# Depth Anything V2

**Lihe Yang**[1]    **Bingyi Kang**[2][†]    **Zilong Huang**[2]
**Zhen Zhao**    **Xiaogang Xu**    **Jiashi Feng**[2]    **Hengshuang Zhao**[1][‡]

[1]HKU          [2]TikTok
[†]project lead    [‡]corresponding author

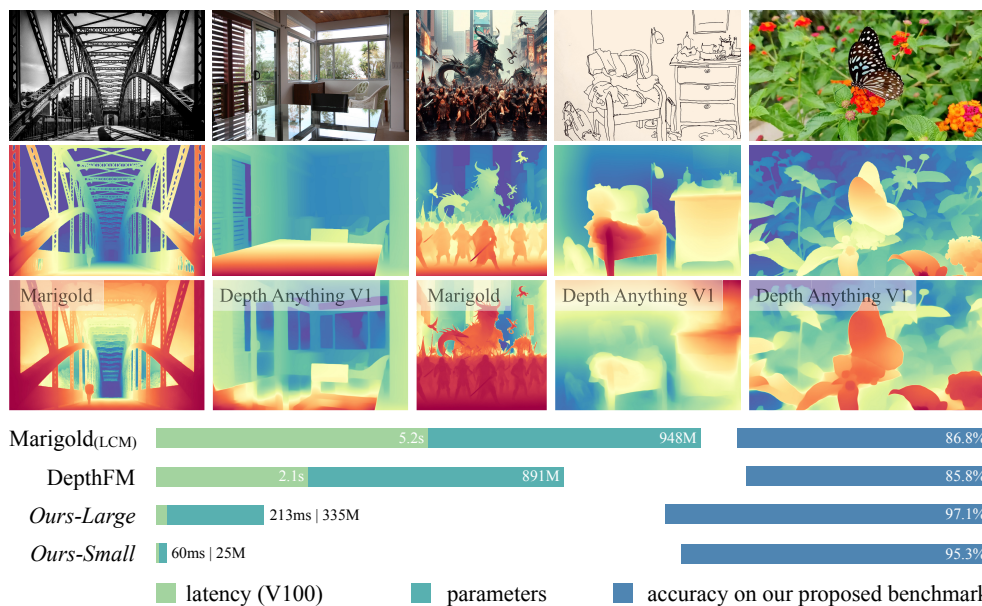https://depth-anything-v2.github.io

Figure 1: Depth Anything V2 significantly outperforms V1 [89] in robustness and fine-grained details. Compared with SD-based models [31, 25], it enjoys faster inference speed, fewer parameters, and higher depth accuracy.

## Abstract

This work presents *Depth Anything V2*. Without pursuing fancy techniques, we aim to reveal crucial findings to pave the way towards building a powerful monocular depth estimation model. Notably, compared with V1 [89], this version produces much finer and more robust depth predictions through three key practices: 1) replacing all labeled real images with synthetic images, 2) scaling up the capacity of our teacher model, and 3) teaching student models via the bridge of large-scale pseudo-labeled real images. Compared with the latest models [31] built on Stable Diffusion, our models are significantly more efficient (more than $10\times$ faster) and more accurate. We offer models of different scales (ranging from 25M to 1.3B params) to support extensive scenarios. Benefiting from their strong generalization capability, we fine-tune them with metric depth labels to obtain our metric depth models. In addition to our models, considering the limited diversity and frequent noise in current test sets, we construct a versatile evaluation benchmark with precise annotations and diverse scenes to facilitate future research.

---

Work done during an internship at TikTok.

Figure 2: *Robustness* (1st row, the misleading room layout) of Depth Anything V1 and *Fine-grained detail* (2nd row, the thin basketball net) of Marigold.

| Preferable Properties | Fine Detail | Transparent Objects | Reflections | Complex Scenes | Efficiency | Transferability |
|---|---|---|---|---|---|---|
| Marigold [31] | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Depth Anything V1 [89] | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ |
| Depth Anything V2 (Ours) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

Table 1: Preferable properties of a powerful monocular depth estimation model.

# 1   Introduction

Monocular depth estimation (MDE) is gaining increasing attention, due to its fundamental role in widespread downstream tasks. Precise depth information is not only favorable in classical applications, such as 3D reconstruction [47, 32, 93], navigation [82], and autonomous driving [80], but is also preferable in modern scenarios, *e.g.*, AI-generated content, including images [101], videos [39], and 3D scenes [87, 64, 68]. Therefore, there have been numerous MDE models [56, 7, 6, 95, 26, 38, 31, 89, 88, 25, 20, 52, 28] emerging recently, which are all capable of addressing open-world images.

From the aspect of model architecture, these works can be divided into two groups. One group [7, 6, 89, 28] is based on discriminative models, *e.g.*, BEiT [4] and DINOv2 [50], while the other [31, 20, 25] is based on generative models, *e.g.*, Stable Diffusion (SD) [59]. In Figure 2, we compare two representative works from the two categories respectively: Depth Anything [89] as a discriminative model and Marigold [31] as a generative model. It can be easily observed that Marigold is superior in modeling the details, while Depth Anything produces more robust predictions for complex scenes. Moreover, as summarized in Table 1, Depth Anything is more efficient and lightweight than Marigold, with different scales to choose from. Meantime, however, Depth Anything is vulnerable to transparent objects and reflections, which are the strengths of Marigold.

In this work, taking all these factors into account, we aim to build a more capable foundation model for monocular depth estimation that can achieve all the strengths listed in Table 1:

- produce robust predictions for complex scenes, including but not limited to complex layouts, transparent objects (*e.g.*, glass), reflective surfaces (*e.g.*, mirrors, screens) [15], *etc.*
- contain fine details (comparable to the details of Marigold) in the predicted depth maps, including but not limited to thin objects (*e.g.*, chair legs) [42], small holes, *etc.*
- provide varied model scales and inference efficiency to support extensive applications [82].
- be generalizable enough to be transferred (*i.e.*, fine-tuned) to downstream tasks, *e.g.*, Depth Anything V1 serves as the pre-trained model for all the leading teams in the 3rd MDEC[1] [72].

Since the nature of MDE is a discriminative task, we start from Depth Anything V1 [89], aiming to maintain its strengths and rectify its weaknesses. Intriguingly, we will demonstrate that, to achieve
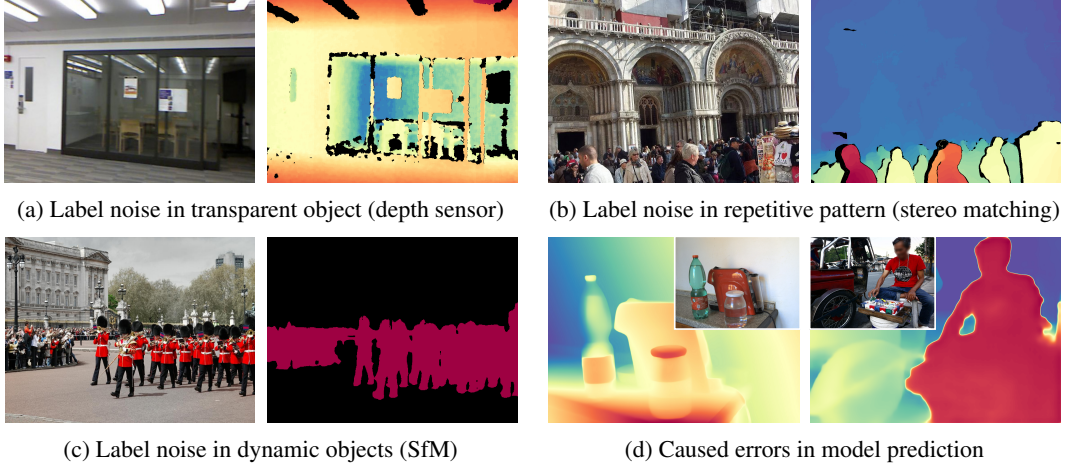
---

[1] https://jspenmar.github.io/MDEC

(a) Label noise in transparent object (depth sensor)   (b) Label noise in repetitive pattern (stereo matching)

(c) Label noise in dynamic objects (SfM)   (d) Caused errors in model prediction

Figure 3: Various noise in "GT" depth labels (a: NYU-D [70], b: HRWSI [83], c: MegaDepth [37]) and prediction errors in correspondingly trained models (d). Black regions are ignored during training.

such a challenging goal, no fancy or sophisticated techniques need to be developed. The most critical part is still **data**. It is indeed the same as the data-driven motivation of V1, which harnesses large-scale unlabeled data to speed up data scaling-up and increase the data coverage. In this work, we instead will first revisit its *labeled data* design, and then highlight the key role of unlabeled data.

We first present three key findings below. We will clarify them in detail in the following three sections.

**Q1 [Section 2]:** *Whether the coarse depth of MiDaS or Depth Anything come from the discriminative modeling itself? Is it a must to adopt the heavy diffusion-based modeling manner for fine details?*
**A1:** No, efficient discriminative models can also produce extremely fine details. The most critical modification is replacing all labeled real images with precise synthetic images.

**Q2 [Section 3]:** *Why do most prior works still stick to real images, if as A1 mentioned, synthetic images are already clearly superior to real images?*
**A2:** Synthetic images have their drawbacks, which are not trivial to address in previous paradigms.

**Q3 [Section 4]:** *How to avoid the drawbacks of synthetic images and also amplify its advantages?*
**A3:** Scale up the teacher model that is solely trained on synthetic images, and then teach (smaller) student models via the bridge of large-scale pseudo-labeled real images.

After the explorations, we successfully build a more capable MDE foundation model. However, we find current test sets [70] are too noisy to reflect the true strengths of MDE models. Thus, we further construct a versatile evaluation benchmark with precise annotations and diverse scenes (Section 6).

## 2   Revisiting the Labeled Data Design of Depth Anything V1

Building on the pioneering work of MiDaS [56, 7] in zero-shot MDE, recent studies tend to construct larger-scale training datasets in an effort to enhance estimation performance. Notably, Depth Anything V1 [89], Metric3D V1 [95] and V2 [28], as well as ZeroDepth [26], have amassed 1.5M, 8M, 16M, and 15M labeled images from various sources for training, respectively. However, few studies have critically examined this trend: *is such a huge amount of labeled images truly advantageous?*

Before answering it, let us first dig into the potentially overlooked drawbacks of **real** labeled images.

**Two disadvantages of real labeled data.** 1) *Label noise*, *i.e.*, *inaccurate labels* in depth maps. Stemming from the limitations inherent in various collection procedures, real labeled data inevitably contain inaccurate estimations. Such inaccuracies can arise from various factors, such as the inability of depth sensors to accurately capture the depth of transparent objects (Figure 3a), the vulnerability of stereo matching algorithms to textureless or repetitive patterns (Figure 3b), and the susceptible nature of SfM methods in handling dynamic objects or outliers (Figure 3c). 2) *Ignored details*. These real datasets often overlook certain details in their depth maps. As depicted in Figure 4a, the depth

3

(a) Coarse depth of real data (HRWSI [83], DIML [14])    (b) Depth of synthetic data (Hypersim [58], vKITTI [9])



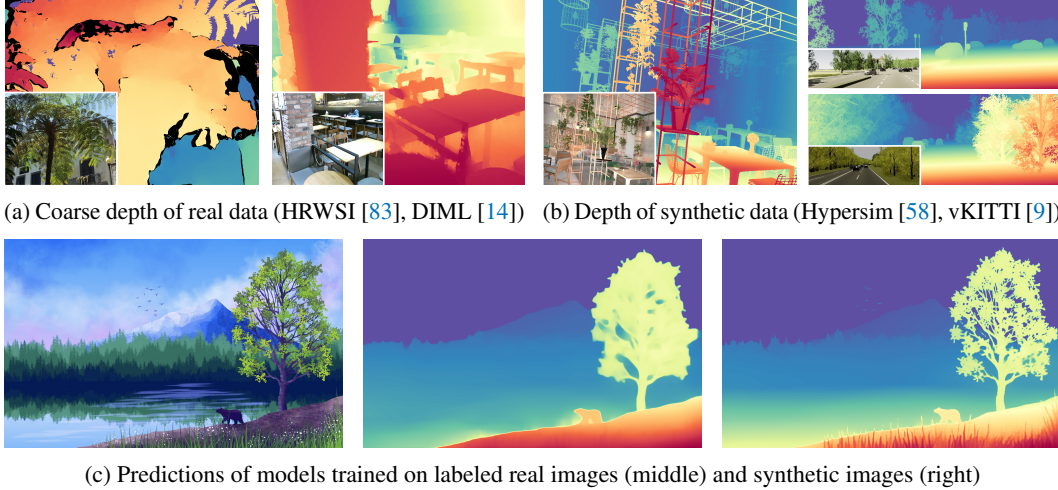(c) Predictions of models trained on labeled real images (middle) and synthetic images (right)

Figure 4: Depth labels of real images (a) and synthetic images (b), and the corresponding model predictions (c). The labels of synthetic images are highly precise, and so are their trained models.

representation of the tree and chair is notably coarse. These datasets struggle to provide detailed supervision at object boundaries or within thin holes, resulting in over-smoothed depth predictions, as seen in the middle of Figure 4c. Hence, these noisy labels are so unreliable that the learned models make similar mistakes (Figure 3d). For example, MiDaS and Depth Anything V1 obtain poor scores of 25.9% and 53.5% respectively in the Transparent Surface Challenge [54] (more details in Table 12: our V2 achieves a competitive score of 83.6% in a zero-shot manner).

To overcome the above problems, we decide to change our training data and seek images with substantially better annotation. Inspired by several recent SD-based studies [31, 20, 25], that exclusively utilize synthetic images with complete depth information for training, we extensively check the label quality of synthetic images and note their potential to mitigate the drawbacks discussed above.

**Advantages of synthetic images.** Their depth labels are highly precise in two folds. 1) All fine details (*e.g.*, boundaries, thin holes, small objects, *etc.*) are correctly labeled. As demonstrated in Figure 4b, even all thin mesh structures and leaves are annotated with true depth. 2) We can obtain the actual depth of challenging transparent objects and reflective surfaces, *e.g.*, the vase on the table in Figure 4b. In a word, the depth of synthetic images is truly "GT". In the right side of Figure 4c, we show the fine-grained prediction of a MDE model trained on synthetic images. Moreover, we can quickly enlarge synthetic training images by collecting from graphics engines [58, 63, 53], which would not cause any privacy or ethical concerns, as compared to real images.

## 3    Challenges in Using Synthetic Data

If synthetic data are so advantageous, why are real data still dominating MDE? In this section, we identify **two limitations of synthetic images** that hinder them from being easily used in reality.

**Limitation 1.** There exists *distribution shift* between synthetic and real images. Although current graphics engines strive for photorealistic effects, their style and color distributions still evidently differ from real images. Synthetic images are too "clean" in color and "ordered" in layout, while real images contain more randomness. For instance, when comparing the images in Figure 4a and 4b, we can immediately distinguish the synthetic ones. Such distribution shift makes models struggle to transfer from synthetic to real images, even if the two data sources share similar layouts [57, 9].

**Limitation 2.** Synthetic images have *restricted scene coverage*. They are iteratively sampled from graphics engines with pre-defined fixed scene types, *e.g.*, "living room" and "street scene". Consequently, despite the astonishing precision of Hypersim [58] or Virtual KITTI [9] (Figure 4b), we cannot expect models trained on them to generalize well in real-world scenes like "crowded people". In contrast, some real datasets constructed from web stereo images (*e.g.*, HRWSI [83]) or monocular videos (*e.g.*, MegaDepth [37]), can cover extensive real-world scenes.
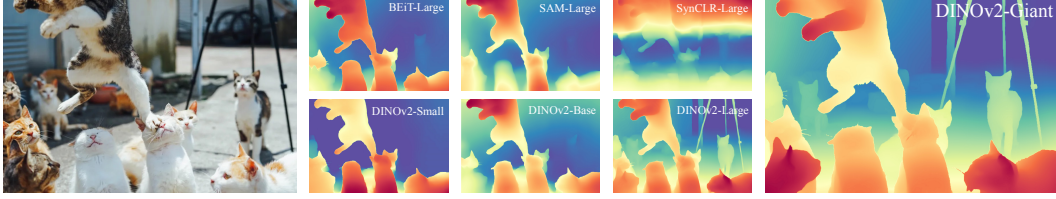
4

Figure 5: Qualitative comparison of different vision encoders on synthetic-to-real transfer. Only DINOv2-G produces a satisfying prediction. For quantitative comparisons, please refer to Section B.6.
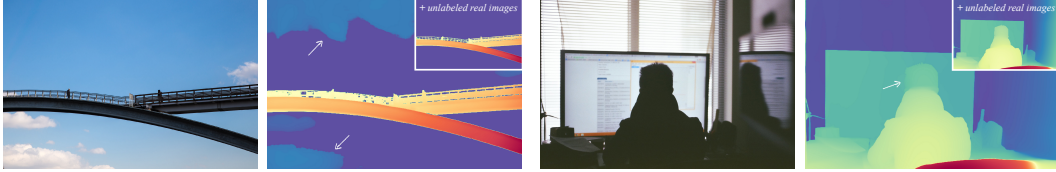


Figure 6: Failure cases of the most capable DINOv2-G model when purely trained on synthetic images. Left: the sky should be ultra far. Right: the depth of the head is not consistent with the body.

**Therefore, synthetic-to-real transfer is non-trivial in MDE.** To validate this claim, we conduct a pilot study to learning MDE models solely on synthetic images with four popular pre-trained encoders, including BEiT [4], SAM [33], SynCLR [75], and DINOv2 [50]. As illustrated in Figure 5, only DINOv2-G achieves satisfying results. All other model serials, as well as smaller DINOv2 models, suffer from severe generalization issues. This pilot study seems to give a straightforward solution to employing synthetic data in MDE, *i.e.*, building on the largest DINOv2 encoder, and relying on its inherent generalization ability. However, this naive solution faces two problems. First, DINOv2-G frequently encounters failure cases when the patterns of real test images are rarely presented in synthetic training images. In Figure 6, we can clearly observe incorrect depth predictions for the sky (cloud) and the human head. Such failures can be expected as our synthetic training sets do not include diverse sky patterns or humans. Moreover, most applications cannot accommodate the resource-intensive DINOv2-G model (1.3B) in terms of storage and inference efficiency. Actually, the smallest model in Depth Anything V1 is used most widely due to its real-time speed.

To alleviate the generalization issue, some works [7, 89, 28] use a combined training set of real and synthetic images. Unfortunately, as shown in Section B.9, the coarse depth map of real images is destructive to fine-grained prediction. Another potential solution is to collect more synthetic images, which is unsustainable as creating graphic engines mimicking every real-world scenario is intractable. Therefore, a reliable solution is demanding in building MDE models with synthetic data. In this paper, we will close this gap and present a roadmap that solves the preciseness and robustness dilemma *without any trade-offs*, and applicable to *any model scale*.

## 4 Key Role of Large-Scale Unlabeled Real Images

Our solution is straightforward: incorporating *unlabeled real* images. Our most capable MDE model, based on DINOv2-G, is initially trained purely on high-quality synthetic images. Then it assigns pseudo depth labels on unlabeled real images. Lastly, our new models are solely trained with large-scale and precisely pseudo-labeled images. Depth Anything V1 [89] has highlighted the importance of large-scale unlabeled real data. Here, in our special context of synthetic labeled images, we will demonstrate its *indispensable* role in more details from three perspectives.

**Bridge the domain gap.** As aforementioned, due to the distribution shift, directly transferring from synthetic training images to real test images is challenging. But if we can leverage extra real images as an intermediate learning target, the process will be more reliable. Intuitively, after explicitly training on pseudo-labeled real images, models can be more familiar with real-world data distribution. Compared with manually annotated images, our auto-generated pseudo labels are much more fine-grained and complete, as visualized in Figure 17.
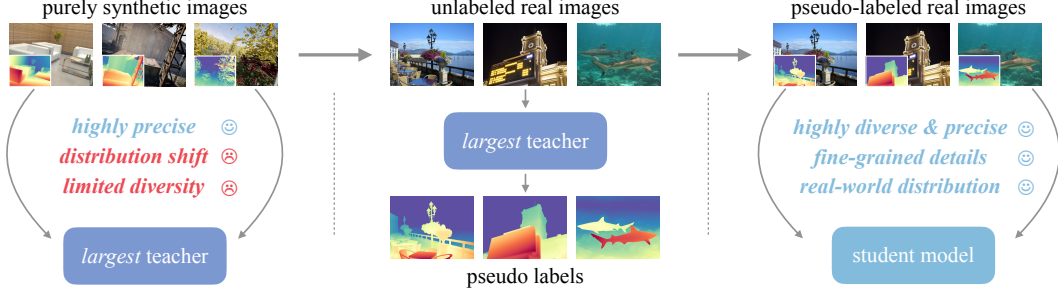
Figure 7: Depth Anything V2. We first train the most capable teacher on precise synthetic images. Then, to mitigate the distribution shift and limited diversity of synthetic data, we annotate unlabeled real images with the teacher. Finally, we train student models on high-quality pseudo-labeled images.

**Enhance the scene coverage.** The diversity of synthetic images is limited, without including enough real-world scenes. Nevertheless, we can easily cover numerous distinct scenes by incorporating large-scale unlabeled images from public datasets. Moreover, synthetic images are indeed very redundant due to being repetitively sampled from pre-defined videos. In comparison, unlabeled real images are clearly distinguished and very informative. By training on sufficient images and scenes, models not only demonstrate stronger zero-shot MDE capability (as shown in Figure 6 "+ *unlabeled real images*"), but they can also serve as better pre-trained sources for downstream related tasks [72].

**Transfer knowledge from the most capable model to smaller ones.** We have shown in Figure 5, that smaller models cannot directly benefit from synthetic-to-real transfer by themselves. However, armed with large-scale unlabeled real images, they can learn to mimic the high-quality predictions of the most capable model, similar to knowledge distillation [27]. But differently, our distillation is enforced at the label level via extra unlabeled real data, instead of at the feature or logit level with original labeled data. This practice is safer because there is evidence showing feature-level distillation is not always beneficial, especially when the teacher-student scale gap is huge [48]. Finally, as supported in Figure 16, unlabeled images boost the robustness of our smaller models tremendously.

# 5 Depth Anything V2

## 5.1 Overall Framework

According to all the above analysis, our final pipeline to train Depth Anything V2 is clear (Figure 7). It consists of three steps:

- train a reliable teacher model based on DINOv2-G *purely* on high-quality *synthetic* images.
- produce precise pseudo depth on large-scale unlabeled *real* images.
- train final student models on *pseudo-labeled real* images for robust generalization (we will show the synthetic images are not necessary in this step).

We will release four student models, based on DINOv2 small, base, large, and giant, respectively.

## 5.2 Details

As shown in Table 7, we use five precise synthetic datasets (595K images) and eight large-scale pseudo-labeled real datasets (62M images) for training. Same as V1 [89], for each pseudo-labeled sample, we ignore its top-$n$-largest-loss regions during training, where $n$ is set as 10%. We consider them as potentially noisy pseudo labels. Similarly, our models produce affine-invariant inverse depth[2]. We use two loss terms for optimization on labeled images: a scale- and shift-invariant loss $\mathcal{L}_{ssi}$ and a gradient matching loss $\mathcal{L}_{gm}$. These two objective functions are not new, as they are proposed by MiDaS [56]. But differently, we find $\mathcal{L}_{gm}$ is super beneficial to the depth sharpness when using synthetic images (Section B.7). On pseudo-labeled images, we follow V1 to add an additional feature alignment loss to preserve informative semantics from pre-trained DINOv2 encoders.

---

[2]To offer capable *metric depth* models, we further fine-tune our basic models with metric depth (Section 7.3).
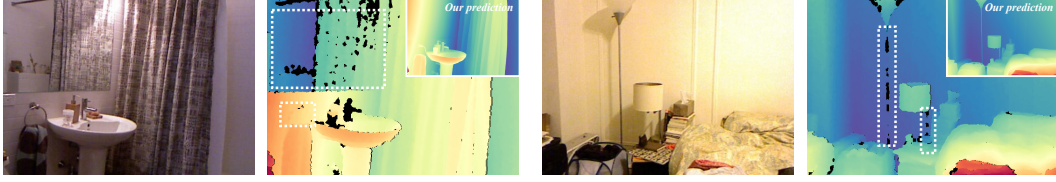
Figure 8: Visualization of widely adopted but indeed noisy test benchmark [70]. As highlighted, the depth of the mirror and thin structures are incorrect (black pixels are ignored). In comparison, our model predictions are accurate. The noise will cause better models instead achieve lower scores.
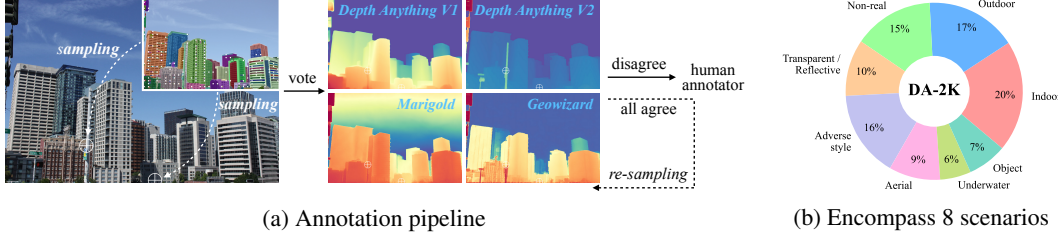


(a) Annotation pipeline

(b) Encompass 8 scenarios

Figure 9: Our proposed evaluation benchmark DA-2K. (a) The annotation pipeline for relative depth between two points. Points are sampled based on SAM [33] mask predictions. Disagreed pairs among four depth models will be popped out for annotators to label. (b) Detail of our scenario coverage.

# 6 A New Evaluation Benchmark: DA-2K

## 6.1 Limitations in Existing Benchmarks

In Section 2, we demonstrated that commonly used real training sets have noisy depth labels. Here, we further argue that widely adopted *test benchmarks* are also noisy. Figure 8 illustrates incorrect annotations for mirrors and thin structures on NYU-D [70] despite using specialized depth sensors. Such frequent label noise makes the reported metrics of powerful MDE models not reliable anymore.

Apart from label noise, another drawback of these benchmarks is *limited diversity*. Most of them were originally proposed for a single scene. For example, NYU-D [70] focuses on a few indoor rooms, while KITTI [24] only contains several street scenes. Performance on these benchmarks may not reflect real-world reliability. Ideally, we expect MDE models can handle any unseen scenes robustly.

The last problem in these existing benchmarks is *low resolution*. They mostly provide images with a resolution of around $500{\times}500$. But with modern cameras, we usually require precise depth estimation for higher-resolution images, *e.g.*, $1000{\times}2000$. It remains unclear whether the conclusions drawn from these low-resolution benchmarks can be safely transferred to high-resolution benchmarks.

## 6.2 DA-2K

Considering the above three limitations, we aim to construct a versatile evaluation benchmark for relative monocular depth estimation, that can 1) provide *precise* depth relationship, 2) cover *extensive* scenes, and 3) contain mostly *high-resolution* images for modern usage. Indeed, it is impractical for humans to annotate the depth of each pixel, especially for in-the-wild images. Thus, following DIW [11], we annotate *sparse* depth pairs for each image. Generally, given an image, we can select two pixels on it, and decide their relative depth between them (*i.e.*, which pixel is closer).

Concretely, we employ two distinct pipelines to select pixel pairs. In the first pipeline, as shown in Figure 9a, we use SAM [33] to automatically predict object masks. Instead of using the masks, we leverage key points (pixels) that prompt out them. We randomly sample two key pixels and query four expert models ([89, 31, 20] and ours) to vote on their relative depth. If there is disagreement, the pair will be sent to human annotators to decide the true relative depth. Due to potential ambiguity, annotators can skip any pair. However, there may be cases where all models incorrectly predict challenging pairs, and they are not flagged. To address this, we introduce a second pipeline, where we carefully analyze images and manually identify challenging pairs.

| Method | Encoder | KITTI [24] | | NYU-D [70] | | Sintel [8] | | ETH3D [62] | | DIODE [76] | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ |
| MiDaS V3.1 [7] | ViT-L | 0.127 | 0.850 | 0.048 | 0.980 | 0.587 | 0.699 | 0.139 | 0.867 | 0.075 | 0.942 |
| Depth Anything V1 [89] | ViT-S | 0.080 | 0.936 | 0.053 | 0.972 | 0.464 | 0.739 | 0.127 | **0.885** | 0.076 | 0.939 |
| | ViT-B | 0.080 | 0.939 | 0.046 | 0.979 | **0.432** | 0.756 | **0.126** | 0.884 | 0.069 | 0.946 |
| | ViT-L | 0.076 | 0.947 | **0.043** | **0.981** | 0.458 | 0.760 | 0.127 | 0.882 | 0.066 | 0.952 |
| **Depth Anything V2** | ViT-S | 0.078 | 0.936 | 0.053 | 0.973 | 0.500 | 0.718 | 0.142 | 0.851 | 0.073 | 0.942 |
| | ViT-B | 0.078 | 0.939 | 0.049 | 0.976 | 0.495 | 0.734 | 0.137 | 0.858 | 0.068 | 0.950 |
| | ViT-L | **0.074** | 0.946 | 0.045 | 0.979 | 0.487 | 0.752 | 0.131 | 0.865 | 0.066 | 0.952 |
| | ViT-G | 0.075 | **0.948** | 0.044 | 0.979 | 0.506 | **0.772** | 0.132 | 0.862 | **0.065** | **0.954** |

Table 2: Zero-shot *relative* depth estimation. Better: AbsRel ↓ , $\delta_1$ ↑. Solely from the metrics, Depth Anything V2 is better than MiDaS, but merely comparable with V1. But indeed, the focus and strengths of our V2 (*e.g.*, fine-grained details, robust to complex layouts, transparent objects, *etc.*) cannot be correctly reflected on these benchmarks. Similar results (*i.e.*, better model but worse score) are also observed in [7, 28].

| Method | Community Models | | | | Depth Anything V2 (Ours) | | | |
|---|---|---|---|---|---|---|---|---|
| | Marigold [31] | Geowizard [20] | DepthFM [25] | Depth Anything V1 [89] | ViT-S | ViT-B | ViT-L | ViT-G |
| Accuracy (%) | 86.8 | 88.1 | 85.8 | 88.5 | 95.3 | 97.0 | 97.1 | **97.4** |

Table 3: Performance on our proposed DA-2K evaluation benchmark, which encompasses eight representative scenarios. Even our most lightweight model is superior to all other community models.

To ensure preciseness, all annotations are triple-checked by the other two annotators. To ensure diversity, we first summarize eight important application scenarios of MDE (Figure 9b), and ask GPT-4 to produce diverse keywords related to each scenario. We then use these keywords to download corresponding images from Flickr. Finally, we annotate 1K images with 2K pixel pairs in total. Limited by space, please refer to Section C for details and comparisons with DIW [11].

**Position of DA-2K.** Despite the advantages, we *do not* expect DA-2K to *replace* current benchmarks. Accurate sparse depth is still far from the precise dense depth required for scene reconstruction. However, DA-2K can be considered a prerequisite for accurate dense depth. As such, we believe DA-2K can serve as *a valuable supplement* to existing benchmarks due to its extensive scene coverage and precision. It can also serve as a quick prior validation for users selecting community models for specific scenarios covered in DA-2K. Lastly, we believe it is also a potential testbed for the 3D awareness of future multimodal LLMs [41, 21, 3].

# 7 Experiment

## 7.1 Implementation details

Follow Depth Anything V1 [89], we use DPT [55] as our depth decoder, built on DINOv2 encoders. All images are trained at the resolution of 518×518 by resizing the shorter size to 518 followed by a random crop. When training the teacher model on synthetic images, we use a batch size of 64 for 160K iterations. In the third stage of training on pseudo-labeled real images, the model is trained with a batch size of 192 for 480K iterations. We use the Adam optimizer and set the learning rate of the encoder and the decoder as 5e-6 and 5e-5, respectively. In both training stages, we do not balance the training datasets, but simply concatenate them. The weight ratio of $\mathcal{L}_{ssi}$ and $\mathcal{L}_{gm}$ is set as 1:2.

## 7.2 Zero-Shot Relative Depth Estimation

**Performance on conventional benchmarks.** Since our model predicts affine-invariant *inverse* depth, for fairness, we compare with Depth Anything V1 [89] and MiDaS V3.1 [7] on five unseen test datasets. As shown in Table 2, our results are superior to MiDaS and comparable to V1 [89]. We are slightly inferior to V1 in *metrics* on two of the datasets. However, the plain metrics on these datasets are not the focus of this paper. This version aims to produce fine-grained predictions for thin structures and robust predictions for complex scenes, transparent objects, *etc.*. Improvement in these dimensions cannot be correctly reflected in current benchmarks.

**Performance on our proposed benchmark DA-2K.** As shown in Table 3, on our proposed benchmark with diverse scenes, even our smallest model is significantly better than other heavy SD-based

Table 4 (a) NYU-D dataset and (b) KITTI dataset:

| Method | Higher is better ↑ | | | Lower is better ↓ | | | Method | Higher is better ↑ | | | Lower is better ↓ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\delta_1$ | $\delta_2$ | $\delta_3$ | AbsRel | RMSE | log10 | | $\delta_1$ | $\delta_2$ | $\delta_3$ | AbsRel | RMSE | RMSE log |
| AdaBins [5] | 0.903 | 0.984 | 0.997 | 0.103 | 0.364 | 0.044 | AdaBins [5] | 0.964 | 0.995 | 0.999 | 0.058 | 2.360 | 0.088 |
| DPT [55] | 0.904 | 0.988 | 0.998 | 0.110 | 0.357 | 0.045 | P3Depth [51] | 0.953 | 0.993 | 0.998 | 0.071 | 2.842 | 0.103 |
| P3Depth [51] | 0.898 | 0.981 | 0.996 | 0.104 | 0.356 | 0.043 | NeWCRFs [99] | 0.974 | 0.997 | 0.999 | 0.052 | 2.129 | 0.079 |
| SwinV2 [44] | 0.949 | 0.994 | 0.999 | 0.083 | 0.287 | 0.035 | SwinV2 [44] | 0.977 | 0.998 | 1.000 | 0.050 | 1.966 | 0.075 |
| AiT [49] | 0.954 | 0.994 | 0.999 | 0.076 | 0.275 | 0.033 | NDDepth [66] | 0.978 | 0.998 | 0.999 | 0.050 | 2.025 | 0.075 |
| VPD [102] | 0.964 | 0.995 | 0.999 | 0.069 | 0.254 | 0.030 | GEDepth [91] | 0.976 | 0.997 | 0.999 | 0.048 | 2.044 | 0.076 |
| IEBins [67] | 0.936 | 0.992 | 0.998 | 0.087 | 0.314 | 0.038 | IEBins [67] | 0.978 | 0.998 | 0.999 | 0.050 | 2.011 | 0.075 |
| ZoeDepth [6] | 0.951 | 0.994 | 0.999 | 0.077 | 0.282 | 0.033 | ZoeDepth [6] | 0.971 | 0.996 | 0.999 | 0.054 | 2.281 | 0.082 |
| Ours (ViT-S) | 0.961 | 0.996 | 0.999 | 0.073 | 0.261 | 0.032 | Ours (ViT-S) | 0.973 | 0.997 | 0.999 | 0.053 | 2.235 | 0.081 |
| Ours (ViT-B) | 0.977 | 0.997 | 1.000 | 0.063 | 0.228 | 0.027 | Ours (ViT-B) | 0.979 | 0.998 | 1.000 | 0.048 | 1.999 | 0.072 |
| Ours (ViT-L) | **0.984** | **0.998** | **1.000** | **0.056** | **0.206** | **0.024** | Ours (ViT-L) | **0.983** | **0.998** | **1.000** | **0.045** | **1.861** | **0.067** |

(a) NYU-D dataset      (b) KITTI dataset

Table 4: Fine-tuning our Depth Anything V2 pre-trained encoder to in-domain metric depth estimation, *i.e.*, training and test images share the same domain. All compared methods use the encoder size close to ViT-L.

| Encoder | $\mathcal{D}^l$ | $\mathcal{D}^u$ | KITTI [24] | | NYU-D [70] | | Sintel [8] | | ETH3D [62] | | DIODE [76] | | DA-2K |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | Acc (%) |
| ViT-S | ✓ | | 0.104 | 0.889 | 0.084 | 0.928 | 0.518 | 0.702 | 0.155 | 0.827 | 0.087 | 0.926 | 89.8 |
| | ✓ | ✓ | 0.085 | 0.928 | 0.054 | 0.971 | **0.491** | **0.723** | 0.143 | 0.849 | 0.074 | 0.941 | 94.1 |
| | | ✓ | **0.078** | **0.936** | **0.053** | **0.973** | 0.500 | 0.718 | **0.142** | **0.851** | **0.073** | **0.942** | **95.3** |
| ViT-B | ✓ | | 0.094 | 0.912 | 0.062 | 0.963 | 0.618 | 0.715 | 0.148 | 0.842 | 0.076 | 0.940 | 92.9 |
| | ✓ | ✓ | 0.080 | 0.938 | **0.049** | **0.976** | 0.515 | 0.732 | **0.137** | **0.859** | 0.068 | 0.950 | 96.7 |
| | | ✓ | **0.078** | **0.939** | **0.049** | **0.976** | **0.495** | **0.734** | **0.137** | 0.858 | 0.068 | 0.950 | **97.0** |
| ViT-L | ✓ | | 0.081 | 0.937 | 0.048 | 0.976 | 0.516 | 0.731 | 0.133 | 0.864 | 0.071 | 0.949 | 96.0 |
| | ✓ | ✓ | 0.075 | **0.947** | **0.045** | **0.979** | 0.542 | 0.741 | **0.130** | **0.866** | **0.066** | **0.953** | **97.3** |
| | | ✓ | **0.074** | 0.946 | **0.045** | **0.979** | **0.487** | **0.752** | 0.131 | 0.865 | **0.066** | 0.952 | 97.1 |
| Teacher model (ViT-G) | | | 0.075 | 0.947 | 0.044 | 0.979 | 0.530 | 0.767 | 0.131 | 0.865 | 0.066 | 0.954 | 97.4 |

Table 5: Importance of pseudo-labeled (unlabeled) real images ($\mathcal{D}^u$). $\mathcal{D}^l$: precisely labeled synthetic images.

models, *e.g.*, Marigold [31] and Geowizard [20]. Our most capable model achieves 10.6% higher accuracy than Margold in terms of relative depth discrimination. Please refer to Table 14 for the comprehensive per-scenario performance of our models.

## 7.3 Fine-tuned to Metric Depth Estimation

To validate the generalization ability of our model, we transfer its encoder to the downstream metric depth estimation task. First, same as V1 [89], we follow the ZoeDepth [6] pipeline, but replace its MiDaS [7] encoder with our pre-trained encoder. As shown in Table 4, we achieve significant improvements over previous methods on both NYU-D and KITTI datasets. Notably, even our most lightweight model which is based on ViT-S, is superior to other models built on ViT-L [6].

Although the reported metrics look impressive, models trained on NYUv2 or KITTI fail to produce fine-grained depth prediction and are not robust to transparent objects, due to the inherent noise in training sets. Therefore, to satisfy real-world applications such as multi-view synthesis, we fine-tune our powerful encoder on Hypersim [58] and Virtual KITTI [9] synthetic datasets, for indoor and outdoor metric depth estimation, respectively. We will release these two metric depth models. Please refer to Figure 15 for qualitative comparisons with the previous ZoeDepth method.

## 7.4 Ablation Study

Limited by space, we defer most of our ablations to the appendix except for two on pseudo labels.

**Importance of large-scale pseudo-labeled real images.** As shown in Table 5, compared with solely trained on synthetic images, our models are greatly enhanced by incorporating pseudo-labeled real images. Different from Depth Anything V1 [89], we further attempt to remove the synthetic images during training student models. We find this can even lead to slightly better results for smaller models (*e.g.*, ViT-S and ViT-B). So we finally choose to train student models purely on pseudo-labeled images. This observation is indeed similar to SAM [33] that only releases its pseudo-labeled masks.

| Label Source | KITTI [24] | | NYU-D [70] | | Sintel [8] | | ETH3D [62] | | DIODE [76] | | DA-2K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | Acc (%) |
| Manual Label | 0.122 | 0.882 | 0.074 | 0.952 | 0.581 | 0.693 | 0.159 | 0.832 | 0.126 | 0.890 | 80.2 |
| Pseudo Label | **0.099** | **0.901** | **0.062** | **0.963** | **0.514** | **0.701** | **0.147** | **0.843** | **0.084** | **0.929** | **89.7** |

Table 6: Comparison between originally manual label and our produced pseudo label on the DIML dataset [14]. Our produced pseudo labels are of much higher quality than the manual labels provided by DIML.

**Pseudo label *vs*. manual label on real labeled images.** We have demonstrated before in Figure 4a that existing labeled real datasets are very noisy. Here we conduct a quantitative comparison. We use real images from the DIML [14] dataset, and compare the transferring performance under its original manual label and our produced pseudo label respectively. We can observe in Table 6 that the model trained with pseudo labels is significantly better than the manual-label counterpart. The huge gap indicates the high quality of our pseudo labels and the rich noise in current labeled real datasets.

## 8   Related Work

**Monocular depth estimation.** Early works [18, 19, 5] focus on the in-domain metric depth estimation, where training and test images must share the same domain [70, 24]. Due to their restricted application scenarios, recently there has been increasing attention on zero-shot relative monocular depth estimation. Among them, some works address this task through better modeling manners, *e.g.*, using Stable Diffusion [59] as a depth denoiser [31, 25, 20]. Other works [94, 96, 89] focus on the data-driven perspective. For example, MiDaS [56, 55, 7] and Metric3D [95] collect 2M and 8M labeled images respectively. Aware of the difficulty of scaling up labeled images, Depth Anything V1 [89] leverages 62M unlabeled images to enhance the model's robustness. In this work, differently, we point out multiple limitations in widely used labeled real images. We thus especially highlight the necessity of resorting to synthetic images to ensure depth preciseness. Meantime, to tackle the generalization issue caused by synthetic images, we adopt both data-driven (large-scale pseudo-labeled real images) and model-driven (scaling up the teacher model) strategies.

**Learning from unlabeled real images.** How to learn informative representations from unlabeled images is widely studied in the field of semi-supervised learning [36, 86, 71, 90]. However, they focus on academic benchmarks [34] which only allow usage of small-scale labeled and unlabeled images. In comparison, we study a real-world application scenario, *i.e.*, how to further boost the baseline of 0.6M labeled images with 62M unlabeled images. Moreover, distinguished from Depth Anything V1 [89], we exhibit the indispensable role of unlabeled real images especially when we replace all labeled real images with synthetic images [22, 23, 61]. We demonstrate "precise synthetic data + pseudo-labeled real data" is a more promising roadmap than labeled real data.

**Knowledge distillation.** We distill transferable knowledge from our most capable teacher model to smaller models. This is similar to the core spirit of knowledge distillation (KD) [27]. But we are also fundamentally different in that we perform distillation at the *prediction level* through extra *unlabeled* real images, while KD [2, 73, 100] typically studies better distillation strategies at the *feature or logit* level through *labeled* images. We aim to reveal the importance of large-scale unlabeled data and larger teacher model, rather than delicate loss designs [43, 69] or distillation pipelines [10]. Moreover, it is indeed non-trivial and risky to directly distill feature representations between two models with a tremendous scale gap [48]. In comparison, our pseudo-label distillation is easier and safer, even from a model of 1.3B parameters to a model of 25M parameters.

## 9   Conclusion

In this work, we present *Depth Anything V2*, a more powerful foundation model for monocular depth estimation. It is capable of 1) providing robust and fine-grained depth prediction, 2) supporting extensive applications with varied model sizes (from 25M to 1.3B parameters), and 3) being easily fine-tuned to downstream tasks as a promising model initialization. We reveal crucial findings to pave the way towards building a strong MDE model. Furthermore, realizing the poor diversity and rich noise in existing test sets, we construct a versatile evaluation benchmark DA-2K, covering diverse high-resolution images with precise and challenging sparse depth labels.

# Appendix

For a thorough understanding and visualization of our Depth Anything V2, we compile a comprehensive appendix. The following table of contents will direct you to specific sections of interest.

# Contents

# A   Sources of Training Data

As listed in Table 7, we replace all labeled real datasets in Depth Anything V1 [89] with five synthetic datasets for label preciseness. Then, to mitigate the issues of distribution shift and limited diversity caused by synthetic images, we further leverage eight large-scale public datasets, comprising 62M real images with great diversity. We only use their raw images, and assign depth to them with our most capable teacher model. Student models are trained purely on these pseudo-labeled real images.

# B   Experiments

## B.1   Fine-tuned to semantic segmentation

Similar to the practice in metric MDE, we further fine-tune our pre-trained encoder to downstream semantic segmentation task to especially examine its semantic awareness. As demonstrated in Table 8,

| Dataset | Indoor | Outdoor | # Images |
|---|---|---|---|
| **Precise *Synthetic* Images (595K)** | | | |
| BlendedMVS [92] | ✓ | ✓ | 115K |
| Hypersim [58] | ✓ | | 60K |
| IRS [77] | ✓ | | 103K |
| TartanAir [79] | ✓ | ✓ | 306K |
| VKITTI 2 [9] | | ✓ | 20K |
| **Pseudo-labeled *Real* Images (62M)** | | | |
| BDD100K [97] | | ✓ | 8.2M |
| Google Landmarks [81] | | ✓ | 4.1M |
| ImageNet-21K [60] | ✓ | ✓ | 13.1M |
| LSUN [98] | ✓ | | 9.8M |
| Objects365 [65] | ✓ | ✓ | 1.7M |
| Open Images V7 [35] | ✓ | ✓ | 7.8M |
| Places365 [103] | ✓ | ✓ | 6.5M |
| SA-1B [33] | ✓ | ✓ | 11.1M |

Table 7: Our training data sources.

| Method | Encoder | mIoU |
|---|---|---|
| DDP [30] | Swin-S [45] | 82.4 |
| Depth Anything V2 | Small | **82.9** |
| DDP [30] | Swin-B [45] | 82.5 |
| Depth Anything V2 | Base | **83.9** |
| Segmenter [74] | ViT-L [17] | 82.2 |
| SegFormer [85] | MiT-B5 [85] | 82.4 |
| Mask2Former [13] | Swin-L [45] | 83.3 |
| OneFormer [29] | Swin-L [45] | 83.0 |
| OneFormer [29] | ConvNeXt-XL [46] | 83.6 |
| DDP [30] | ConvNeXt-L [46] | 83.2 |
| Depth Anything V2 | Large | **85.6** |

(a) Cityscapes dataset

| Method | Encoder | mIoU |
|---|---|---|
| UperNet [84] | InternImage-S [78] | 50.1 |
| Depth Anything V2 | Small | **53.9** |
| UperNet [84] | InternImage-B [78] | 50.8 |
| Depth Anything V2 | Base | **57.1** |
| UperNet [84] | InternImage-XL [78] | 55.0 |
| UperNet [84] | BEiT-L [4] | 56.3 |
| Mask2Former [13] | Swin-L [45] | 56.4 |
| ViT-Adapter [12] | BEiT-L [4] | 58.3 |
| OneFormer [29] | Swin-L [45] | 57.4 |
| OneFormer [29] | ConNeXt-XL [46] | 57.4 |
| Depth Anything V2 | Large | **58.6** |

(b) ADE20K dataset

Table 8: Transferring our Depth Anything V2 encoders to semantic segmentation. We adopt Mask2Former as our segmentation model. We achieve the results *without* Mapillary [1] or COCO [40] pre-training.

our models of various scales consistently achieve the best performance, outperforming other methods remarkably. These promising results indicate the potential of our model to serve as the initialization for diverse downstream semantic-related tasks.

## B.2 Transferring performance of each *labeled* dataset

We totally use five synthetic datasets to train our teacher model for pseudo labeling. Here we examine their individual effect on the model generalization capability. As demonstrated in Table 9, among them, the two purely indoor datasets Hypersim [58] and IRS [77] surprisingly fuel the most generalization ability. Although VKITTI 2 [9] has poor metric results, we find it is highly beneficial to the prediction sharpness, due to the large number of fine-grained structures (*e.g.*, leaves) in its training samples. Moreover, BlendedMVS [92] is critical to the capability of dealing with the bird's-eye view. Overall, each dataset has its own good properties to benefit the combined performance.

| Labeled Dataset | KITTI [24] | | NYU-D [70] | | Sintel [8] | | ETH3D [62] | | DIODE [76] | |
|---|---|---|---|---|---|---|---|---|---|---|
| | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ |
| BlendedMVS [92] | 0.088 | 0.919 | 0.069 | 0.957 | 0.538 | 0.661 | 0.150 | 0.839 | 0.095 | 0.915 |
| Hypersim [58] | <u>0.086</u> | <u>0.928</u> | <u>0.054</u> | 0.972 | 0.550 | 0.711 | **0.123** | **0.884** | 0.088 | <u>0.937</u> |
| IRS [77] | 0.100 | 0.900 | 0.055 | <u>0.973</u> | **0.435** | **0.738** | 0.149 | 0.831 | <u>0.084</u> | 0.931 |
| TartanAir [79] | 0.094 | 0.913 | 0.063 | 0.963 | 0.618 | 0.710 | 0.159 | 0.820 | 0.088 | 0.929 |
| VKITTI 2 [9] | 0.102 | 0.896 | 0.127 | 0.842 | 0.887 | 0.663 | 0.215 | 0.714 | 0.134 | 0.867 |
| All labeled data | **0.081** | **0.937** | **0.048** | **0.976** | <u>0.516</u> | <u>0.731</u> | <u>0.133</u> | <u>0.864</u> | **0.071** | **0.949** |

Table 9: Transferring performance of each *labeled* dataset with ViT-L. **Best results**, <u>second best results</u>.

## B.3 Transferring performance of each *unlabeled* dataset

We further analyze the benefit of each unlabeled source in Table 10. Accordingly, we present three observations. 1) Except the Sintel [8] synthetic game test set, unlabeled real images benefit all test sets tremendously. 2) When unlabeled images and test images share the same domain, the test results are improved most, *e.g.*, LSUN (indoor) improves the $\delta_1$ metric on NYU-D (indoor) from 0.928 $\rightarrow$ 0.970. 3) Even if unlabeled images and test images belong to contradictory domains, unlabeled images are still beneficial, *e.g.*, LSUN improves the $\delta_1$ on KITTI (street scene) from 0.889 $\rightarrow$ 0.913.

| Dataset | KITTI [24] | | NYU-D [70] | | Sintel [8] | | ETH3D [62] | | DIODE [76] | |
|---|---|---|---|---|---|---|---|---|---|---|
| | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ |
| Labeled datasets | 0.104 | 0.889 | 0.084 | 0.928 | <u>0.518</u> | 0.702 | 0.155 | 0.827 | 0.087 | 0.926 |
| + BDD100K | 0.091 | 0.916 | 0.071 | 0.951 | 0.600 | 0.708 | 0.153 | 0.834 | 0.087 | 0.927 |
| + Google Landmarks | 0.091 | 0.918 | 0.063 | 0.963 | 0.566 | 0.704 | 0.145 | 0.844 | <u>0.078</u> | <u>0.938</u> |
| + ImageNet-21K | <u>0.089</u> | <u>0.923</u> | 0.060 | 0.965 | 0.579 | 0.703 | 0.148 | 0.840 | 0.083 | 0.932 |
| + LSUN | 0.093 | 0.913 | <u>0.055</u> | <u>0.970</u> | 0.529 | 0.707 | 0.148 | 0.839 | 0.084 | 0.931 |
| + Objects365 | 0.089 | 0.920 | 0.058 | 0.967 | 0.551 | 0.701 | 0.145 | 0.846 | 0.080 | 0.937 |
| + Open Images V7 | 0.089 | 0.921 | 0.060 | 0.965 | 0.606 | 0.712 | 0.144 | 0.847 | 0.080 | 0.937 |
| + Places365 | 0.090 | 0.919 | 0.059 | 0.967 | 0.539 | 0.705 | 0.150 | 0.839 | 0.080 | 0.937 |
| + SA-1B | 0.092 | 0.915 | 0.067 | 0.956 | 0.652 | 0.708 | **0.142** | **0.850** | 0.080 | 0.935 |
| + All unlabeled data | **0.085** | **0.928** | **0.054** | **0.971** | **0.491** | **0.723** | <u>0.143</u> | <u>0.849</u> | **0.074** | **0.941** |

Table 10: Transferring performance by incorporating each *unlabeled* dataset with ViT-S. **Best**, <u>second best</u>.

## B.4 Are such large-scale unlabeled images really necessary?

We have proved that our used 62M unlabeled images are critical to model performance. However, we question that, is such a huge scale really necessary? What if we only use part of unlabeled sets and iterate the model for more epochs on it? To validate this, we solely use the SA-1B [33] dataset as our unlabeled source and train a model on it for the same iterations we use for 62M unlabeled images. As shown in Table 11, data diversity (*i.e.*, more datasets) is still highly important, which cannot be bridged by simply iterating a single dataset for more cycles. So we believe our large-scale unlabeled real images are necessary to ensure open-world generalization.

| Unlabeled Sets | # Images | Iterations | KITTI [24] | | NYU-D [70] | | Sintel [8] | | ETH3D [62] | | DIODE [76] | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ |
| SA-1B [33] | 11M | 480K | 0.090 | 0.915 | 0.073 | 0.948 | 0.588 | 0.707 | **0.141** | **0.852** | **0.073** | **0.942** |
| All eight sets | 62M | | **0.085** | **0.928** | **0.054** | **0.971** | **0.491** | **0.723** | 0.143 | 0.849 | 0.074 | 0.941 |

Table 11: Training the model solely on SA-1B for the same iterations as all sets (thus more cycles) with ViT-S.

## B.5 Performance on transparent or reflective surfaces

As aforementioned, one advantage of synthetic samples is the precise depth of the challenging transparent and reflective surfaces, which is important in navigation applications [82]. To validate the performance of our V2 in this specific domain, we compare different model predictions in the latest NTIRE 2024 Transparent Surface Challenge[3] [54]. Validation results are summarized in Table 12. Our V2 model achieves a remarkable boost over MiDaS [56] and Depth Anything V1 [89] in a zero-shot manner. Further, by simply fine-tuning our model with the challenge training data, we can nearly achieve the first-place score (0.912 *vs.* 0.917). Compared with the DINOv2 [50] encoder, our pre-trained model acts as a much stronger initialization (0.758 *vs.* 0.912).

| Method | Zero-shot (no fine-tuning) | | | Simple fine-tuning | | First place |
|---|---|---|---|---|---|---|
| | MiDaS V3.1 [7] | Depth Anything V1 [89] | V2 (Ours) | DINOv2 [50] | Depth Anything V2 (Ours) | |
| $\delta_1$ ($\uparrow$) | 0.259 | 0.535 | 0.836 | 0.758 | **0.912** | **0.917** |

Table 12: Results under different models and strategies in the NTIRE 2024 Transparent Surface Challenge [54].

---

[3] https://cvlab-unibo.github.io/booster-web/ntire24.html

## B.6 Comparison among various pre-trained encoders

We compare several currently most powerful pre-trained encoders in our MDE task, including BEiT [4], SAM [33], SynCLR [75], DINOv2 [50], and DINOv2 with registers [16]. As shown in Table 13, at the ViT-large scale, we find DINOv2 serial [50, 16] is remarkably superior to all other encoders. The success of DINOv2 further reflects the promising future of the data-driven roadmap, since it carefully collects 142M pre-training data without designing fancy algorithms or architectures.

When scaling up the ViT-large encoder to ViT-giant, we surprisingly observe DINOv2-G Reg [16] is much inferior to the non-register initial version [50]. This is the same as the findings in Probe3D [3]. Thus, we choose to build our teacher and student models on the original DINOv2 encoders.

| Encoder | KITTI [24] | | NYU-D [70] | | Sintel [8] | | ETH3D [62] | | DIODE [76] | |
|---|---|---|---|---|---|---|---|---|---|---|
| | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ | AbsRel | $\delta_1$ |
| BEiT-L [4] | 0.149 | 0.814 | 0.068 | 0.950 | 0.777 | 0.627 | 0.145 | 0.846 | 0.103 | 0.912 |
| SAM-L [33] | 0.104 | 0.893 | 0.186 | 0.745 | 0.703 | 0.688 | 0.143 | 0.849 | 0.108 | 0.907 |
| SynCLR-L [75] | 0.278 | 0.650 | 0.344 | 0.469 | 1.608 | 0.493 | 0.301 | 0.638 | 0.262 | 0.712 |
| DINOv2-L [50] | 0.081 | 0.937 | **0.048** | **0.976** | **0.516** | 0.731 | **0.133** | **0.864** | 0.071 | 0.949 |
| DINOv2-L Reg [16] | **0.078** | **0.942** | 0.049 | 0.975 | 0.522 | **0.734** | 0.138 | 0.856 | **0.068** | **0.952** |
| DINOv2-G [50] | **0.075** | **0.947** | **0.044** | **0.979** | **0.530** | **0.767** | **0.131** | **0.865** | **0.066** | **0.954** |
| DINOv2-G Reg [16] | 0.084 | 0.926 | 0.061 | 0.964 | 0.753 | 0.729 | 0.141 | 0.852 | 0.086 | 0.931 |

Table 13: Comparison among various pre-trained encoders when purely trained on synthetic images.

## B.7 Benefit of gradient matching loss to fine-grained predictions

MiDaS [56] proposes a gradient matching loss $\mathcal{L}_{gm}$ to enhance the depth sharpness. Unfortunately, we find this loss term fails to bring evident improvement when the model is trained on labeled real datasets. We speculate that, the *sparse* and *coarse* groundtruth label in real datasets cannot provide fine-grained supervision, even with this explicit regularization. To check this, we further apply and ablate this loss term on synthetic training datasets, whose labels are complete and highly precise. We gradually increase the loss weight of $\mathcal{L}_{gm}$ and observe the corresponding depth sharpness. As shown in Figure 10, when the weight is increased from the default 0.5 to 4.0, the sharpness is steadily improved. We finally set the weight as 2.0 to trade off between the metric results and sharpness.
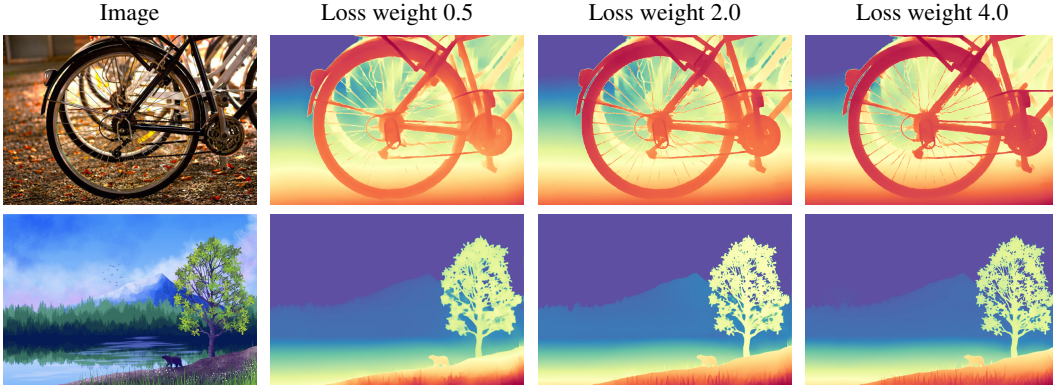


Figure 10: Effect of the gradient matching loss $\mathcal{L}_{gm}$ in terms of fine-grained details.

## B.8 Test-time resolution scaling up

By default, we test images at the same resolution as that used in training, *i.e.*, resizing the shorter size to 518 with the aspect ratio kept. This is a common practice to achieve the optimal performance. However, we surprisingly find that our model has the property of "test-time resolution scaling up". It means we can almost freely increase the image resolution at test time to produce more fine-grained depth maps. As shown in Figure 11, when gradually adjusting the resolution by $2\times$ and $4\times$ of the base resolution (518), the depth sharpness is also gradually improved.
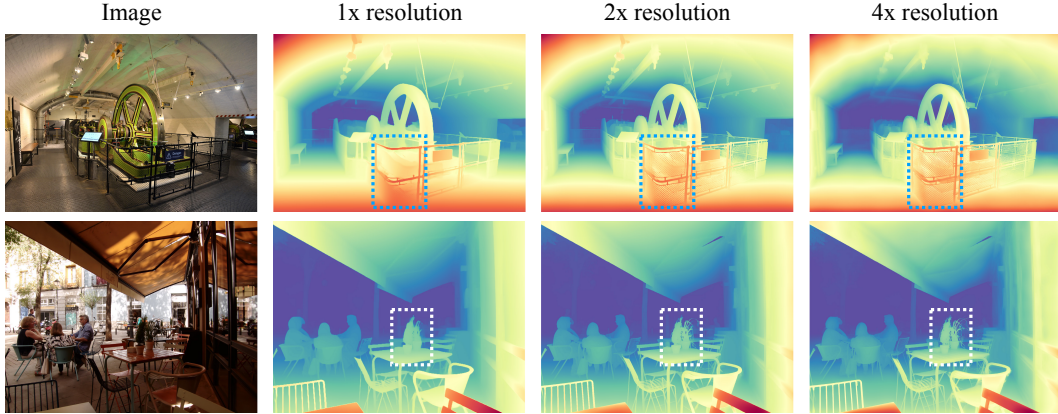
Figure 11: Test-time resolution scaling up can further improve the prediction sharpness.

### B.9 Harm of real labeled images to fine-grained predictions

According to the ablation study in Depth Anything V1 [89], HRWSI [83] is the best-performed real training dataset. We attempt to add it to our synthetic training sets. However, as shown in Figure 12, we find although it only accounts for 5% of the total training images, its coarse depth labels have a huge negative impact on the original fine-grained predictions. So we choose to use purely synthetic images to train our largest teacher model to ensure the supervision preciseness.
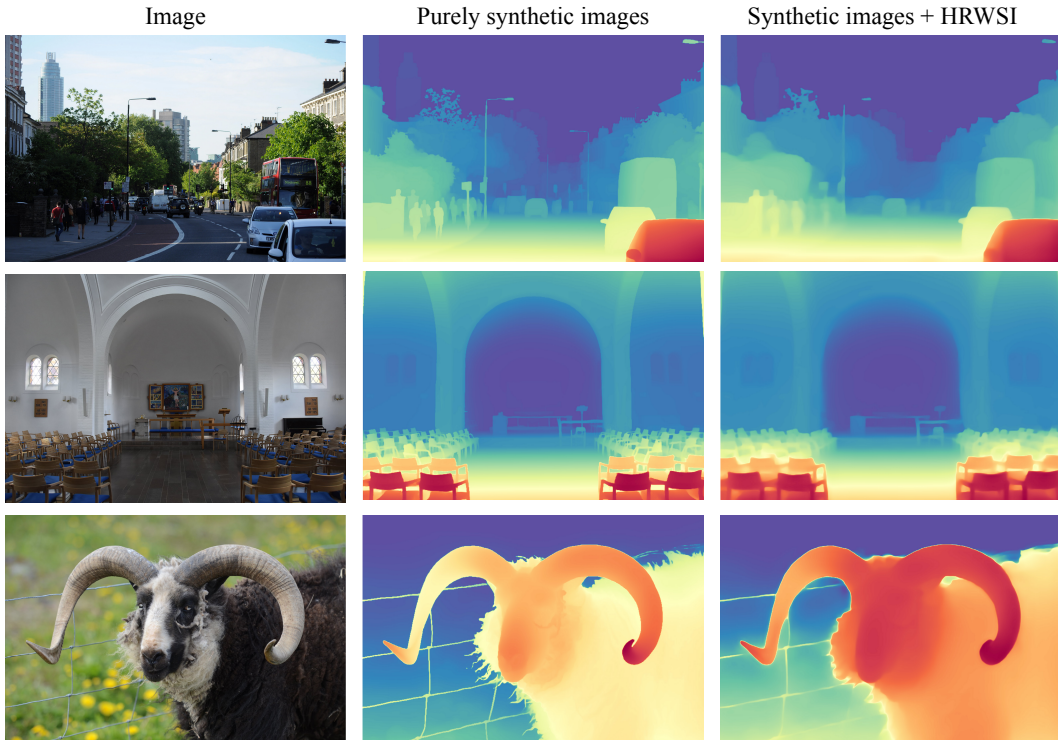


Figure 12: Adding real training dataset, *e.g.*, HRWSI, to synthetic training datasets, will ruin the original fine-grained depth predictions.

### B.10 Qualitative comparison between Depth Anything V1 and V2

Please refer to Figure 13. Our Depth Anything V2 produces much more fine-grained depth predictions than V1 [89]. Ours are also highly robust to transparent objects.

### B.11 Qualitative comparison between Marigold and Depth Anything V2

Please refer to Figure 14. Our Depth Anything V2 is significantly more robust than Marigold [31].

### B.12 Qualitative comparison between our metric depth models and ZoeDepth

We fine-tune our finally released metric depth models purely on synthetic datasets, such as Hypersim [58] and Virtual KITTI [9]. In Figure 15, we compare our metric depth predictions with ZoeDepth, which is trained on real datasets like NYUv2 [70].

### B.13 Qualitative comparison between w/ and w/o pseudo-labeled real images

Please refer to Figure 16. As shown, purely trained on precise synthetic images, the DINOv2-small-based model suffers severe generalization problem. However, when trained on the high-quality and diverse pseudo-labeled real images, even the small model (25M parameters) exhibits powerful generalization capability to complex scenes.

### B.14 Qualitative results of produced pseudo labels

Please refer to Figure 17. Our teacher produces highly precise pseudo labels on diverse real images.

### B.15 Qualitative results on test benchmarks

Please refer to Figure 18. Our model is consistently better than V1 [89] on standard benchmarks.

## C  DA-2K Evaluation Benchmark

### C.1 Per-scenario accuracy

We report the per-scenario accuracy on our DA-2K evaluation benchmark. By comparing the results of training on labeled synthetic images ($\mathcal{D}^l$) and pseudo-labeled real images ($\mathcal{D}^u$), we can clearly see the value of large-scale unlabeled data and also the preciseness of our pseudo labels.

| Encoder | $\mathcal{D}^l$ | $\mathcal{L}^u$ | Indoor | Outdoor | Non-real | Transparent | Adverse style | Aerial | Underwater | Object | **Mean** |
|---------|-----------------|-----------------|--------|---------|----------|-------------|---------------|--------|------------|--------|----------|
| ViT-S   | ✓ |   | 88.1 | 87.8 | 90.8 | 86.9 | 90.6 | 93.8 | 94.9 | 89.9 | 89.8 |
|         |   | ✓ | 92.9 | 93.0 | 98.4 | 94.4 | 95.7 | 96.4 | 99.2 | 96.6 | **95.3** |
| ViT-B   | ✓ |   | 91.2 | 91.9 | 95.7 | 90.2 | 90.9 | 96.4 | 94.9 | 96.6 | 92.9 |
|         |   | ✓ | 96.2 | 94.8 | 98.7 | 96.3 | 96.7 | 99.0 | 100 | 97.3 | **97.0** |
| ViT-L   | ✓ |   | 94.5 | 93.9 | 98.4 | 93.9 | 96.3 | 97.4 | 99.2 | 98.0 | 96.0 |
|         |   | ✓ | 96.4 | 93.9 | 99.0 | 96.3 | 97.3 | 99.5 | 99.2 | 98.0 | **97.1** |

Table 14: Per-scenario accuracy (%) of Depth Anything V2 on our proposed benchmark DA-2K.

### C.2 Comparison with the DIW dataset

Although DIW [11] and our DA-2K use the same annotation format (both sparse pixel pairs, we are inspired by DIW), our proposed DA-2K dataset is better in four aspects:

- **(more precise)** DIW is very noisy. For *most* pairs in DIW, we cannot decide the relative depth or hold the opposite opinion as the provided label. This can also be supported by MiDaS [7] that, better and larger models instead perform worse on DIW. In comparison, our DA-2K is precise, because we exclude many hard-to-decide or controversial pairs.

- **(better organized)** DIW randomly downloads images from Flickr, without carefully organizing. This would make users struggle to obtain straightforward insights from the evaluation results. In comparison, our DA-2K organizes all images by application scenarios, and thus can provide results for each individual application scenario.

16

- **(more diverse)** DIW images are typically collected from real life. However, considering the widespread application of MDE models in AIGC [101, 39], we provide additional non-real images, such as AI-generated images, cartoon images, *etc.*.
- **(high-resolution)** Most images in DIW have a low resolution of around $300 \times 500$, while we provide mostly $1500 \times 2000$ high-resolution images.

## C.3  Annotation details

To alleviate the burden of human annotators and avoid hard-to-decide pairs, we only pop out pixel pairs whose predicted depth ratio is larger than 3. For the evaluation scenarios of "transparent" and "object", we do not rely on model disagreement to pop out pairs. We simply manually analyze the images and select challenging pairs suited to the scenario. For other scenarios, we adopt both selection pipelines (*i.e.*, automatic disagreement-based selection and manual selection). In Table 15, we list the keywords we use to download images for each evaluation scenario.

| Evaluation scenario | Keywords |
| --- | --- |
| Indoor | room, home, living room, kitchen, bedroom, office, store, library, restaurant, museum, hall |
| Outdoor | road, outdoor, street, urban, rural, park, beach, mountain, downtown, alley, skyscraper, traffic, bridge, construction, parade, fireworks, festival, sporting event |
| Non-real (*e.g.*, AIGC, painting, *etc.*) | AI-generated, computer-generated, artwork, oil painting, impressionism, realism, abstract art, cartoon, animation, comic, caricature, illustration, fantasy, sci-fi, cyberpunk, alien, mythology |
| Transparent / reflective surfaces | glass, window, crystal, ice, water, transparent, clear, acrylic, plastic, reflective, mirror, see-through |
| Adverse style (*e.g.*, foggy, dark, *etc.*) | fog, dark, night, mid-night, overexposed, blur, snow, rain |
| Aerial | aerial, landscape, drone view, bird's eye view, city, cityscape, satellite view, top-down view |
| Underwater | underwater, ocean, sea, coral reef, diving, submarine, aquarium, marine life, shipwreck |
| Object | car, bicycle, motorcycle, airplane, bus, train, truck, boat, traffic light, fire hydrant, stop sign, parking meter, bench, bird, cat, dog, horse, sheep, cow, elephant, bear, zebra, giraffe, backpack, umbrella, sports ball, kite, baseball bat, cup, fork, knife, spoon, bowl, banana, apple, chair, bed, dining table, microwave, oven, toaster, sink, refrigerator, vase, scissors, teddy bear |

Table 15: Eight evaluation scenarios encompassed in our DA-2K. We use the keywords generated by GPT-4 to download images of corresponding scenarios on Flickr.

## C.4  Visualization

In Figure 19, we visualize some samples in our proposed DA-2K benchmark. They cover diverse representative scenarios and are of precise sparse annotations.

# D  Limitations

Currently, we use 62M unlabeled images for training. The computational burden is very heavy. Thus, in the future, we will study how to leverage such large-scale visual data more efficiently. Moreover, the current synthetic training sets are not diverse enough. We will attempt to collect synthetic images from more sources to train a more capable teacher model for better pseudo labeling.
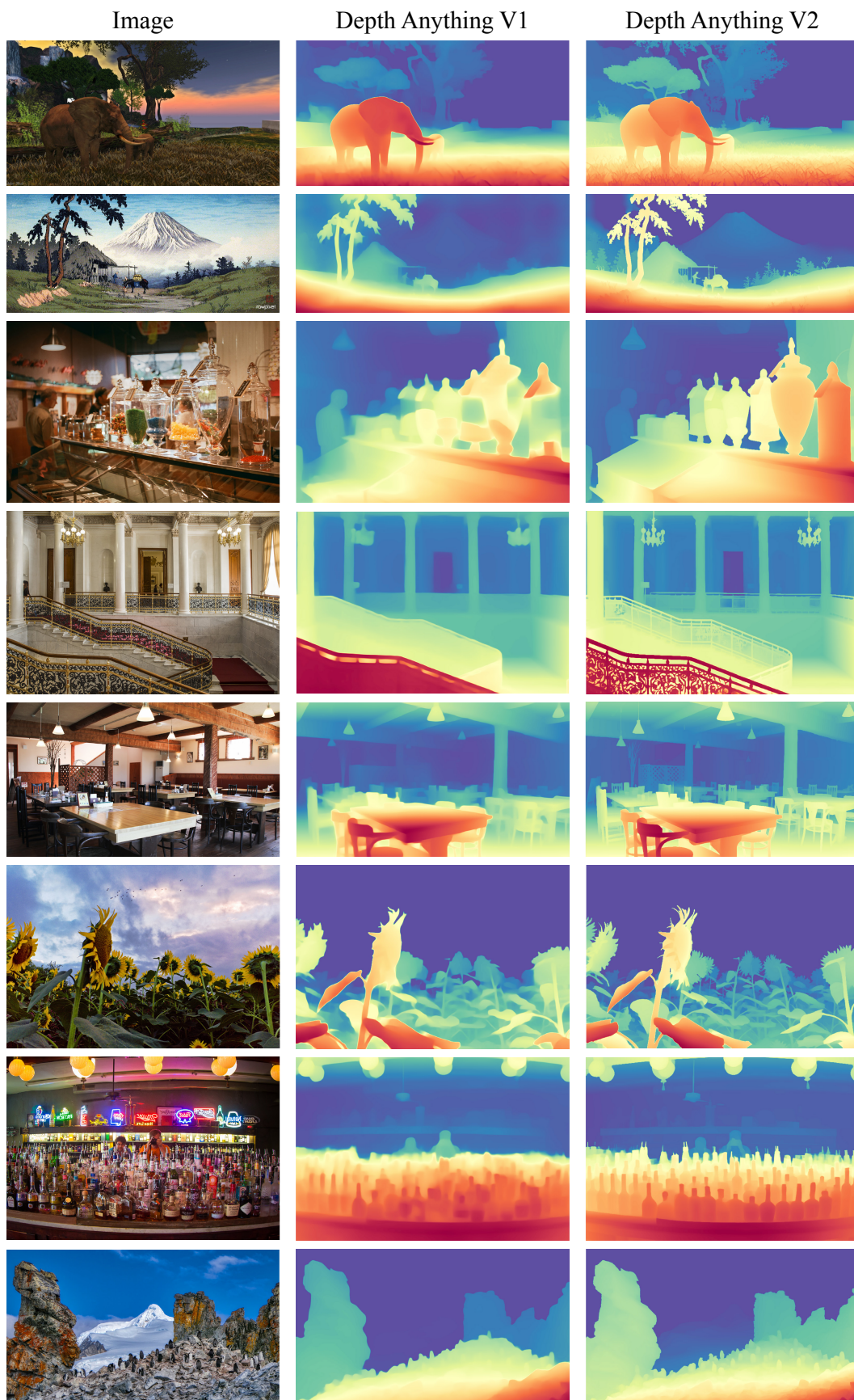
Figure 13: Comparison between Depth Anything V1 [89] and our V2 in open-world images.
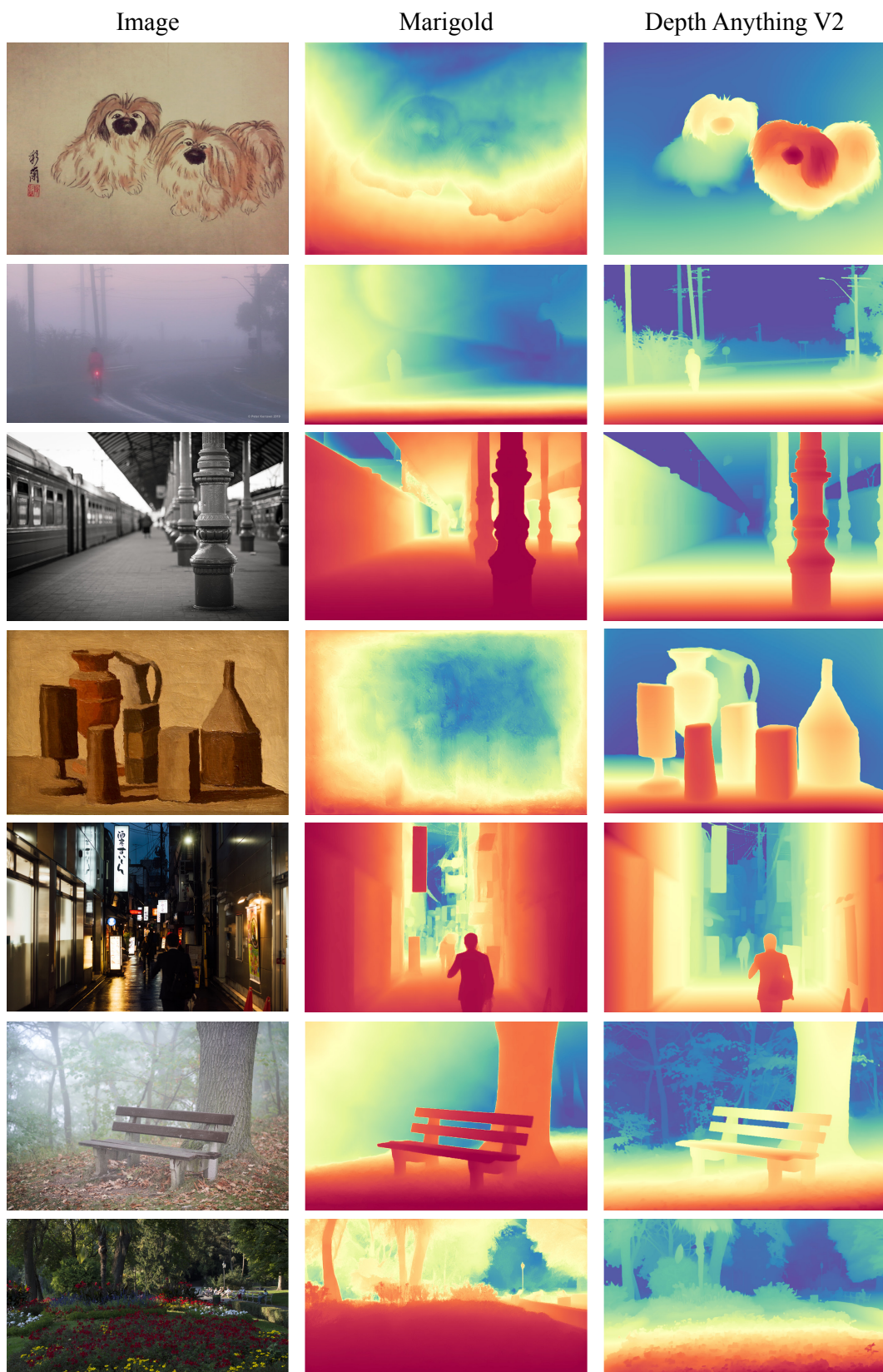
|  | Image | Marigold | Depth Anything V2 |

Figure 14: Comparison between Marigold [31] and our V2 in open-world images.

Figure 15: Comparison between ZoeDepth [6] and our fine-tuned metric depth model.

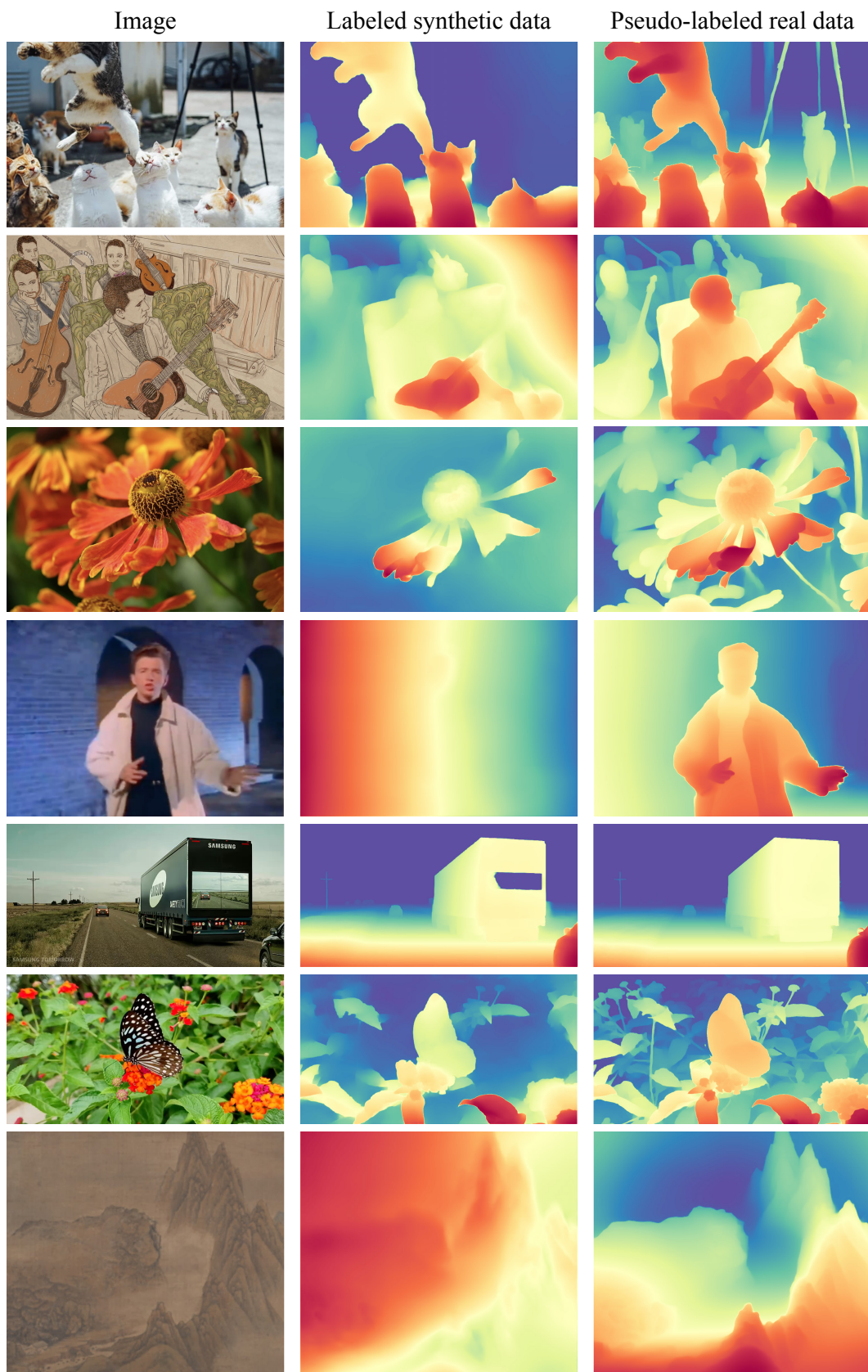| Image | Labeled synthetic data | Pseudo-labeled real data |

Figure 16: Qualitative comparison of the DINOv2-small-based depth model trained solely on labeled synthetic images and solely pseudo-labeled real images. The robustness is tremendously enhanced.

| Unlabeled image | Pseudo label | Unlabeled image | Pseudo label |

Figure 17: Visualization of our produced pseudo depth labels. From top to bottom, the highly diverse images are sampled from BDD100K [97], Google Landmarks [81], ImageNet-21K [60], LSUN [98], Objects365 [65], Open Images V7 [35], Places365 [103], and SA-1B [33] datasets, respectively.

| Image | Depth Anything V1 | Depth Anything V2 |
|---|---|---|



Figure 18: Qualitative results on widely adopted test benchmarks, *e.g.*, KITTI, NYU, and DIODE.

Figure 19: Visualization of images and precise sparse annotations on our benchmark DA-2K. Please **zoom in** to better view the annotated pairs. The green point is annotated as closer than the red point. From top to bottom, the images are sampled from indoor, outdoor, non-real, transparent/reflective, adverse style, aerial, underwater, and object scenarios, respectively.

# References

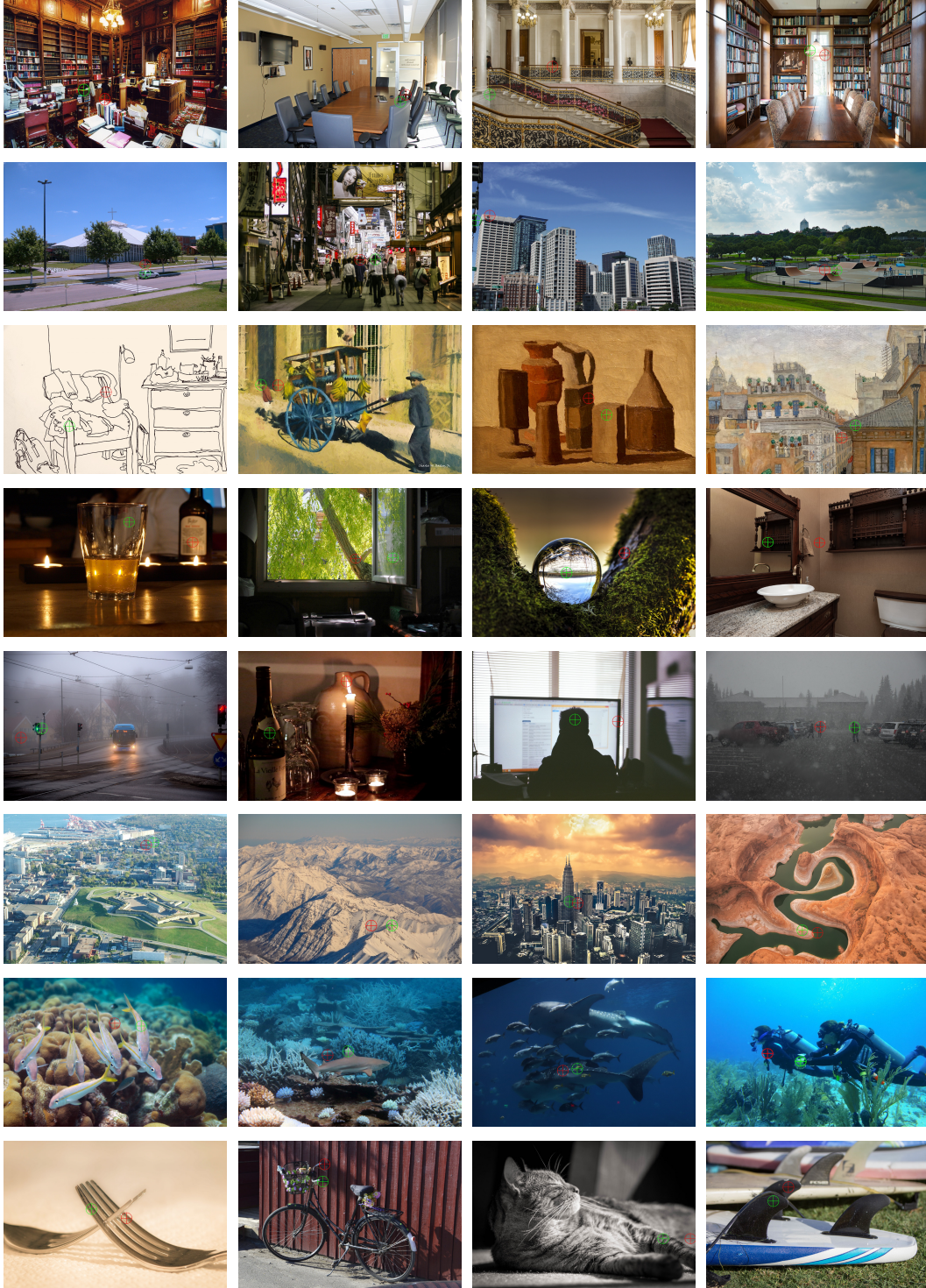[1] Manuel López Antequera, Pau Gargallo, Markus Hofinger, Samuel Rota Bulò, Yubin Kuang, and Peter Kontschieder. Mapillary planet-scale depth dataset. In *ECCV*, 2020. 12

[2] Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *NeurIPS*, 2014. 10

[3] Mohamed El Banani, Amit Raj, Kevis-Kokitsi Maninis, Abhishek Kar, Yuanzhen Li, Michael Rubinstein, Deqing Sun, Leonidas Guibas, Justin Johnson, and Varun Jampani. Probing the 3d awareness of visual foundation models. In *CVPR*, 2024. 8, 14

[4] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. In *ICLR*, 2022. 2, 5, 12, 14

[5] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *CVPR*, 2021. 9, 10

[6] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv:2302.12288*, 2023. 2, 9, 20

[7] Reiner Birkl, Diana Wofk, and Matthias Müller. Midas v3. 1–a model zoo for robust monocular relative depth estimation. *arXiv:2307.14460*, 2023. 2, 3, 5, 8, 9, 10, 13, 16

[8] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source movie for optical flow evaluation. In *ECCV*, 2012. 8, 9, 10, 12, 13, 14

[9] Yohann Cabon, Naila Murray, and Martin Humenberger. Virtual kitti 2. *arXiv:2001.10773*, 2020. 4, 9, 12, 16

[10] Shengcao Cao, Mengtian Li, James Hays, Deva Ramanan, Yu-Xiong Wang, and Liangyan Gui. Learning lightweight object detectors via multi-teacher progressive distillation. In *ICML*, 2023. 10

[11] Weifeng Chen, Zhao Fu, Dawei Yang, and Jia Deng. Single-image depth perception in the wild. In *NeurIPS*, 2016. 7, 8, 16

[12] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. In *ICLR*, 2023. 12

[13] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022. 12

[14] Jaehoon Cho, Dongbo Min, Youngjung Kim, and Kwanghoon Sohn. Diml/cvl rgb-d dataset: 2m rgb-d images of natural indoor and outdoor scenes. *arXiv:2110.11590*, 2021. 4, 10

[15] Alex Costanzino, Pierluigi Zama Ramirez, Matteo Poggi, Fabio Tosi, Stefano Mattoccia, and Luigi Di Stefano. Learning depth estimation for transparent and mirror surfaces. In *ICCV*, 2023. 2

[16] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. In *ICLR*, 2024. 14

[17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 12

[18] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NeurIPS*, 2014. 10

[19] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *CVPR*, 2018. 10

[20] Xiao Fu, Wei Yin, Mu Hu, Kaixuan Wang, Yuexin Ma, Ping Tan, Shaojie Shen, Dahua Lin, and Xiaoxiao Long. Geowizard: Unleashing the diffusion priors for 3d geometry estimation from a single image. *arXiv:2403.12013*, 2024. 2, 4, 7, 8, 9, 10

[21] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. *arXiv:2404.12390*, 2024. 8

[22] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, 2015. 10

[23] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 2016. 10

[24] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 2013. 7, 8, 9, 10, 12, 13, 14

[25] Ming Gui, Johannes S Fischer, Ulrich Prestel, Pingchuan Ma, Dmytro Kotovenko, Olga Grebenkova, Stefan Andreas Baumann, Vincent Tao Hu, and Björn Ommer. Depthfm: Fast monocular depth estimation with flow matching. *arXiv:2403.13788*, 2024. 1, 2, 4, 8, 10

[26] Vitor Guizilini, Igor Vasiljevic, Dian Chen, Rareș Ambruș, and Adrien Gaidon. Towards zero-shot scale-aware monocular depth estimation. In *ICCV*, 2023. 2, 3

[27] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv:1503.02531*, 2015. 6, 10

[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:2404.15506*, 2024. 2, 3, 5, 8

[29] Jitesh Jain, Jiachen Li, Mang Tik Chiu, Ali Hassani, Nikita Orlov, and Humphrey Shi. One-former: One transformer to rule universal image segmentation. In *CVPR*, 2023. 12

[30] Yuanfeng Ji, Zhe Chen, Enze Xie, Lanqing Hong, Xihui Liu, Zhaoqiang Liu, Tong Lu, Zhenguo Li, and Ping Luo. Ddp: Diffusion model for dense visual prediction. In *ICCV*, 2023. 12

[31] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing diffusion-based image generators for monocular depth estimation. In *CVPR*, 2024. 1, 2, 4, 7, 8, 9, 10, 16, 19

[32] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *TOG*, 2023. 2

[33] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *ICCV*, 2023. 5, 7, 9, 12, 13, 14, 22

[34] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 10

[35] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*, 2020. 12, 22

[36] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *ICMLW*, 2013. 10

[37] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *CVPR*, 2018. 3, 4

[38] Zhenyu Li, Shariq Farooq Bhat, and Peter Wonka. Patchfusion: An end-to-end tile-based framework for high-resolution monocular metric depth estimation. In *CVPR*, 2024. 2

[39] Jun Hao Liew, Hanshu Yan, Jianfeng Zhang, Zhongcong Xu, and Jiashi Feng. Magicedit: High-fidelity and temporally coherent video editing. *arXiv:2308.14749*, 2023. 2, 17

[40] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 12

[41] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 8

[42] Lingjie Liu, Nenglun Chen, Duygu Ceylan, Christian Theobalt, Wenping Wang, and Niloy J Mitra. Curvefusion: reconstructing thin structures from rgbd sequences. *TOG*, 2018. 2

[43] Yifan Liu, Ke Chen, Chris Liu, Zengchang Qin, Zhenbo Luo, and Jingdong Wang. Structured knowledge distillation for semantic segmentation. In *CVPR*, 2019. 10

[44] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *CVPR*, 2022. 9

[45] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 12

[46] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, 2022. 12

[47] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2

[48] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant. In *AAAI*, 2020. 6, 10

[49] Jia Ning, Chen Li, Zheng Zhang, Chunyu Wang, Zigang Geng, Qi Dai, Kun He, and Han Hu. All in tokens: Unifying output space of visual tasks via soft token. In *ICCV*, 2023. 9

[50] 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. *TMLR*, 2023. 2, 5, 13, 14

[51] Vaishakh Patil, Christos Sakaridis, Alexander Liniger, and Luc Van Gool. P3depth: Monocular depth estimation with a piecewise planarity prior. In *CVPR*, 2022. 9

[52] 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. 2

[53] Weichao Qiu, Fangwei Zhong, Yi Zhang, Siyuan Qiao, Zihao Xiao, Tae Soo Kim, and Yizhou Wang. Unrealcv: Virtual worlds for computer vision. In *ACM MM*, 2017. 4

[54] Pierluigi Zama Ramirez, Fabio Tosi, Matteo Poggi, Samuele Salti, Stefano Mattoccia, and Luigi Di Stefano. Open challenges in deep stereo: the booster dataset. In *CVPR*, 2022. 4, 13

[55] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *ICCV*, 2021. 8, 9, 10

[56] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *TPAMI*, 2022. 2, 3, 6, 10, 13, 14

[57] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *ECCV*, 2016. 4

[58] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *ICCV*, 2021. 4, 9, 12, 16

[59] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 10

[60] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 12, 22

[61] Swami Sankaranarayanan, Yogesh Balaji, Arpit Jain, Ser Nam Lim, and Rama Chellappa. Learning from synthetic data: Addressing domain shift for semantic segmentation. In *CVPR*, 2018. 10

[62] Thomas Schops, Johannes L Schonberger, 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 *CVPR*, 2017. 8, 9, 10, 12, 13, 14

[63] Shital Shah, Debadeepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In *Field and Service Robotics*, 2017. 4

[64] Mohamad Shahbazi, Liesbeth Claessens, Michael Niemeyer, Edo Collins, Alessio Tonioni, Luc Van Gool, and Federico Tombari. Inserf: Text-driven generative object insertion in neural 3d scenes. *arXiv:2401.05335*, 2024. 2

[65] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *ICCV*, 2019. 12, 22

[66] Shuwei Shao, Zhongcai Pei, Weihai Chen, Xingming Wu, and Zhengguo Li. Nddepth: Normal-distance assisted monocular depth estimation. In *ICCV*, 2023. 9

[67] Shuwei Shao, Zhongcai Pei, Xingming Wu, Zhong Liu, Weihai Chen, and Zhengguo Li. Iebins: Iterative elastic bins for monocular depth estimation. In *NeurIPS*, 2023. 9

[68] Jaidev Shriram, Alex Trevithick, Lingjie Liu, and Ravi Ramamoorthi. Realmdreamer: Text-driven 3d scene generation with inpainting and depth diffusion. *arXiv:2404.07199*, 2024. 2

[69] Changyong Shu, Yifan Liu, Jianfei Gao, Zheng Yan, and Chunhua Shen. Channel-wise knowledge distillation for dense prediction. In *ICCV*, 2021. 10

[70] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, 2012. 3, 7, 8, 9, 10, 12, 13, 14, 16

[71] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. In *NeurIPS*, 2020. 10

[72] Jaime Spencer, Fabio Tosi, Matteo Poggi, Ripudaman Singh Arora, Chris Russell, Simon Hadfield, Richard Bowden, GuangYuan Zhou, ZhengXin Li, Qiang Rao, et al. The third monocular depth estimation challenge. *arXiv:2404.16831*, 2024. 2, 6

[73] Rupesh K Srivastava, Klaus Greff, and Jürgen Schmidhuber. Training very deep networks. In *NeurIPS*, 2015. 10

[74] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *ICCV*, 2021. 12

[75] Yonglong Tian, Lijie Fan, Kaifeng Chen, Dina Katabi, Dilip Krishnan, and Phillip Isola. Learning vision from models rivals learning vision from data. In *CVPR*, 2024. 5, 14

[76] Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z Dai, Andrea F Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R Walter, et al. Diode: A dense indoor and outdoor depth dataset. *arXiv:1908.00463*, 2019. 8, 9, 10, 12, 13, 14

[77] Qiang Wang, Shizhen Zheng, Qingsong Yan, Fei Deng, Kaiyong Zhao, and Xiaowen Chu. Irs: A large naturalistic indoor robotics stereo dataset to train deep models for disparity and surface normal estimation. In *ICME*, 2021. 12

[78] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions. In *CVPR*, 2023. 12

[79] 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 *IROS*, 2020. 12

[80] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *CVPR*, 2019. 2

[81] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval. In *CVPR*, 2020. 12, 22

[82] Diana Wofk, Fangchang Ma, Tien-Ju Yang, Sertac Karaman, and Vivienne Sze. Fastdepth: Fast monocular depth estimation on embedded systems. In *ICRA*, 2019. 2, 13

[83] Ke Xian, Jianming Zhang, Oliver Wang, Long Mai, Zhe Lin, and Zhiguo Cao. Structure-guided ranking loss for single image depth prediction. In *CVPR*, 2020. 3, 4, 15

[84] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, 2018. 12

[85] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *NeurIPS*, 2021. 12

[86] Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation for consistency training. *NeurIPS*, 2020. 10

[87] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Yi Wang, and Zhangyang Wang. Neurallift-360: Lifting an in-the-wild 2d photo to a 3d object with 360deg views. In *CVPR*, 2023. 2

[88] Guangkai Xu, Yongtao Ge, Mingyu Liu, Chengxiang Fan, Kangyang Xie, Zhiyue Zhao, Hao Chen, and Chunhua Shen. Diffusion models trained with large data are transferable visual models. *arXiv:2403.06090*, 2024. 2

[89] 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. 1, 2, 3, 5, 6, 7, 8, 9, 10, 11, 13, 15, 16, 18

[90] Lihe Yang, Zhen Zhao, Lei Qi, Yu Qiao, Yinghuan Shi, and Hengshuang Zhao. Shrinking class space for enhanced certainty in semi-supervised learning. In *ICCV*, 2023. 10

[91] Xiaodong Yang, Zhuang Ma, Zhiyu Ji, and Zhe Ren. Gedepth: Ground embedding for monocular depth estimation. In *ICCV*, 2023. 9

[92] 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 *CVPR*, 2020. 12

[93] Chongjie Ye, Yinyu Nie, Jiahao Chang, Yuantao Chen, Yihao Zhi, and Xiaoguang Han. Gaustudio: A modular framework for 3d gaussian splatting and beyond. *arXiv:2403.19632*, 2024. 2

[94] Wei Yin, Yifan Liu, Chunhua Shen, and Youliang Yan. Enforcing geometric constraints of virtual normal for depth prediction. In *ICCV*, 2019. 10

[95] 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. 2, 3, 10

[96] Wei Yin, Jianming Zhang, Oliver Wang, Simon Niklaus, Long Mai, Simon Chen, and Chunhua Shen. Learning to recover 3d scene shape from a single image. In *CVPR*, 2021. 10

[97] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *CVPR*, 2020. 12, 22

[98] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv:1506.03365*, 2015. 12, 22

[99] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. New crfs: Neural window fully-connected crfs for monocular depth estimation. *arXiv:2203.01502*, 2022. 9

[100] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR*, 2017. 10

[101] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *ICCV*, 2023. 2, 17

[102] Wenliang Zhao, Yongming Rao, Zuyan Liu, Benlin Liu, Jie Zhou, and Jiwen Lu. Unleashing text-to-image diffusion models for visual perception. In *ICCV*, 2023. 9

[103] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *TPAMI*, 2017. 12, 22