# Depth Any Panoramas: A Foundation Model for Panoramic Depth Estimation

Xin Lin[2]    Meixi Song[1]    Dizhe Zhang[1‡]    Wenxuan Lu[1]    Haodong Li[2]

Bo Du[3]    Ming-Hsuan Yang[4]    Truong Nguyen[2*]    Lu Qi[1,3*†]

[1] Insta360 Research    [2] University of California, San Diego    [3] Wuhan University    [4] University of California, Merced

Figure 1. Metric depth visualizations generated by DAP from diverse panoramic inputs. For clarity, each depth map is displayed using its own adaptive truncation range. DAP achieves robust, metrically consistent panoramic depth across diverse real-world scenes, highlighting the power of large-scale data and model designing.

## Abstract

*In this work, we present a panoramic metric depth foundation model that generalizes across diverse scene distances. We explore a data-in-the-loop paradigm from the view of both data construction and framework design. We collect a large-scale dataset by combining public datasets, high-quality synthetic data from our UE5 simulator and text-to-image models, and real panoramic images from the web. To reduce domain gaps between indoor/outdoor and synthetic/real data, we introduce a three-stage pseudo-label curation pipeline to generate reliable ground truth for unlabeled images. For the model, we adopt*

---

\* Equal advising    † Corresponding author    ‡ Project leader

*DINOv3-Large as the backbone for its strong pre-trained generalization, and introduce a plug-and-play range mask head, sharpness-centric optimization, and geometry-centric optimization to improve robustness to varying distances and enforce geometric consistency across views. Experiments on multiple benchmarks (e.g., Stanford2D3D, Matterport3D, and Deep360) demonstrate strong performance and zero-shot generalization, with particularly robust and stable metric predictions in diverse real-world scenes. The project page can be found at: https://insta360-research-team.github.io/DAP_website/*

## 1. Introduction

The estimation of panoramic depth has attracted increasing attention due to the full $360° \times 180°$ coverage of the surrounding environment for spatial intelligence. It has benefited various robotic applications, such as omnidirectional obstacle avoidance during navigation tasks.

Despite its importance, the panoramic depth estimation still lags behind. Both panorama-specific relative/scale-invariant methods (e.g., Panda [9], Depth Anywhere [44], $DA^2$ [25]) and unified metric-depth (DAC [15], Unik3D [32]) frameworks struggle to generalize to diverse real-world scenes, particularly outdoors. A possible reason is the limited scale and diversity of existing data, due to the high cost of data collection and annotation.

Motivated by this, we explore the data-in-the-loop paradigm for panoramic depth estimation, which raises two key challenges for data scaling: constructing large-scale datasets with reliable, high-quality ground truth, and designing models that can effectively adapt to such data scaling. Addressing both challenges is essential for a geometry-consistent and generalizable panoramic foundation model.

At the data level, we integrate the indoor dataset Structured3D [52], synthesize 90K high-quality outdoor samples using UE5-based AirSim360 simulator [14], and collect about 1.7M unlabeled panoramic images from internet and generate 200k indoor panoramas by DiT360 [13]. To reduce performance degradation caused by domain gaps between indoor/outdoor and synthetic/real data, we propose a two-stage pseudo-label curation pipeline. In the first stage, we train a Scene-Invariant Labeler on a balanced mix of Structured3D and AirSim360, then generate pseudo labels for the 1.9M unlabeled images. In the second stage, we select the top 60k highest-confidence pseudo-labeled samples in both indoor and outdoor scenes using a discriminator, and train the Realism-Invariant Labeler on this expanded dataset to refine the pseudo labels. Finally, in the third stage, our foundation model DAP is trained jointly on all labeled data and the refined pseudo-labeled data produced by the Realism-Invariant Labeler.

For the foundation model design, we adopt metric depth estimation for its scalability to panoramas across arbitrary distances. We use DINOv3-Large as the encoder to leverage a strong pre-trained visual priors and conventional depth head module for dense estimation [49]. We further introduce a set of losses to enforce robustness to distances and geometric consistency across viewpoints. Specifically, a lightweight depth-range mask head is presented to mitigate uneven depth distributions by filtering regions using distance thresholds (e.g., 10/20/50/100 m) in a plug-and-play manner. Moreover, Silog, DF-Gram, Gradient, Normal, and Point-Cloud losses are designed by a distortion-aware map to compensate for non-uniform pixel geometry in equirect-angular projection.

The extensive experiments on both indoor and outdoor test sets (e.g., Stanford2D3D, Matterport3D, and Deep360) demonstrate the strong generalization and effectiveness of our foundation model across diverse benchmarks. Beyond quantitative superiority, as shown in Figure 1, our model exhibits excellent visual consistency and scale-awareness, producing realistic depth predictions in challenging real-world scenarios with complex geometry, distant regions, and sky areas. These results highlight the model's robustness and strong adaptability to both synthetic and real-world environments. In summary, the main contributions can be summarized as follows:

- We construct a large-scale panoramic dataset with more than 2M samples across synthetic and real domains. It includes high-quality indoor data from Structured3D, 90K photorealistic outdoor panoramas rendered with the UE5-based AirSim360 simulator, and 1.9M filtered panoramas collected from internet and DiT360, enabling diverse and scalable depth supervision.

- We design a three-stage pipeline that progressively refines pseudo labels and bridges both the synthetic–real and indoor–outdoor domain gaps. This is achieved through multiple curation techniques and large-scale semi-supervised learning to enhance cross-domain generalization.

- We incorporate a plug-and-play range mask together with geometry- and sharpness-oriented optimization, where the associated loss terms ($\mathcal{L}_{SILog}$, $\mathcal{L}_{DF}$, $\mathcal{L}_{grad}$, $\mathcal{L}_{normal}$, and $\mathcal{L}_{pts}$), ensuring metric consistency and structural fidelity across diverse panoramas.

- Comprehensive evaluations show that our model generalizes well to open real-world scenarios, producing scale-consistent and perceptually coherent depth maps. It achieves state-of-the-art performance both quantitatively on synthetic benchmarks and qualitatively on diverse real-world panoramas.

## 2. Related work

**Perspective Depth Estimation.** With the rapid progress of deep learning and large-scale perspective depth datasets,

Table 1. Comparison of training data compositions used by recent panoramic depth estimation methods. Unlike previous approaches, which rely on limited or domain-specific datasets, our DAP data engine scales up to 2M panoramas across both indoor/outdoor and synthetic/real domains, providing a unified and comprehensive data foundation for panoramic depth modeling. * in DA$^2$ refers to pseudo-panoramic data generated from perspective images through P2E projection and out-painting model.

| Method | Metric Ability | Dataset Type | Number | Scene | | Panorama | |
| | | | | Indoor | Outdoor | Synthetic | Real World |
|---|---|---|---|---|---|---|---|
| PanDA [9] | ✗ | Panorama | 122k | ✓ | ✓ | ✓(20k) | ✓(102k) |
| DA$^2$ [25] | ✗ | Mixed* | 606k | ✓ | ✓ | ✓(63k) | ✗ |
| Unik3D [32] | ✓ | Mixed | 694k | ✓ | ✓ | ✓(1k) | ✗ |
| DAC [15] | ✓ | Mixed | 800k | ✓(670k) | ✓(130k) | ✗ | ✗ |
| DAP (Ours) | ✓ | Panorama | 2M | ✓(500k) | ✓(1.5M) | ✓(300k) | ✓(1.7M) |

perspective depth estimation has advanced rapidly, with recent metric and scale-invariant models achieving strong performance, such as UniDepth [31, 33], Metric3D [18, 50], DepthPro [8], and MoGe [46, 47]. Some relative depth estimation methods have greatly benefited from data scaling, with models like DepthAnything [48, 49] showing impressive zero-shot generalization. Recently, some methods fine-tune large pre-trained generative models with strong prior capabilities, such as Stable Diffusion [17, 35] and FLUX [7], on limited but high-quality datasets, achieving competitive results [16, 21, 24, 42]. Nevertheless, the perspective paradigm inherently restricts perception to a limited field of view (FoV), failing to capture the complete 360° spatial geometry of a scene.

**Panoramic Depth Estimation.** *In-domain.* Early methods focus on in-domain settings, where models are trained and evaluated on the same dataset. To address the severe distortions of the equirectangular projection (ERP), there are two main directions: distortion-aware designs [23, 29, 36, 39, 51, 53] and projection-driven strategies [1–3, 5, 6, 11, 20, 27, 30, 34, 37, 40]. However, over-reliance on in-domain training often leads to overfitting and limited generalization, motivating recent efforts toward zero-shot and cross-domain panoramic depth estimation.

*Zero-shot.* Compared with in-domain training, zero-shot panoramic depth estimation is more practical for cross-domain applications due to its stronger generalization ability. Existing approaches can be grouped into three categories: The first leverages pretrained perspective depth models to generate pseudo labels for panoramas, as in Depth Anywhere [43] and PanDA [9], which distill knowledge via cube projection or semi-supervised learning on large-scale unlabeled panoramas. More recently, DA$^2$ [25] improves zero-shot performance by expanding training data through perspective-to-ERP conversion and diffusion-based out-painting, combined with a distortion-aware transformer. A third line of work aims at universal, metric-capable

camera modeling: Depth Any Camera [15] unifies diverse imaging geometries under an ERP representation with geometry-driven augmentation, while UniK3D [32] reformulates depth in spherical coordinates using harmonic ray representations to enhance wide-FoV generalization.

## 3. Method

We introduce a scalable panoramic depth foundation model that unifies metric estimation across diverse domains and scene types. Our approach consists of three main components: a large-scale data engine for data scaling across synthetic and real panoramas (Sec. 3.1), a three-stage pipeline for effectively exploiting large-scale unlabeled data (Sec. 3.2), and a geometry-consistent network design with multi-loss optimization for high-quality metric depth estimation (Sec. 3.3).

### 3.1. Data Engine

**Overview.** To enable the scaling of the panoramic depth estimation, we construct a comprehensive data engine that unifies diverse sources in synthetic and real domains, as summarized in Table 1. Compared to the datasets used in previous panoramic depth methods such as PanDA [9], UniK3D [32], and DAC [15], our data engine achieves the largest panorama scale with about 2 million samples and the broadest domain coverage, spanning indoor/outdoor and synthetic/real-world panoramas. It serves as the cornerstone for building our foundation model by supporting both large-scale supervised and semi-supervised training.

**Simulated Outdoor Scene Data.** To address the scarcity of outdoor panoramic supervision, we construct a synthetic outdoor dataset named DAP-2M-Labeled using the high-fidelity simulation platform Airsim360 [14]. We simulate drone flights following mid- and low-altitude trajectories across diverse environments to capture panoramic imagery and corresponding depth maps under realistic illumination and environmental conditions. In total, over 90K panoramic
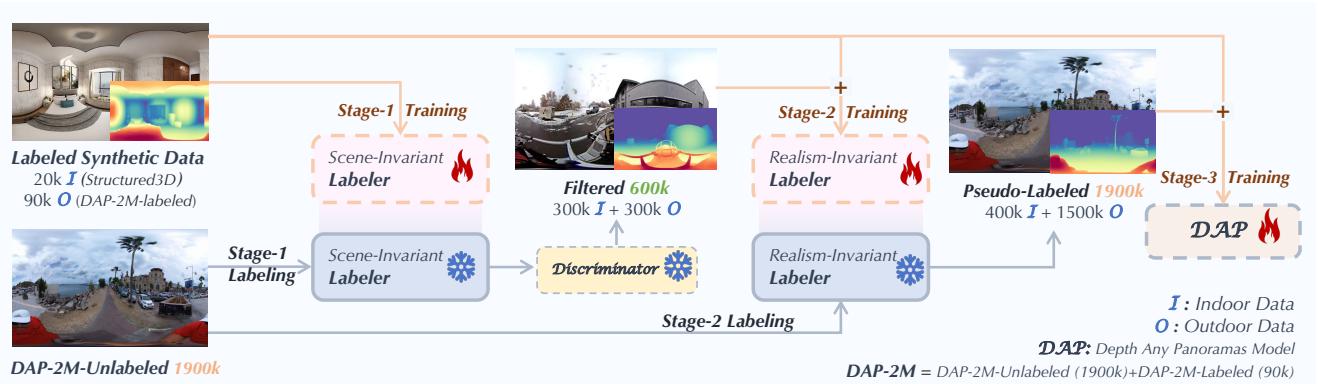
Figure 2. Overview of the proposed progressive three-stage pipeline. Stage 1 trains a **Scene-Invariant Labeler** on high-quality synthetic indoor and outdoor data to provide strong initialization. Stage 2 introduces a **Realism-Invariant Labeler**, where a PatchGAN-based discriminator selects 300K indoor and 300K outdoor high-confidence pseudo-labeled samples to mitigate domain gaps between synthetic and real data. Stage 3 performs **DAP** training on all labeled and pseudo-labeled data, enabling large-scale semi-supervised learning and strong generalization across real-world panoramic scenes.

frames with pixel-aligned depth annotations are collected from five representative outdoor scenes, including *New York City*, *SF City*, *Downtown West*, *City Park*, and *Rome*, covering more than 26,600 complete panoramic sequences.

**Unlabeled Data,** We have collected 250K panoramic videos from the internet and processed them into image frames. After a careful curation procedure and filtering out samples with unreasonable horizons, we obtain a total of 1.7 million high-quality panoramic images. We then employ the large multimodal model Qwen2-VL [45] to automatically categorize these panoramas into indoor and outdoor scenes, yielding approximately 250K indoor and 1.45M outdoor samples. Since indoor panoramas are relatively scarce in real-world data, we further supplement this domain by generating an additional 200K indoor samples using the state-of-the-art panoramic generation model DiT-360 [13]. Together, these collections form our DAP-2M-Unlabeled, which provides abundant coverage of diverse environments for pretraining and semi-supervised learning.

### 3.2. Three-Stage Pipeline

As shown in Figure 2, we adopt a three-stage pipeline to efficiently exploit unlabeled large-scale data in the real world, while enhancing the learning capacity of the network through geometric and dense-detail supervision paradigms. Our pipeline consists of three stages:

**Stage 1. Scene-Invariant Labeler Training.** We first train a Scene-Invariant Labeler on 20k synthetic indoor and 90k synthetic outdoor datasets with accurate metric depth annotations. The goal of this stage is to learn a labeling model that generalizes across both indoor and outdoor environments rather than overfitting to specific scene layouts or lighting conditions. Training on geometrically and photometrically diverse synthetic scenes enables the labeler to

learn consistent, physically grounded depth cues that generalize across domains, providing a robust initialization for generating reliable pseudo-depth labels on real-world panoramic data.

**Stage 2. Realism-Invariant Labeler Training.** We first pre-train a depth quality discriminator to assess the reliability of depth predictions, where synthetic ground-truth depth maps are treated as real samples and the Scene-Invariant Labeler outputs as fake ones, this step enables the discriminator to learn a scene-agnostic quality prior for subsequent filtering. For more details, please refer to the supplementary material. Next, we apply the Scene-Invariant Labeler to all unlabeled real images and estimate depth quality using the trained discriminator. The top-ranked 300K indoor and 300K outdoor samples are selected as high-confidence pseudo-labeled data, which are combined with the synthetic datasets from Stage 1 to train a Realism-Invariant Labeler. By learning from reliable pseudo labels across diverse real domains, this labeler becomes robust to appearance variations and realism-specific differences, enabling it to generalize beyond synthetic textures and lighting conditions.

**Stage 3. DAP Training.** As shown in Table 2, our Depth Any Panorama (DAP) model is trained on all 1.9M pseudo-labeled data generated by the Realism-Invariant Labeler and previous labeled data. This enables the DAP model to effectively benefit from dense features and large-scale pseudo supervision, leading to improved generalization on real-world panoramic depth estimation.

### 3.3. Model Design

As illustrated in Figure 3, our DAP network takes a panoramic image as input and uses DINOv3-Large [38] as a visual backbone for powerful feature extraction. We introduce two task-specific heads: a range mask head and a
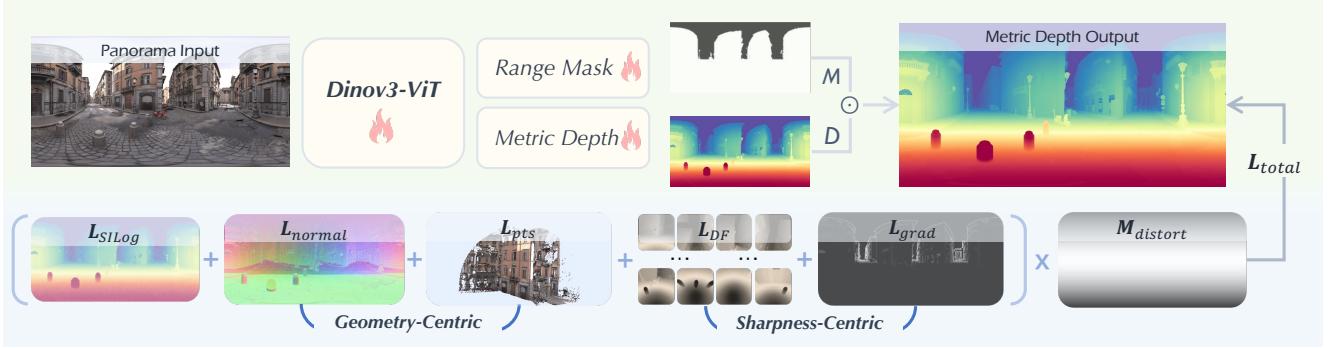
Figure 3. Architecture of the proposed **DAP network**. Built upon DINOv3-Large [38] as the visual backbone, our model adopts a distortion-aware depth decoder and a plug-and-play range mask head for adaptive distance control across diverse scenes. Training is guided by multi-level geometric and sharpness-aware losses, including $\mathcal{L}_{SILog}$, $\mathcal{L}_{DF}$, $\mathcal{L}_{grad}$, $\mathcal{L}_{normal}$, and $\mathcal{L}_{pts}$ losses, ensuring metric accuracy, edge fidelity, and geometric consistency in panoramic depth estimation.

Table 2. Overview of datasets used for training DAP, covering synthetic and real, labeled and unlabeled panoramic data.

| Datasets | Indoor | Outdoor | Label | Samples |
|---|---|---|---|---|
| Structured3D [52] | ✓ | | ✓ | 18,298 |
| DAP-2M-Labeled | | ✓ | ✓ | 90k |
| DAP-2M-Unlabeled | ✓ | ✓ | ✗ | 1.9M |

metric depth head. The range mask head outputs a binary mask $M$ that defines valid spatial regions under different distance thresholds, while the depth head predicts a dense metric depth map $D$. To accommodate diverse environments from confined indoor spaces to large-scale outdoor scenes, we provide four plug-and-play range mask heads with distance thresholds of 10 m, 20 m, 50 m, and 100 m, allowing flexible adaptation to different spatial scales. Each range mask head is independently optimized with a combination of weighted BCE and Dice losses:

$$\mathcal{L}_{mask} = \|M - M_{gt}\|^2 + 0.5\,\mathcal{L}_{Dice}(M, M_{gt}), \quad (1)$$

where $M$ and $M_{gt}$ denote the predicted and ground-truth range masks, respectively. The final metric depth output is obtained through element-wise multiplication of $M$ and $D$, ensuring that the predictions remain physically valid and scale-consistent in varying depth ranges. This dual-head design enables DAP to robustly adapt to diverse scene geometries and maintain high-quality metric depth estimation across a wide spectrum of spatial conditions. For optimization, besides adopting the SILog loss $\mathcal{L}_{SILog}$ as in previous works [9], we introduce a set of complementary loss functions to enhance geometric consistency and edge fidelity.

**Sharpness-centric Optimization** Inspired by [38, 46], we introduce a dense fidelity constraint, termed $\mathcal{L}_{DF}$, to enhance local sharpness and structural consistency. To mitigate geometric distortion in ERP, we first decompose each

depth map into 12 perspective patches using virtual cameras positioned at the vertices of an icosahedron. These perspective views preserve fine-grained details and avoid the stretching artifacts near poles, thus providing higher-fidelity supervision. For each view, we apply a valid mask, normalize depth values, and compute a Gram-based similarity between the predicted and ground-truth depth maps. The final $\mathcal{L}_{DF}$ is defined as the average loss across all $N = 12$ views:

$$\mathcal{L}_{DF} = \frac{1}{N} \sum_{k=1}^{N} \left\| D_{pred}^{(k)} \odot D_{pred}^{(k)}{}^{\top} - D_{gt}^{(k)} \odot D_{gt}^{(k)}{}^{\top} \right\|_F^2,$$
(2)

where $N = 12$ denotes the number of views.

**Sharpness-oriented Gradient Refinement.** To further enhance the sharpness of object boundaries, we introduce a gradient-based loss $\mathcal{L}_{grad}$ that explicitly focuses on high-frequency edge regions in the ERP domain. While the $\mathcal{L}_{DF}$ strengthens dense fidelity on distortion-free perspective patches, $\mathcal{L}_{grad}$ complements it by preserving local discontinuities directly in the ERP representation. Specifically, we compute gradient magnitude maps using Sobel operators along both $x$ and $y$ directions and derive an edge mask $M_E$ by thresholding the ground-truth Sobel gradient magnitudes. The gradient loss is then applied only within these masked regions using SILog loss [12]:

$$\mathcal{L}_{grad} = \mathcal{L}_{SILog}(M_E \odot D_{pred}, M_E \odot D_{gt}). \quad (3)$$

This design improves the consistency and sharpness of depth boundaries, effectively complementing the dense fidelity constraint for recovering fine geometric details.

**Geometry-centric optimization** To improve geometric consistency, we incorporate a normal loss $\mathcal{L}_{normal}$ [23, 25]. Both predicted and ground-truth depth maps are converted into surface normal fields $\mathbf{n}_{pred}, \mathbf{n}_{gt} \in \mathbb{R}^{H \times W \times 3}$. The $\mathcal{L}_{normal}$ is then defined as the L1 distance between pre-
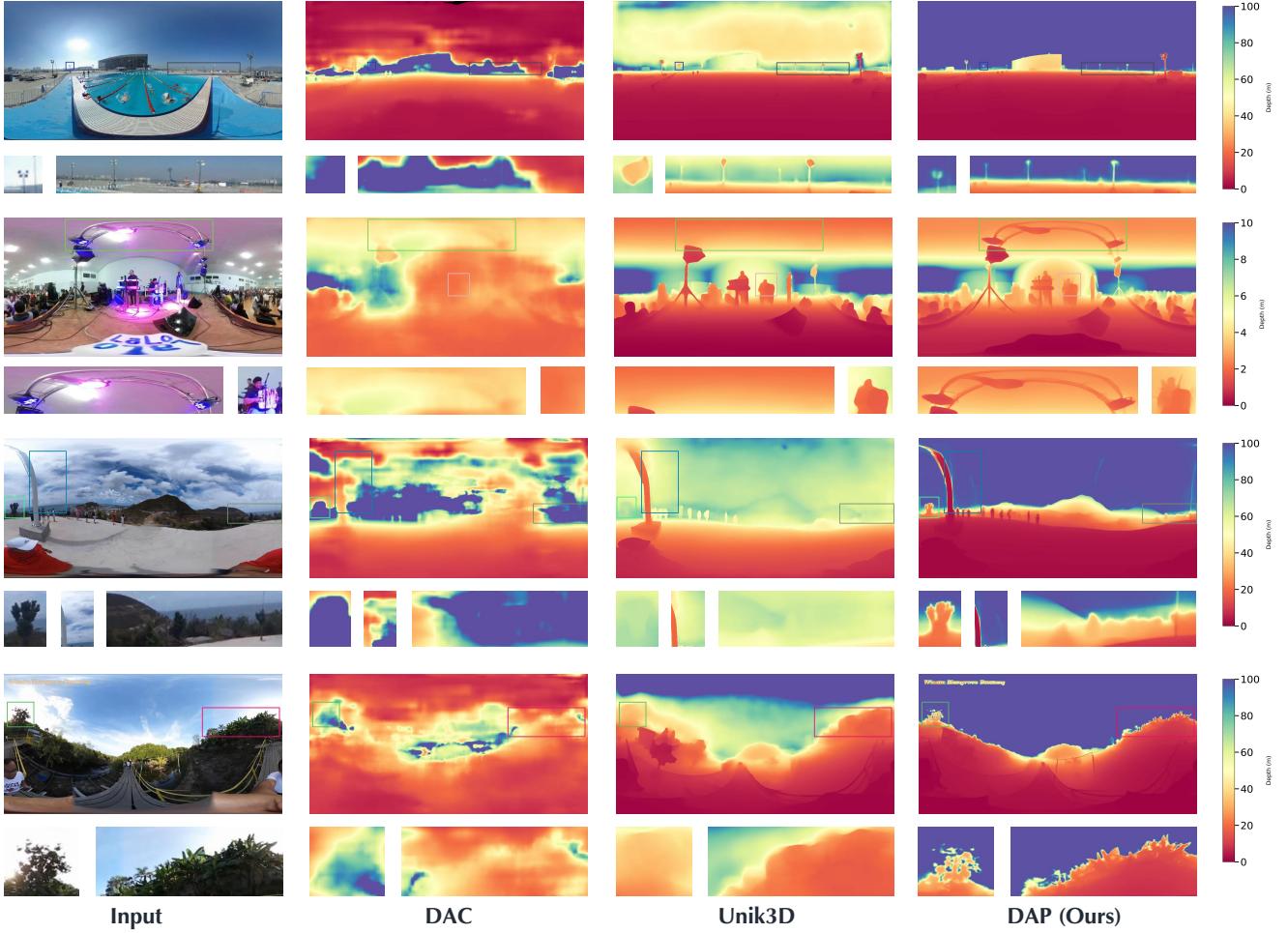
Figure 4. Qualitative comparison across diverse real-world indoor and outdoor scenes. Our DAP produces sharper object boundaries, smoother global geometry, and superior robustness in distant and sky regions compared to DAC [15] and Unik3D [32].

dicted and ground-truth normals:

$$\mathcal{L}_{normal} = \|\mathbf{n}_{pred}(i,j) - \mathbf{n}_{gt}(i,j)\|_1. \quad (4)$$

We further use a point cloud loss $\mathcal{L}_{pts}$. The depth maps are projected onto the spherical coordinate system to obtain 3D point clouds $\mathbf{P}_{pred}, \mathbf{P}_{gt} \in \mathbb{R}^{H \times W \times 3}$, and the loss is defined as:

$$\mathcal{L}_{pts} = \|\mathbf{P}_{pred}(i,j) - \mathbf{P}_{gt}(i,j)\|_1. \quad (5)$$

**Overall Objective** The overall training objective is a weighted combination of all the above losses with distortion map to address the non-uniform pixel distribution [28] in equirectangular projections. The distortion map compensates for the over-representation of pixels near the poles, ensuring balanced gradient contributions across the spherical domain:

$$\mathcal{L}_{total} = M_{distort} \odot (\lambda_1 \mathcal{L}_{SILog} + \lambda_2 \mathcal{L}_{DF} + \lambda_3 \mathcal{L}_{grad}$$
$$+ \lambda_4 \mathcal{L}_{normal} + \lambda_5 \mathcal{L}_{pts} + \lambda_6 \mathcal{L}_{mask}), \quad (6)$$

where $\lambda_i, \ i \in \{1, 2, 3, 4, 5, 6\}$, denotes a weight parameters.

## 4. Experiment

**Training Datasets.** As discussed in Sec. 3.1, the datasets used for training are summarized in Table 2. We utilize both labeled and unlabeled data from diverse indoor and outdoor environments to enhance the generalization ability of our model. Structured3D [52] provides high-quality synthetic indoor scenes with ground-truth depth supervision. DAP-2M-Labeled contains 90K real-world samples with depth annotations, while DAP-2M-Unlabeled includes 1.9M real images without depth labels, which are leveraged for pseudo-label learning in our pipeline. For our DAP model, the training resolution is $512 \times 1024$.

**Evaluation Datasets & Metrics.** For evaluation, following prior work [9], we assess our method on two widely used indoor datasets, Matterport3D [10] and Stanford2D3D [4],
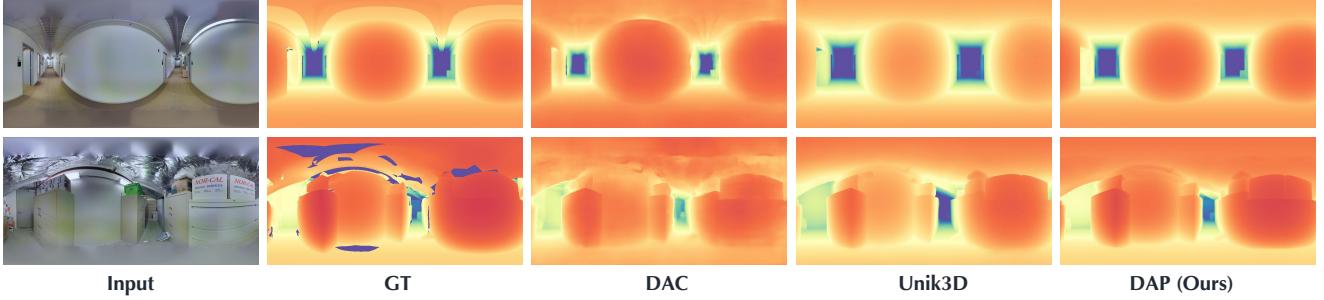
Figure 5. Qualitative comparison on Stanford2D3D. Our method preserves fine structural details and demonstrates superior scale-awareness.

Table 3. Zero-shot comparison of panoramic metric depth estimation on three benchmarks. The best and second best metric depth performances are highlighted. Our DAP consistently achieves the best results across all datasets, demonstrating strong generalization without fine-tuning. We also include several scale-invariant methods only for reference from $DA^2$ [25].

| Methods | | Stanford2D3D (Indoor) | | | Matterport3D (Indoor) | | | Deep360 (Outdoor) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | AbsRel($\downarrow$) | RMSE($\downarrow$) | $\delta_1(\uparrow$) | AbsRel($\downarrow$) | RMSE($\downarrow$) | $\delta_1(\uparrow$) | AbsRel($\downarrow$) | RMSE($\downarrow$) | $\delta_1(\uparrow$) |
| Scale-invariant | MoGe [46] | 0.1581 | 0.2576 | 0.7902 | 0.1004 | 0.3591 | 0.9080 | - | - | - |
| | VGGT [41] | 0.1870 | 0.3350 | 0.7408 | 0.1078 | 0.3880 | 0.8870 | - | - | - |
| | DepthAnythingv2 [49] | 0.3679 | 0.4339 | 0.4766 | 0.2585 | 0.7067 | 0.5842 | - | - | - |
| | PanDA [9] | 0.1648 | 0.2364 | 0.7326 | 0.0888 | 0.3325 | 0.9209 | - | - | - |
| | $DA^2$ [25] | 0.0723 | 0.1400 | 0.9545 | 0.0667 | 0.2882 | 0.9561 | - | - | - |
| Metric | Unik3D [32] | 0.1795 | 0.4850 | 0.7823 | 0.2224 | 0.6680 | 0.6634 | 0.0885 | 6.148 | 0.9293 |
| | DAC [15] | 0.1366 | 0.4509 | 0.8393 | 0.1803 | 0.9390 | 0.7203 | 0.2611 | 8.371 | 0.6311 |
| | DAP (Ours) | 0.0921 | 0.3820 | 0.9135 | 0.1186 | 0.7510 | 0.8518 | 0.0659 | 5.224 | 0.9525 |

to evaluate its zero-shot performance. For outdoor scenes, we use Deep360 [26] as the test set, also under a zero-shot setting. Since existing outdoor data are limited, we further introduce a new benchmark, DAP-Test, which consists of 1,343 high-quality outdoor images with accurate depth annotations. Following [19], we evaluate depth estimation performance with metrics including Absolute Relative Error (*AbsRel*), Root Mean Squared Error (*RMSE*), and a percentage metrics $\delta_1$, where $i = 1.25$.

**Implementation Details.** All experiments are conducted on H20 GPUs. For model training, the learning rate is set to 5e-6 for the ViT backbone and 5e-5 for the decoders, using the Adam optimizer [22]. The loss weight $\lambda_1$, $\lambda_2$, $\lambda_3$, $\lambda_4$, $\lambda_5$ and $\lambda_6$ are set to 1.0, 0.4, 5.0, 2.0, 2.0 and 2.0, respectively. Data augmentation includes color jittering, horizontal translation, and flipping following [9]. Please refer to supplementary material for more details.

## 4.1. Qualitative and Quantitative Evaluation

**Zero-shot Performance.** Table 3 reports the zero-shot metric depth results on Stanford2D3D, Matterport3D, and Deep360. We compare our method with recent metric approaches and also list several scale-invariant methods only as reference. Note that scale-invariant models require

Table 4. Quantitative comparison on the proposed DAP-Test benchmark. Our DAP achieves the best performance across all metrics, demonstrating the effectiveness of large-scale data scaling and domain-consistent training.

| Method | AbsRel($\downarrow$) | RMSE($\downarrow$) | $\delta_1(\uparrow$) |
|---|---|---|---|
| DAC [15] | 0.3197 | 8.799 | 0.5193 |
| Unik3D [32] | 0.2517 | 10.56 | 0.6086 |
| DAP (Ours) | 0.0781 | 6.804 | 0.9370 |

ground-truth depth to obtain scale during evaluation, while our method predicts absolute metric scale directly from the input panorama without any post-alignment. Nevertheless, our results remain comparably good with them.

Across all three benchmarks, DAP consistently delivers the best performance. On Stanford2D3D and Matterport3D, it significantly lowers AbsRel while achieving markedly higher $\delta_1$, demonstrating strong generalization to unseen indoor scenes. On Deep360, DAP obtains the lowest AbsRel and RMSE and the highest $\delta_1$. These results show that DAP unifies indoor and outdoor panoramic depth estimation within a single foundation model and achieves state-of-the-art performance without any fine-tuning, demonstrating

Table 5. Ablation study on the proposed components. The **best** and **second best** performances are highlighted.

| Distortion Map | Geometry Loss | Sharpness Loss | Stanford2D3D (Indoor) | | | Deep360 (Outdoor) | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | | AbsRel ($\downarrow$) | RMSE ($\downarrow$) | $\delta_1$ ($\uparrow$) | AbsRel ($\downarrow$) | RMSE ($\downarrow$) | $\delta_1$ ($\uparrow$) |
| ✗ | ✗ | ✗ | 0.1166 | 0.449 | 0.8409 | 0.0942 | 6.374 | 0.8396 |
| ✓ | ✗ | ✗ | 0.1149 | 0.449 | 0.8440 | 0.0926 | 6.297 | 0.8423 |
| ✓ | ✓ | ✗ | 0.1112 | 0.448 | 0.8509 | 0.0880 | 6.373 | 0.8592 |
| ✓ | ✓ | ✓ | 0.1084 | 0.442 | 0.8576 | 0.0862 | 6.212 | 0.8719 |

Table 6. Ablation on the range mask head (m). The **best** and **second best** results are highlighted.

| Mask Depth | DAP-2M-Labeled | | Deep360 | |
|:---:|:---:|:---:|:---:|:---:|
| | AbsRel ($\downarrow$) | $\delta_1$ ($\uparrow$) | AbsRel ($\downarrow$) | $\delta_1$ ($\uparrow$) |
| 10 | 0.0801 | 0.9315 | 0.0934 | 0.8493 |
| 20 | 0.0823 | 0.9164 | 0.0873 | 0.8668 |
| 50 | 0.0864 | 0.9104 | 0.0843 | 0.8594 |
| 100 | 0.0793 | 0.9353 | 0.0862 | 0.8719 |
| ✗ | 0.0832 | 0.9042 | 0.0938 | 0.8411 |

strong robustness and metric consistency even under such challenging conditions, achieving smooth and realistic representations of both near and far spatial structures. These visual results further confirm that our large-scale data training and range-aware design effectively improve both metric fidelity and visual realism in panoramic depth estimation.

Figure 5 presents a qualitative comparison on the Stanford2D3D dataset. Compared with existing methods, our approach more accurately reconstructs distant structures and maintains fine geometric details that other models tend to blur or oversmooth. The results also show that our method exhibits stronger scale-awareness, producing depth and illumination distributions that are visually closer to the ground truth. In particular, the distant wall and ceiling regions in our predictions retain consistent color gradients and structural integrity, demonstrating the effectiveness of our degradation-perceptive design in preserving global scale consistency while recovering high-frequency details.

### 4.2. Ablation Studies

To fairly evaluate the contribution of each component in our framework while minimizing the influence of external priors, we adopt DINOv3-Large as the encoder and fully fine-tune all variants. We conduct comprehensive ablation studies to analyze the effectiveness of key modules in DAP. Additional visual comparisons are provided in the supplementary material.

**Model Design.** Table 5 summarizes the ablation results. Starting from a baseline trained only with the standard $\mathcal{L}_{SILog}$ loss, we progressively add the distortion map, geometry-consistent losses, and sharpness-related losses. The distortion map improves optimization stability under ERP distortion, while the geometry-centric losses ($\mathcal{L}_{normal}$ and $\mathcal{L}_{pts}$) further enhance structural consistency. Adding the sharpness-centric losses ($\mathcal{L}_{DF}$ and $\mathcal{L}_{grad}$) achieves the best performance, yielding the lowest AbsRel (0.1084/0.0862) and highest $\delta_1$ (0.8576/0.8719) on Stanford2D3D and Deep360.

**Range Mask Head.** We evaluate the plug-and-play range mask head using thresholds of 10 m, 20 m, 50 m, and 100 m (Table 6), with ground-truth depth truncated accordingly for fair comparison. Smaller thresholds (10 m and 20 m) emphasize near-range geometry and achieve $\delta_1$ above 0.91

both the scalability of our data engine and the robustness of the three-stage pipeline.

**DAP-Test Performance.** Table 4 presents the quantitative results on our proposed DAP-Test benchmark. Although this dataset is considered in-domain for our model, it serves as an essential evaluation of the effectiveness of large-scale data scaling and training strategies. Compared with previous state-of-the-art methods, DAP performs better than both DAC and Unik3D across all metrics. Specifically, DAP reduces the AbsRel from 0.2517 to 0.0781 and the RMSE from 10.563 to 6.804, while increasing $\delta_1$ from 0.6086 to 0.9307. These substantial improvements demonstrate the advantages of our data engine and unified framework in achieving accurate, metrically consistent, and robust depth estimation for panoramic outdoor scenes.

**Qualitative Comparison.** Figure 4 presents qualitative comparisons among DAC [15], Unik3D [32], and our DAP across diverse indoor and outdoor real-world scenes. Compared with existing approaches, DAP generates sharper object boundaries and more coherent global geometry, particularly in regions with large depth discontinuities or complex scene layouts. For indoor environments, our model accurately preserves fine structural details such as furniture edges and wall boundaries, while Unik3D and DAC often exhibit over-smoothed or distorted depth transitions. Notably, although Unik3D performs well on Deep360, its generalization to diverse real-world outdoor scenes is limited. In particular, existing methods fail to maintain stable depth structures in distant regions and sky areas, leading to distorted or collapsed predictions. In contrast, DAP maintains

on DAP-2M-Labeled. The 100 m setting offers the best overall balance, achieving AbsRel of 0.0793/0.0862 and $\delta_1$ of 0.9353/0.8719 on DAP-2M-Labeled and Deep360. Removing the mask notably degrades performance, confirming its effectiveness in filtering unreliable far-depth predictions and stabilizing training. Overall, the range mask head provides a flexible and reliable mechanism for maintaining metric consistency across diverse scene scales.

## 5. Conclusion

In this paper, we propose a panoramic metric-depth foundation model DAP, which is built through large-scale data scaling and a unified three-stage training pipeline. By combining reliable pseudo-labeling, geometry-aware design, and a plug-and-play range mask head, DAP achieves strong zero-shot generalization and state-of-the-art performance across indoor–outdoor benchmarks, with particularly robust and stable metric predictions in diverse real-world outdoor environments.

## References

[1] Hao Ai and Lin Wang. Elite360d: Towards efficient 360 depth estimation via semantic-and distance-aware bi-projection fusion. In *CVPR*, 2024. 3

[2] Hao Ai and Lin Wang. Elite360m: Efficient 360 multi-task learning via bi-projection fusion and cross-task collaboration. *arXiv preprint arXiv:2408.09336*, 2024.

[3] Hao Ai, Zidong Cao, Yan-Pei Cao, Ying Shan, and Lin Wang. Hrdfuse: Monocular 360deg depth estimation by collaboratively learning holistic-with-regional depth distributions. In *CVPR*, 2023. 3

[4] 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

[5] Jiayang Bai, Haoyu Qin, Shuichang Lai, Jie Guo, and Yanwen Guo. Glpanodepth: Global-to-local panoramic depth estimation. *TIP*, 2024. 3

[6] Yaniv Benny and Lior Wolf. Sphereuformer: A u-shaped transformer for spherical 360 perception. In *CVPR*, 2025. 3

[7] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024. 3

[8] Aleksei Bochkovskiy, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R Richter, and Vladlen Koltun. Depth pro: Sharp monocular metric depth in less than a second. In *ICLR*, 2025. 3

[9] Zidong Cao, Jinjing Zhu, Weiming Zhang, Hao Ai, Haotian Bai, Hengshuang Zhao, and Lin Wang. Panda: Towards panoramic depth anything with unlabeled panoramas and mobius spatial augmentation. In *CVPR*, 2025. 2, 3, 5, 6, 7

[10] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niebner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. In *3DV*, 2017. 6

[11] Jiaxi Deng, Yushen Wang, Haitao Meng, Zuoxun Hou, Yi Chang, and Gang Chen. Omnistereo: Real-time omnidirectional depth estimation with multiview fisheye cameras. In *CVPR*, 2025. 3

[12] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014. 5

[13] Haoran Feng, Dizhe Zhang, Xiangtai Li, Bo Du, and Lu Qi. Dit360: High-fidelity panoramic image generation via hybrid training. *arXiv preprint arXiv:2510.11712*, 2025. 2, 4

[14] Xian Ge, Yuling Pan, Yuhang Zhang, Xiang Li, Weijun Zhang, Dizhe Zhang, Zhaoliang Wan, Xin Lin, Xiangkai Zhang, Juntao Liang, et al. Airsim360: A panoramic simulation platform within drone view. *arXiv preprint arXiv:2512.02009*, 2025. 2, 3

[15] Yuliang Guo, Sparsh Garg, S Mahdi H Miangoleh, Xinyu Huang, and Liu Ren. Depth any camera: Zero-shot metric depth estimation from any camera. In *CVPR*, 2025. 2, 3, 6, 7, 8

[16] Jing He, Haodong Li, Wei Yin, Yixun Liang, Leheng Li, Kaiqiang Zhou, Hongbo Zhang, Bingbing Liu, and Ying-Cong Chen. Lotus: Diffusion-based visual foundation model for high-quality dense prediction. *arXiv preprint arXiv:2409.18124*, 2024. 3

[17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 3

[18] 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. *PAMI*, 2024. 3

[19] Hualie Jiang, Zhe Sheng, Siyu Zhu, Zilong Dong, and Rui Huang. Unifuse: Unidirectional fusion for 360° panorama depth estimation. *IEEE Robotics and Automation Letters*, 6: 1519–1526, 2021. 7

[20] Hualie Jiang, Zhe Sheng, Siyu Zhu, Zilong Dong, and Rui Huang. Unifuse: Unidirectional fusion for 360 panorama depth estimation. *RAL*, 2021. 3

[21] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing diffusion-based image generators for monocular depth estimation. In *CVPR*, 2024. 3

[22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 7

[23] Jongsung Lee, Harin Park, Byeong-Uk Lee, and Kyungdon Joo. Hush: Holistic panoramic 3d scene understanding using spherical harmonics. In *CVPR*, 2025. 3, 5

[24] Haodong Li, Chen Wang, Jiahui Lei, Kostas Daniilidis, and Lingjie Liu. Stereodiff: Stereo-diffusion synergy for video depth estimation. *arXiv preprint arXiv:2506.20756*, 2025. 3

[25] Haodong Li, Wangguangdong Zheng, Jing He, Yuhao Liu, Xin Lin, Xin Yang, Ying-Cong Chen, and Chunchao Guo. Da$^2$: Depth anything in any direction. *arXiv preprint arXiv:2509.26618*, 2025. 2, 3, 5, 7

[26] Ming Li, Xueqian Jin, Xuejiao Hu, Jingzhao Dai, Sidan Du, and Yang Li. Mode: Multi-view omnidirectional depth estimation with 360 cameras. In *European Conference on Computer Vision*, pages 197–213. Springer, 2022. 7

[27] Yuyan Li, Yuliang Guo, Zhixin Yan, Xinyu Huang, Ye Duan, and Liu Ren. Omnifusion: 360 monocular depth estimation via geometry-aware fusion. In *CVPR*, 2022. 3

[28] Xin Lin, Xian Ge, Dizhe Zhang, Zhaoliang Wan, Xianshun Wang, Xiangtai Li, Wenjie Jiang, Bo Du, Dacheng Tao, Ming-Hsuan Yang, et al. One flight over the gap: A survey from perspective to panoramic vision. *arXiv preprint arXiv:2509.04444*, 2025. 6

[29] Payal Mohadikar and Ye Duan. Omnidiffusion: Reformulating 360 monocular depth estimation using semantic and surface normal conditioned diffusion. In *WACV*, 2025. 3

[30] Chi-Han Peng and Jiayao Zhang. High-resolution depth estimation for 360deg panoramas through perspective and panoramic depth images registration. In *WACV*, 2023. 3

[31] Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. Unidepth: Universal monocular metric depth estimation. In *CVPR*, 2024. 3

[32] Luigi Piccinelli, Christos Sakaridis, Mattia Segu, Yung-Hsu Yang, Siyuan Li, Wim Abbeloos, and Luc Van Gool. Unik3d: Universal camera monocular 3d estimation. In *CVPR*, 2025. 2, 3, 6, 7, 8

[33] Luigi Piccinelli, Christos Sakaridis, Yung-Hsu Yang, Mattia Segu, Siyuan Li, Wim Abbeloos, and Luc Van Gool. Unidepthv2: Universal monocular metric depth estimation made simpler. *arXiv preprint arXiv:2502.20110*, 2025. 3

[34] Manuel Rey, Mingze Yuan Area, and Christian Richardt. 360monodepth: High-resolution 360 monocular depth estimation. In *CVPR*, 2022. 3

[35] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 3

[36] Zhijie Shen, Chunyu Lin, Kang Liao, Lang Nie, Zishuo Zheng, and Yao Zhao. Panoformer: panorama transformer for indoor 360° depth estimation. In *ECCV*, 2022. 3

[37] Zhijie Shen, Chunyu Lin, Lang Nie, Kang Liao, Weisi Lin, and Yao Zhao. Revisiting 360 depth estimation with panogabor: A new fusion perspective. *arXiv preprint arXiv:2408.16227*, 2024. 3

[38] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv preprint arXiv:2508.10104*, 2025. 4, 5

[39] Keisuke Tateno, Nassir Navab, and Federico Tombari. Distortion-aware convolutional filters for dense prediction in panoramic images. In *ECCV*, 2018. 3

[40] 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 *CVPR*, 2020. 3

[41] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Visual geometry grounded transformer. In *CVPR*, 2025. 7

[42] Jiyuan Wang, Chunyu Lin, Cheng Guan, Lang Nie, Jing He, Haodong Li, Kang Liao, and Yao Zhao. Jasmine: Harnessing diffusion prior for self-supervised depth estimation. *arXiv preprint arXiv:2503.15905*, 2025. 3

[43] Ning-Hsu Albert Wang and Yu-Lun Liu. Depth anywhere: Enhancing 360 monocular depth estimation via perspective distillation and unlabeled data augmentation. *NeurIPS*, 2024. 3

[44] Ning-Hsu Albert Wang and Yu-Lun Liu. Depth anywhere: Enhancing 360 monocular depth estimation via perspective distillation and unlabeled data augmentation. *NeurIPS*, 2024. 2

[45] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 4

[46] 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. In *CVPR*, 2025. 3, 5, 7

[47] Ruicheng Wang, Sicheng Xu, Yue Dong, Yu Deng, Jianfeng Xiang, Zelong Lv, Guangzhong Sun, Xin Tong, and Jiaolong Yang. Moge-2: Accurate monocular geometry with metric scale and sharp details. *arXiv preprint arXiv:2507.02546*, 2025. 3

[48] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024. 3

[49] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *NeurIPS*, 2024. 2, 3, 7

[50] 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 *ICCV*, 2023. 3

[51] Ilwi Yun, Chanyong Shin, Hyunku Lee, Hyuk-Jae Lee, and Chae Eun Rhee. Egformer: Equirectangular geometry-biased transformer for 360 depth estimation. In *ICCV*, 2023. 3

[52] Jia Zheng, Junfei Zhang, Jing Li, Rui Tang, Shenghua Gao, and Zihan Zhou. Structured3d: A large photo-realistic dataset for structured 3d modeling. In *ECCV*, 2020. 2, 5, 6

[53] Chuanqing Zhuang, Zhengda Lu, Yiqun Wang, Jun Xiao, and Ying Wang. Acdnet: Adaptively combined dilated convolution for monocular panorama depth estimation. In *AAAI*, 2022. 3