [60]	30.2	34.8	1.85	94.7	37.2	74.8
MASt3R [39]	20.6	46.0	2.41	89.3	29.5	83.0
DepthPro [9]	24.7	56.5	2.31	86.0	26.5	75.7
UniK3D-Small	27.6	33.4	1.48	95.0	33.0	82.6
UniK3D-Base	<u>31.6</u>	<u>30.0</u>	<u>1.35</u>	<u>96.1</u>	<u>37.0</u>	<u>85.1</u>
UniK3D-Large	26.9	30.0	1.33	97.5	36.1	85.4

Table 27. **Comparison on zero-shot evaluation for ScanNet++ (DSLRL)**. Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel ↓	RMSE ↓	$\delta_1^{SSI} \uparrow$	F _A ↑	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	51.4	-	-
DepthAnythingv2 [82]	-	-	-	52.3	-	-
Metric3D ^{†‡} [85]	16.5	180.5	1.83	51.2	-	-
Metric3Dv2 ^{†‡} [28]	5.2	237.0	2.51	71.3	-	-
ZoeDepth [†] [7]	2.0	158.5	1.45	71.2	-	-
UniDepth [60]	0.6	162.9	1.59	71.0	9.1	20.2
MASt3R [39]	5.8	114.8	1.07	73.0	21.0	16.6
DepthPro [9]	9.6	95.8	0.928	74.1	24.4	30.9
UniK3D-Small	6.2	92.8	0.931	78.1	23.5	35.1
UniK3D-Base	<u>55.4</u>	<u>33.1</u>	<u>0.340</u>	<u>86.6</u>	<u>53.9</u>	<u>65.1</u>
UniK3D-Large	65.1	25.3	0.285	90.8	59.1	70.0

D. Q&A

Here we list possible questions that might arise after reading the paper. We structure the section in a discursive question-and-answer fashion.

- **What is the importance of data for generalization w.r.t. scene scale?**

Data diversity is crucial for generalizing depth estimation, especially for monocular methods that heavily rely on semantic cues and are sensitive to domain gaps. Scale predic-

Table 28. **Comparison on zero-shot evaluation for ADT.** Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel \downarrow	RMSE \downarrow	$\delta_1^{SSI} \uparrow$	F _A \uparrow	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	81.7	-	-
DepthAnythingv2 [82]	-	-	-	82.6	-	-
Metric3D ^{†‡} [85]	72.5	26.2	0.560	85.3	-	-
Metric3Dv2 ^{†‡} [28]	75.6	21.9	0.433	92.4	-	-
ZoeDepth [†] [7]	11.0	81.4	1.36	83.5	-	-
UniDepth [60]	13.3	76.0	1.37	90.8	27.1	32.1
MASt3R [39]	44.8	40.1	0.717	86.7	52.5	51.4
DepthPro [9]	33.6	45.1	0.902	81.3	47.9	48.0
UniK3D-Small	89.8	13.4	0.323	93.8	82.9	92.2
UniK3D-Base	93.5	10.3	0.288	95.0	88.1	93.8
UniK3D-Large	94.6	9.3	0.275	95.6	89.5	93.7

Table 29. **Comparison on zero-shot evaluation for KITTI360.** Missing values (-) indicate the model’s inability to produce the respective output. †: ground-truth camera for 3D reconstruction. ‡: ground-truth camera for 2D depth map inference.

Method	$\delta_1 \uparrow$	A.Rel \downarrow	RMSE \downarrow	$\delta_1^{SSI} \uparrow$	F _A \uparrow	$\rho_A \uparrow$
DepthAnything [81]	-	-	-	9.5	-	-
DepthAnythingv2 [82]	-	-	-	11.3	-	-
Metric3D ^{†‡} [85]	0.2	1366.2	34.78	39.7	-	-
Metric3Dv2 ^{†‡} [28]	0.1	1655.3	40.32	43.9	-	-
ZoeDepth [†] [7]	0.7	1200.2	24.71	41.2	-	-
UniDepth [60]	29.4	152.2	4.23	44.0	14.6	7.1
MASt3R [39]	16.5	312.8	7.17	41.7	15.7	7.4
DepthPro [9]	5.5	103.8	7.35	38.0	5.9	17.5
UniK3D-Small	74.9	39.8	2.58	81.6	59.5	82.8
UniK3D-Base	73.3	33.8	2.62	80.8	61.2	80.9
UniK3D-Large	81.7	24.4	2.40	85.3	66.4	82.5

Table 30. **Comparison with UniDepth.** All models use ViT-S backbone and the same training data. Test set grouping as in the main paper. Best viewed on a screen and zoomed in.

Method	Small FoV			Small FoV _{best}			Large FoV			Panoramic		
	$\delta_1^{SSI} \uparrow$	F _A \uparrow	$\rho_A \uparrow$	$\delta_1^{SSI} \uparrow$	F _A \uparrow	$\rho_A \uparrow$	$\delta_1^{SSI} \uparrow$	F _A \uparrow	$\rho_A \uparrow$	$\delta_1^{SSI} \uparrow$	F _A \uparrow	$\rho_A \uparrow$
UniDepth [60]	89.0	54.7	77.8	92.7	35.4	45.6	71.8	41.9	48.8	34.9	1.5	1.2
UniK3D	89.1	57.3	79.8	93.1	44.6	59.3	79.8	53.5	64.8	64.3	58.6	26.3

tion in monocular metric depth estimation is inherently ill-posed, making it highly dependent on the training domain and its distribution coverage. Excessive diversity can hurt performance in narrow, specialized domains like KITTI, where models trained on large, diverse datasets often underperform compared to those trained on domain-specific data. Conversely, these models perform better in broader domains like NYU. Scale prediction is typically noisy and sensitive to domain shifts, but this issue can be mitigated through in-domain fine-tuning. For example, a few hundred optimization steps can largely resolve the “scale gap” when fine-tuning on KITTI.

- **The camera representation is superior to pinhole or fully non-parametric camera model, but you did not compare it to some common camera models, why so?**

We initially experimented with explicit parametric camera models but encountered significant drawbacks. Most standard camera models rely on backprojection operations

which are not differentiable and, thus cannot be used in a standard deep learning pipeline. Addressing this limitation requires either (i) using differentiable parametric models, such as EUCM [33] or DoubleSphere [69], (ii) approximating polynomial inversions with differentiable functions, or (iii) supervising only the model parameters without direct camera supervision. All these approaches suffer from the inherent instability of parametric models, where parameter variations need to be considered jointly on their actual output, namely the pencil of rays. This compounding effect, where small compounded changes lead to large output variations, often leads to unstable optimization. Furthermore, parametric models limit the expressiveness of the backprojection operation and constrain applicability to only those cameras the model can represent. In contrast, our representation avoids these limitations and provides greater flexibility and stability.

- **DUST3R / MAST3R architecture directly predicts point maps, are they unable to work with generic cameras?** While DUST3R and MAST3R networks can theoretically represent any camera model, our studies revealed that fully non-parametric approaches struggle when trained on diverse datasets and tested on edge cases or distribution tails. Additionally, the test-time point cloud global alignment technique used in DUST3R [74] and MAST3R [39] explicitly requires a pinhole camera, further limiting their applicability to generic cameras.
- **What is the role of the confidence prediction?** Confidence prediction is included primarily for its utility in downstream tasks and also for legacy reasons. It is worth noting that, like most regression tasks, confidence prediction is vulnerable to domain gaps, which can render it unreliable in strong out-of-domain scenarios.
- **What is the rationale of camera augmentations?** Camera augmentations were employed to address the lack of diverse real-camera data. While our simple augmentation pipeline resulted in minor improvements, we observed that many generated cameras are unrealistic and fall outside the distribution of real-world cameras. However, softmax-based warping proved effective in creating realistic images. We hypothesize that a more sophisticated camera sampling procedure, considering the realism of the output rays instead of the singled-out parameters, could significantly enhance the robustness and generalization across real and practical camera models.
- **What are the differences with UniDepth?** UniDepth [60] and UniK3D differ in camera modeling and 3D representation, both ablated in Tabs. 3, 4, and 5. UniDepth relies on the *pinhole model* by predicting the calibration matrix (cf. [60, Sec. 3.2]), thus not being able to predict *any* camera. In addition, [60] represents the 3rd dimension as *depth* (z) [60, Sec. 3.1]. These two aspects force [60] to model *only* pinhole and to output FoV

$< 180^\circ$. In contrast, UniK3D uses spherical harmonics (SH) to approximate *any* camera model and it exploits *radial distance* (r) as 3rd dimension. UniDepth projects the predicted *pinhole ray map* [60, Sec. 3.1] onto a high-dimensional space \mathbf{E} using SH, whereas UniK3D directly *predicts the SH coefficient* used to generate the ray map \mathbf{C} via inverse transform (L230-262). This key methodological difference leads to modeling any camera. Table 30 (row 1 vs. row 3) shows its impact, as UniK3D consistently outperforms [60] also when trained on identical data.

- **Has someone done something similar before?**

Yes, there are a few works [27, 35] which tried to remove the pinhole assumption for depth estimation. However, they are different for two important reasons: (i) those works focused on single-domain scenarios, leading to a simpler setting and (ii) the task is self-supervised depth estimation, where the camera is needed to define the warping-based photometric loss, inherently needing the camera, rather than supervised large-scale monocular 3D estimation.

- We provide here the δ_1^{SSI} scores of row 3 and 4 of Tab. 5: 92.1 and 92.2, respectively. This score similarity, along with F_A and ρ_A drops (Tab. 5), spotlights angular module’s role. In fact, radial- and SH-based model (row 4) overestimates FoV of images with lens distortions. Retraining with stronger distortion augmentation for small FoV leads to $(F_A, \rho_A) = (43.1, 62.3)$, validating our assumption.

E. Additional Qualitative Results

We provide here more qualitative comparisons, in particular from validation domains not presented in the main paper and with distorted cameras, namely ScanNet++ (DSLR), IBims-1_{Dist}, and Diode_{Dist} (Indoor), in Fig. 5. In addition, we test our model on complete in-the-wild scenarios, for instance, frames from movies, TV series, YouTube, or animes. All images depicted in Fig. 6 and Fig. 7 present deformed cameras or unusual points of view. The visualization here presented, both from the validation sets and the in-the-wild ones are casually selected and not cherry-picked.

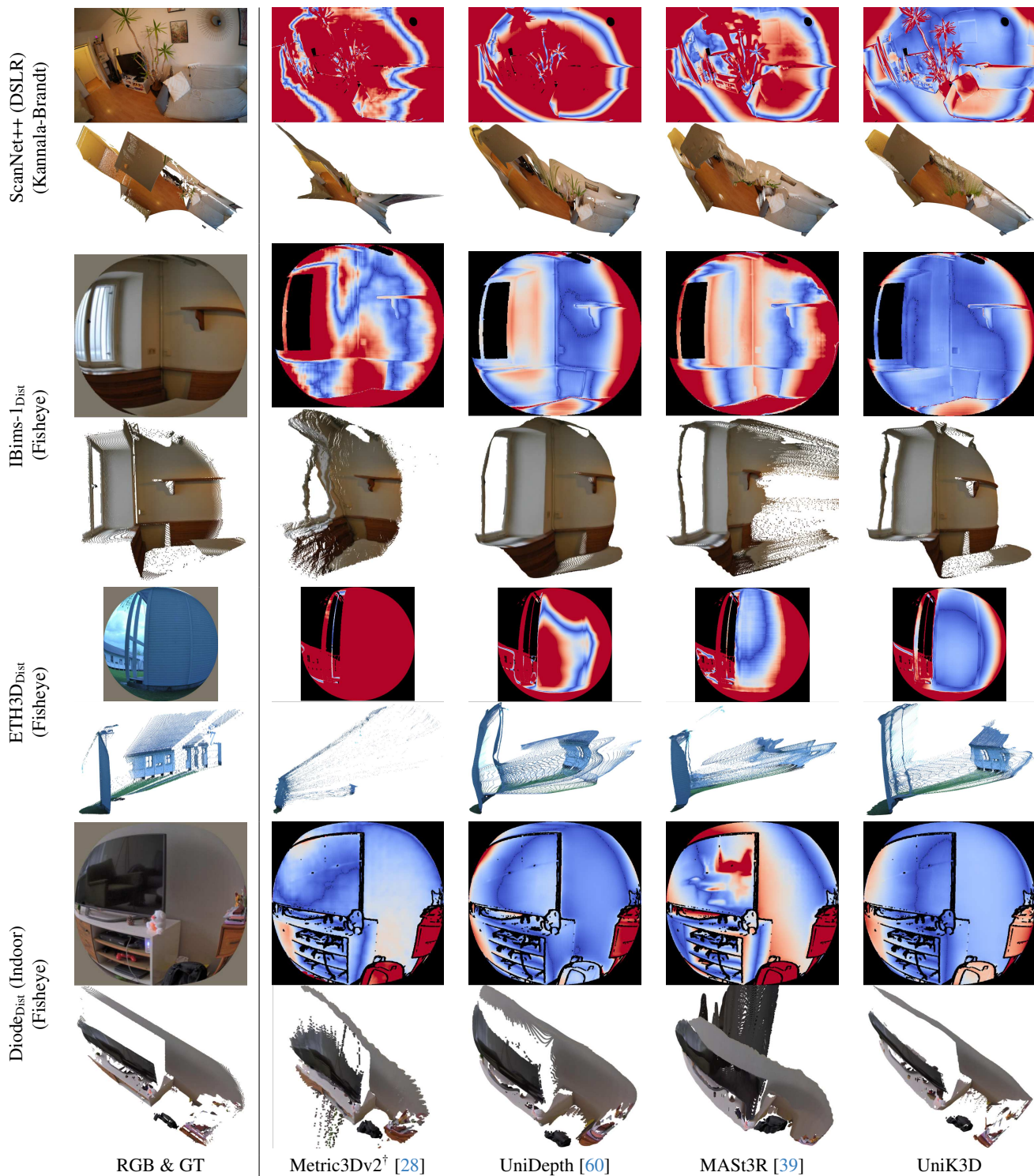


Figure 5. **Qualitative comparisons.** Each pair of consecutive rows represents one test sample. Each odd row displays the input RGB image and the 2D error map, color-coded with the *coolwarm* colormap based on absolute relative error with blue corresponding to 0% error and red to 25%. To ensure a fair comparison, errors are calculated on GT-based shifted and scaled outputs for all models. Each even row shows the ground truth and predictions of the 3D point cloud. All samples are randomly selected and not picked. †: GT-camera unprojection.

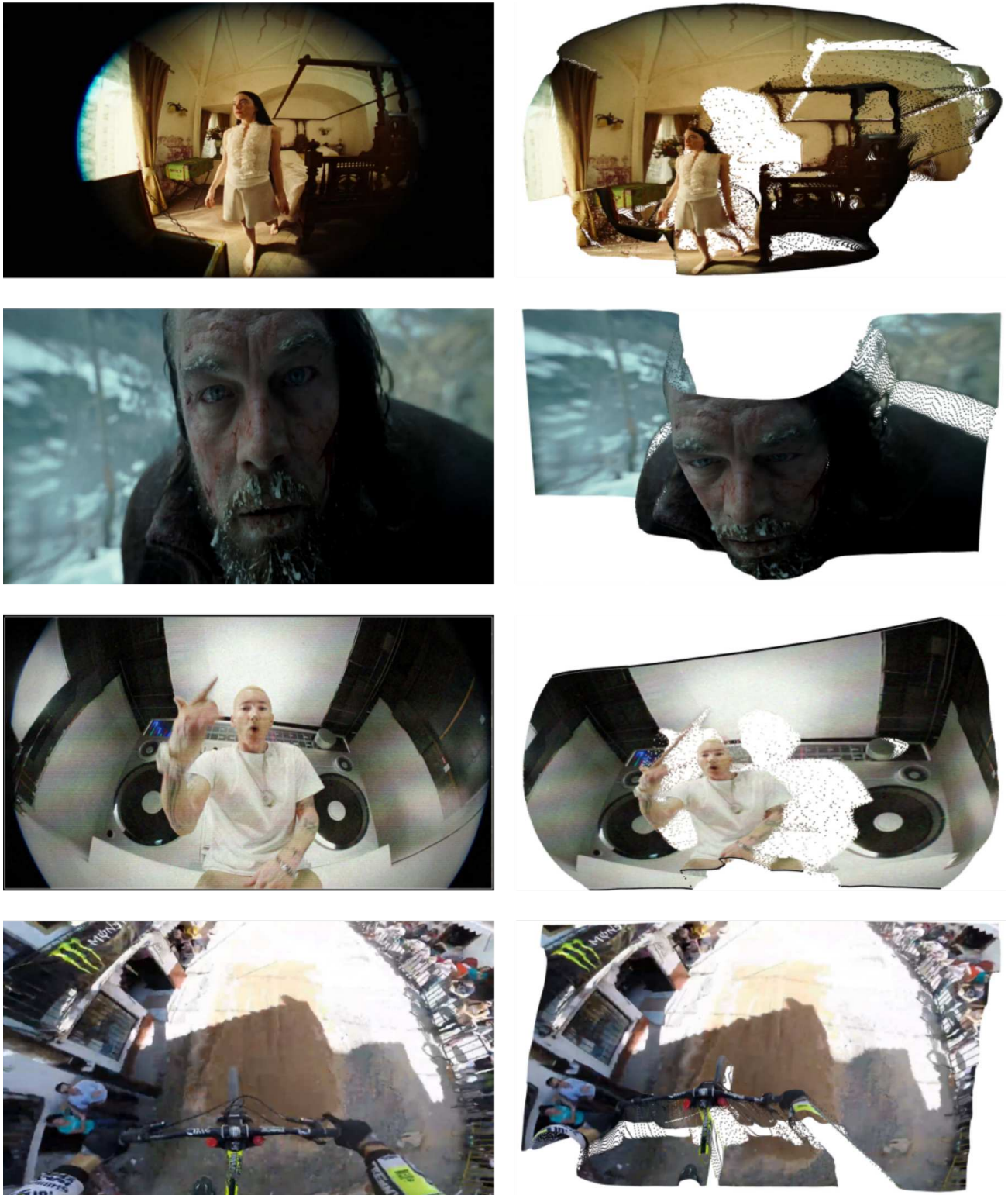


Figure 6. **Qualitative in-the-wild 3D results.** UniK3D is fed solely each single image in the left column and it outputs the corresponding point cloud in the right column, the point of view is slightly tilted to better appreciate the 3D. The images are video frames respectively from Poor Things (movie), The Revenant (movie), Eminem (music video), and YouTube (egocentric GoPro). The frames present a variety of camera types and unusual viewpoints.



Figure 7. **Qualitative in-the-wild 3D results.** UniK3D is fed solely each single image in the left column and it outputs the corresponding point cloud in the right column, the point of view is slightly tilted to better appreciate the 3D. The images are video frames respectively from Trainspotting (movie), YouTube (doorbell camera), Naruto (anime), and Breaking Bad (TV series). The frames present a variety of camera types and unusual viewpoints.