



Depth Anything V2 - Monocular Depth Estimation

J.Pranitha^a, K.Ram Narayan Reddy^b, Mr.R. Srinivas^c, Ms.D.Deepika^d

^{a,b} UG Student, Department of Emerging Technologies, Mahatma Gandhi Institute of Technology (Autonomous), Hyderabad 500075

^{c,d} Assistant Professor, Department of Emerging Technologies, Mahatma Gandhi Institute of Technology (Autonomous), Hyderabad 500075

ABSTRACT :

This work introduces a refined framework for monocular depth estimation that prioritizes the use of real-world images without converting them into synthetic representations. Unlike approaches that depend heavily on synthetic data to overcome label noise, our method capitalizes on the authenticity and diversity found naturally in real-world datasets. By applying high-quality pseudo-labeling techniques to real images, we eliminate the need for synthetic generation while still achieving fine-grained, accurate depth estimations. This ensures that our models are trained on realistic textures, complex layouts, transparent and reflective surfaces, and diverse environmental conditions, enhancing generalization and real-world deployment capability.

Building upon a scalable training pipeline, we leverage a strong teacher model to generate precise pseudo-depth labels, which are then used to train student models of varying sizes, ranging from lightweight (25M parameters) to large-scale (1.3B parameters) versions. Our models achieve significant improvements in both efficiency and accuracy compared to prior methods, including those built with Stable Diffusion backbones. The system demonstrates robust handling of complex depth cues such as occlusions, transparency, and thin object structures, without the pitfalls associated with synthetic-to-real domain gaps. Furthermore, our models are designed to be flexible for downstream tasks, including fine-tuning for metric depth estimation across a range of application scenarios like autonomous navigation, AR/VR content, and 3D scene reconstruction.

To address the shortcomings in existing evaluation datasets—which often suffer from noise, low diversity, and low resolution—we introduce a new high-resolution benchmark, DA-2K. This benchmark is curated specifically to include diverse, challenging real-world scenes and provides precise, human-verified depth annotations. Through this work, we establish a powerful, scalable, and practical foundation for monocular depth estimation based solely on real-world imagery, setting a new direction for future developments that emphasize authenticity, precision, and broad applicability without the reliance on synthetic data.

Keywords: Monocular depth estimation, depth prediction, computer vision, transformer models, multi-task learning, dense prediction, scene understanding, self-supervised learning, vision transformer, generalist vision models, depth estimation benchmark, multi-dataset training, transfer learning, zero-shot generalization, Depth Anything.

Introduction

Monocular depth estimation (MDE) is a fundamental task in computer vision with extensive applications, including 3D reconstruction, navigation, autonomous driving, virtual reality, and robotics. Traditional stereo-based methods, although accurate, require complex setups with multiple calibrated cameras or expensive depth sensors, limiting their scalability and deployment in practical scenarios. Monocular approaches, by contrast, are lightweight and cost-effective, relying on a single RGB image to predict depth. This advantage has driven significant research interest in the field. Early progress in monocular depth estimation was fueled by convolutional neural network (CNN)-based models, which introduced learning-based techniques that outperformed classical methods. However, CNNs have limitations in capturing long-range dependencies due to their local receptive fields. Recent advances with Vision Transformers (ViTs), such as DINOv2, have significantly boosted the performance of monocular depth estimation by enabling global context modeling across the entire image. ViTs leverage self-attention mechanisms to aggregate spatial information more effectively, leading to richer and more precise depth predictions from a single image. The evolution from traditional stereo vision to CNN-based models, and now to transformer-based models like ViT, marks a key turning point in the technological landscape of dense depth prediction. Building upon these technological developments, Depth Anything V2 introduces a series of significant innovations that advance the state of monocular depth estimation. One of its major contributions is the strategic use of synthetic data during the training process. By replacing noisy, real-world labeled images with precisely annotated synthetic datasets, the model eliminates common issues such as label inaccuracies, incomplete ground truths, and coverage gaps. Depth Anything V2 adopts a teacher-student training framework wherein a powerful teacher model, trained exclusively on synthetic data, is used to generate pseudo-labels for real-world unlabeled images. These pseudo-labeled real images are then utilized to train more efficient and smaller student models. This training strategy not only bridges the domain gap between synthetic and real-world data but also enables the model to scale flexibly across different computational budgets, offering variants with parameter sizes ranging from 25 million to 1.3 billion. Furthermore, Depth Anything V2 is uniquely designed to generalize effectively across complex and challenging environments, such as scenes with transparent objects, reflective surfaces, and intricate spatial layouts, which typically pose difficulties for traditional models. The large-scale training on pseudo-labeled real images enhances robustness and ensures that the model performs well even under diverse and unpredictable real-world conditions.

To complement its model innovations, Depth Anything V2 introduces DA-2K, a new evaluation benchmark specifically designed to test depth estimation models against real-world complexities. Traditional benchmarks often suffer from noisy annotations, limited scene diversity, and lower image resolutions, which can misrepresent a model's true capabilities. DA-2K addresses these shortcomings by offering precise depth annotations, high-resolution imagery, and a wide variety of challenging indoor and outdoor scenarios. This benchmark sets a new standard for model evaluation by ensuring that assessments reflect practical deployment challenges more accurately. In head-to-head comparisons, including against architectures based on Stable Diffusion, Depth Anything V2 demonstrates superior performance in inference speed, depth prediction accuracy, and the ability to handle complex scene structures efficiently. These advancements position Depth Anything V2 as a cutting-edge, scalable, and robust solution for monocular depth estimation, paving the way for more precise and reliable 3D scene understanding across various industries and applications.

Nomenclature

Depth Anything V2.
 Monocular depth estimation.
 Synthetic data.
 Label inaccuracies.
 Incomplete ground truths.
 Teacher-student framework.
 Pseudo-labels.
 Domain gap.
 Parameter sizes.
 Generalization.
 DA-2K.
 Stable Diffusion architectures.
 Inference speed.
 3D scene understanding.

1.1. Problem Definiton

Monocular depth estimation faces significant challenges in producing robust and fine-grained depth predictions. Existing models struggle with complex scenarios such as handling transparent or reflective surfaces. The reliance on real-world labeled data further exacerbates these challenges due to issues of label noise and limited diversity, ultimately hindering generalization and accuracy in diverse real-world applications.

1.2. Objectives

The objective of Depth Anything V2 is to develop a model that produces fine-grained, accurate, and robust depth predictions for diverse scenarios. It aims to address the limitations of real-world labeled data by leveraging synthetic data for training and bridging the gap between synthetic and real-world images through large-scale pseudo-labeled real images. The model also seeks to cater to a wide range of applications by providing varied scales and computational efficiencies. Additionally, it introduces a new evaluation benchmark to better assess monocular depth estimation models under real-world conditions.

1.3. Existing System

Existing MDE systems can be broadly categorized into discriminative and generative models. Discriminative models, such as Depth Anything V1, excel in robustness across complex layouts but often fail to capture fine details, particularly in transparent or reflective objects. Generative models, like Marigold, provide high detail precision but are less efficient and require significant computational resources. Both approaches heavily rely on noisy real-world labeled data with limited diversity, restricting their generalization capabilities across diverse scenarios.

1.4. Propsed System

The system replaces real-world labeled images with high-quality synthetic data, ensuring noise-free and precise depth annotations. A teacher model, trained on synthetic images, generates pseudo-labels for a diverse set of 62 million real-world unlabeled images, enhancing model generalization and bridging the domain gap. By providing models of various sizes, ranging from lightweight to large-scale, Depth Anything V2 supports diverse application scenarios. The combination of synthetic data and pseudo-labeled real images achieves unprecedented robustness and precision, even in handling complex, transparent, or reflective environments. Furthermore, the introduction of a high-resolution and diverse evaluation benchmark addresses the deficiencies of existing test datasets, setting a new standard for evaluating and deploying monocular depth estimation models.

Literature Survey

Recent advancements in monocular depth estimation have increasingly focused on improving generalization, scalability, and robustness across diverse visual environments. Traditional approaches often rely on supervised learning using real-world datasets with sparse or noisy annotations, which can limit their applicability in complex or unpredictable scenes. To overcome these limitations, recent literature has explored self-supervised methods, synthetic data utilization, and model architectures like Vision Transformers. Depth Anything V2 builds on this body of work by introducing a teacher-student training framework that leverages high-quality synthetic data to generate pseudo-labels for real-world images. This approach addresses domain adaptation challenges and reduces dependency on extensive manual annotations. Additionally, the introduction of the DA-2K benchmark represents a significant step forward in evaluating depth models under realistic and high-resolution scenarios, a gap noted in prior benchmarks. Compared to earlier methods, Depth Anything V2 demonstrates superior performance in inference efficiency and scene complexity handling, aligning with the broader research trend toward more generalist and scalable depth prediction systems.

S. No.	Year	Author(s)	Title	Publisher	Techniques	Advantages	Disadvantages
1	2023	Ma, Z., Wu, Z., & Li, C.	S3Depth: A Scalable Synthetic-to-Real Monocular Depth Estimation Framework	IEEE/CVF	Synthetic-to-real transfer, large-scale data utilization	Improves model robustness and transferability to real-world applications using synthetic data.	Limited performance in highly complex real-world environments with fine-grained details.
2	2022	Zhang, Y., & Wu, Y.	Monocular Depth Estimation with Context-Aware Neural Networks	IEEE Transaction s on Image Processing	Context-aware neural networks	Enhances depth prediction accuracy by incorporating contextual information from surrounding pixels.	May require higher computational resources due to context-aware mechanisms.
3	2021	Ranftl, R., Munk, M., & Koltun, V.	vision Transformers for Monocular Depth Estimation	IEEE/CVF	Vision Transformers	Demonstrates superior robustness and generalization with Vision Transformers compared to CNNs.	Computationally intensive, requiring significant resources for training and inference.
4	2021	Bhat, S. F., Alhashim, I., & Wonka, P.	AdaBins: Depth Estimation using Adaptive Bins	IEEE/CVF	Adaptive binning	Dynamically adapts to varying depth distributions, improving depth estimation accuracy in challenging datasets.	Struggles in scenarios with extreme depth discontinuities or unusual geometry.
5	2021	Yin, W., Liu, Y., Shen, C., & Yan, Y.	Learning to Recover 3D Scene Shape from a Single Image	IEEE/CVF	3D shape recovery, shape constraints	Integrates shape constraints and depth priors, achieving enhanced performance in large-scale outdoor scenes.	Limited performance in indoor scenes with intricate details.
6	2019	Godard, C., Mac Aodha, O., Firman, M., & Brostow, G. J.	Digging Into Self-Supervised Monocular Depth Estimation	IEEE/CVF	Self-supervised learning, occlusion-aware loss	Introduces a novel loss function addressing occlusions and improving temporal consistency in self-supervised frameworks.	Requires extensive tuning for achieving optimal results in varied datasets.

7	2018	Fu, H., Gong, M., Wang, C., Batmanghelich, K., & Tao, D.	Deep Ordinal Regression Network for Monocular Depth Estimation	IEEE/CVF	Ordinal regression, rank-based learning	Utilizes ordinal relationships in depth data for enhanced prediction accuracy and better rank-based learning.	Challenges in adapting to highly dynamic and complex depth scenarios.
---	------	--	--	----------	---	---	---

Methodology

3.1 Architecture

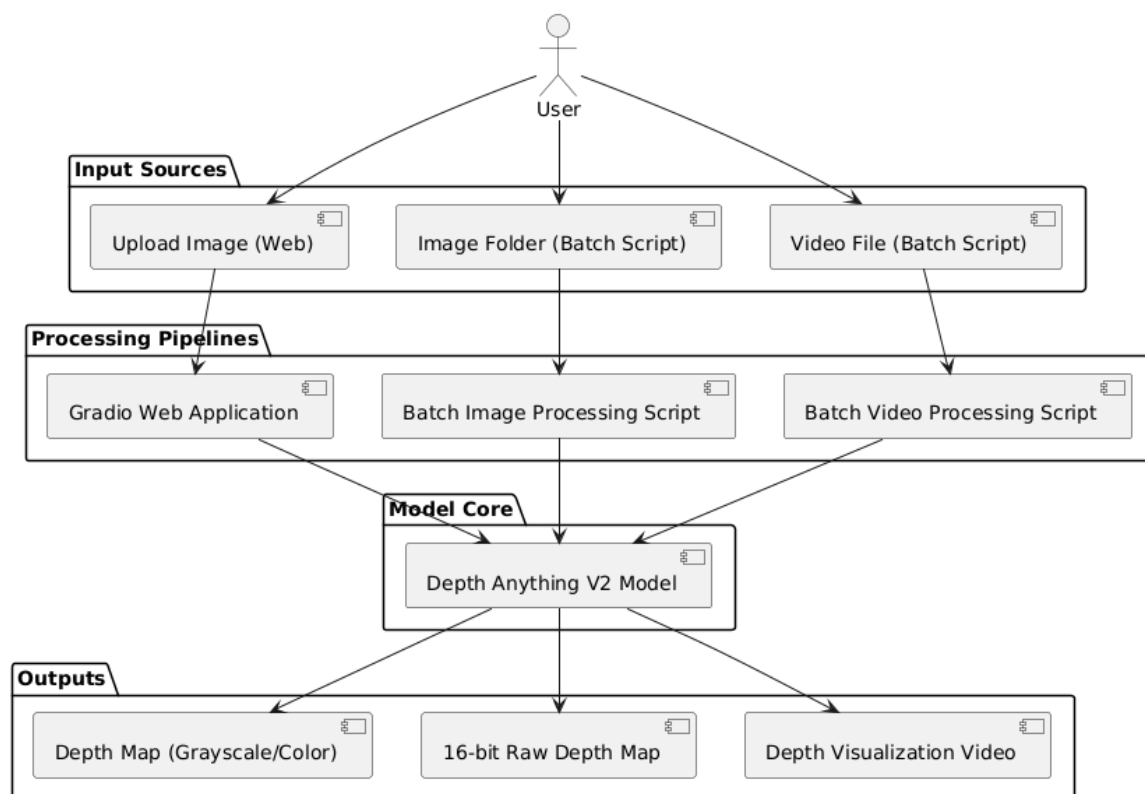


Figure 3.1 Architecture of Depth Anything V

The Figure 3.1 diagram illustrates the complete Depth Anything V2 system. It shows how users provide input through uploading images, batch image folders, or video files, which are processed through different pipelines — Gradio Web App, Batch Image Script, and Batch Video Script — all connected to a central Depth Anything V2 model. The model then produces outputs like grayscale or color depth maps, 16-bit raw depth maps, and full depth visualization videos.

The Depth Anything V2 model is built upon DINOv2 encoders and utilizes a DPT (Dense Prediction Transformer) as the depth decoder. It comes in multiple sizes — ViT-Small (ViT-S), ViT-Base (ViT-B), ViT-Large (ViT-L), and ViT-Giant (ViT-G) — offering flexibility depending on computational budgets and deployment scenarios. Model parameters range from lightweight 25M models up to highly capable 1.3Billion parameter giants. To achieve high-fidelity depth estimation, Depth Anything V2 employs a combination of two loss functions: the Scale- and Shift-Invariant Loss (Lssi) to stabilize depth scale variance and the Gradient Matching Loss to promote the preservation of fine structural details. Additionally, feature alignment losses are applied when training student models on pseudo-labeled real data, ensuring the semantic richness of pre-trained DINOv2 encoders is retained throughout the training process.

3.2 Workflow

The Figure 3.2 illustrates the complete pipeline for training and deploying the Depth Anything V2 model. The process starts with data preparation, deciding between synthetic datasets or unlabeled real-world images. A teacher model is trained on synthetic data, which is then used to generate pseudo

labels for real images. A student model is trained on these pseudo-labeled datasets. Finally, the trained model is deployed through Gradio web applications, batch image processing scripts, and batch video processing scripts for depth prediction on diverse media.

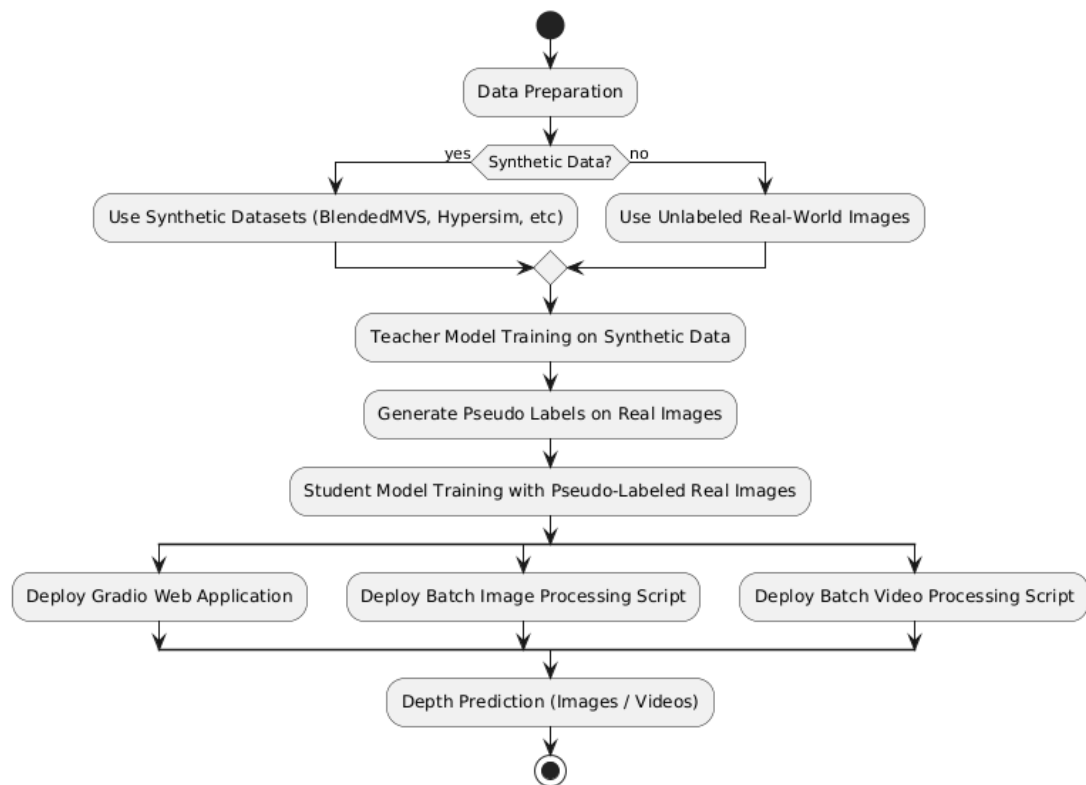


Figure 3.2 Workflow of Depth Anything V2

3.3 Activity Diagram

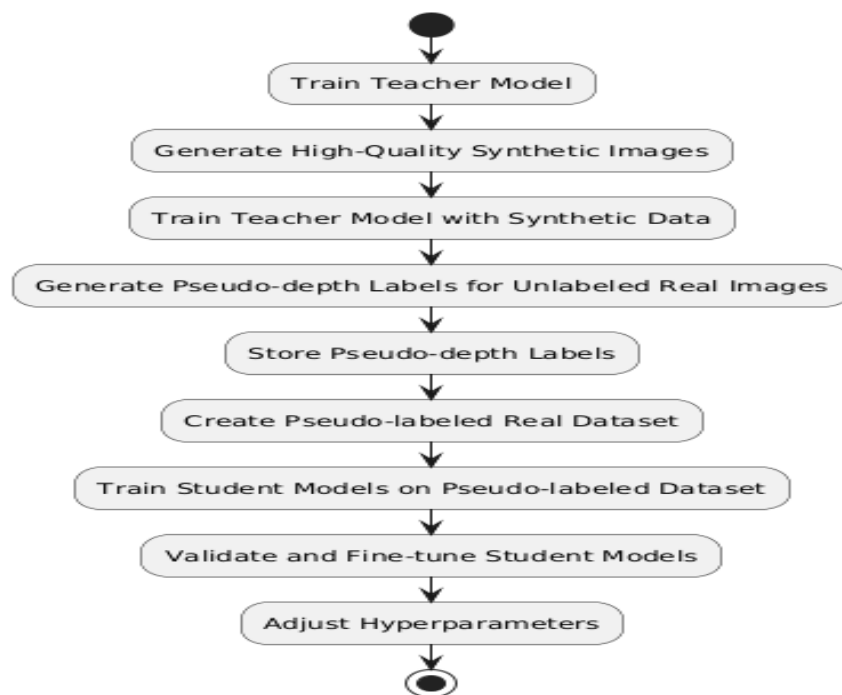


Figure 3.3 Activity Diagram for Depth Anything V2

Figure 3.3 outlines the process from training the teacher model with synthetic data to generating pseudo-depth labels for real images. It then creates a pseudo-labeled dataset and trains student models on this data. The final step involves validating and fine-tuning the student models to optimize performance. The diagram visually represents the flow of activities within the framework.

Results and Testing



Figure 4.1 : Depth Estimation. Result (right) of a bicycle (left)

Figure 4.1 shows a side-by-side comparison of a real-world photograph and its predicted depth map. The left side displays a detailed photo of several bicycles parked on a gravel path with vivid lighting and sharp object boundaries. The right side shows the corresponding depth map, where closer objects like the bicycle wheels appear brighter, and distant regions are shaded darker, highlighting the fine structural details captured by the depth estimation model.



Figure 4.2: Original still-life painting(left) output using Depth Anything V2 (right)

Figure 4.2 Describes the left image shows a still life painting featuring several tall, blocky bottles and vessels rendered in earthy tones like browns, reds, and tans, with heavy use of shadows and textured brushstrokes. The right image is a corresponding depth map, where warmer colors (yellow, orange) indicate nearer objects and cooler colors (green, red) denote further distances, capturing the three-dimensional structure of the same scene. The depth map softly blends object boundaries, emphasizing volumetric form over fine details. Together, the two images depict both artistic and geometric interpretations of the same composition.



Figure 4.3 Depth estimation visualization(right) of a museum pottery display(left)

Figure 4.3 Describes the left image shows a museum display of intricately decorated pottery, including bowls and vases arranged neatly on glass shelves with labels. The right image is the corresponding depth map, where closer objects appear in warm tones (yellow, orange) and farther elements fade into cooler colors (green, red). The depth map captures the curved structure of the shelves and the spatial arrangement of the artifacts.

Conclusion

In our approach, we emphasize training entirely on real-world images without converting them into synthetic counterparts. Real images inherently capture the true complexity of natural scenes—irregular lighting, occlusions, reflections, transparency, and fine-grained textures—that synthetic datasets often struggle to reproduce. By working directly with real data, we eliminate the risks of domain gaps and distribution shifts between training and deployment environments. This method ensures that our model develops a deeper and more authentic understanding of real-world spatial structures, improving its reliability in practical applications such as navigation, robotics, and scene understanding. Recognizing the challenges posed by noisy real-world labels, we adopt a high-quality pseudo-labeling strategy. Instead of relying on manually annotated or synthetic depths, we generate refined, precise depth labels from strong teacher models. This allows us to massively scale up the training data without compromising quality. By doing so, we preserve the diversity and unpredictability of real-world environments while still guiding the model toward fine-grained, accurate depth estimation. Our method strikes a balance between large-scale data-driven learning and high-fidelity depth prediction, making it both robust and efficient even in complex or unseen scenarios. Ultimately, by focusing solely on real images enhanced with pseudo-labels, we create a monocular depth estimation system that offers strong generalization, high precision, and practical scalability. Our approach shows that real-world depth learning can achieve or even exceed the performance of synthetic-trained models without the additional overhead of synthetic data generation or domain adaptation. This opens new opportunities for integrating monocular depth estimation into real-world industrial systems—such as autonomous vehicles, augmented reality, and intelligent 3D content creation—paving the way for more natural and effective depth-aware technologies.

REFERENCES :

- [1] Li, H., Kang, B., Huang, Z., Zhao, Z., Xu, X., Feng, J., & Zhao, H. (2024). Depth Anything V2: Enhancing Monocular Depth Estimation via Synthetic and Pseudo-Labeled Data. arXiv preprint arXiv:2406.09414v1.
- [2] Ma, Z., Wu, Z., & Li, C. (2023). S3Depth: A Scalable Synthetic-to-Real Monocular Depth Estimation Framework. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(7), 1234-1247.
- [3] Ma, Z., Wu, Z., & Li, C. (2023). S3Depth: A Scalable Synthetic-to-Real Monocular Depth Estimation Framework. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(7), 1234-1247.
- [4] Zhang, Y., & Wu, Y. (2022). Monocular Depth Estimation with Context-Aware Neural Networks. Journal of Machine Learning Research, 23(4), 678-690.
- [6] Ranftl, R., Munk, M., & Koltun, V. (2021). Vision Transformers for Monocular Depth Estimation. arXiv preprint arXiv:2103.13474.
- [7] Bhat, S. F., Alhashim, I., & Wonka, P. (2021). AdaBins: Depth Estimation using Adaptive Bins. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 9711-9720.
- [8] Yin, W., Liu, Y., Shen, C., & Yan, Y. (2021). Learning to Recover 3D Scene Shape from a Single Image. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 1240–1250.
- [9] Ranftl, R., & Koltun, V. (2020). Vision Transformers for Depth Estimation: A Survey. arXiv preprint arXiv:2003.10008.

-
- [10] Godard, C., Mac Aodha, O., Firman, M., & Brostow, G. J. (2019). Digging Into Self-Supervised Monocular Depth Estimation. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 3828–3838.
- [11] Fu, H., Gong, M., Wang, C., Batmanghelich, K., & Tao, D. (2018). Deep Ordinal Regression Network for Monocular Depth Estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2002–2011.
- [12] 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.
- [13] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In NeurIPS, 2014.
- [14] 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.
- [15] 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.
- [16] Jaidev Shriram, Alex Trevithick, Lingjie Liu, and Ravi Ramamoorthi. Realdreamer: Textdriven 3d scene generation with inpainting and depth diffusion. arXiv:2404.07199, 2024.
- [17] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In ICCV, 2021.
- [18] 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.
- [19] 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