

# Towards Depth Foundation Model: Recent Trends in Vision-Based Depth Estimation

Zhen Xu\*, Hongyu Zhou\*, Sida Peng, Haotong Lin, Haoyu Guo, Jiahao Shao, Peishan Yang, Qinglin Yang  
Sheng Miao, Xingyi He, Yifan Wang, Yue Wang, Ruizhen Hu, Yiyi Liao, Xiaowei Zhou, Hujun Bao

**Abstract**—Depth estimation is a fundamental task in 3D computer vision, crucial for applications such as 3D reconstruction, free-viewpoint rendering, robotics, autonomous driving, and AR/VR technologies. Traditional methods relying on hardware sensors like LiDAR are often limited by high costs, low resolution, and environmental sensitivity, limiting their applicability in real-world scenarios. Recent advances in vision-based methods offer a promising alternative, yet they face challenges in generalization and stability due to either the low-capacity model architectures or the reliance on domain-specific and small-scale datasets. The emergence of scaling laws and foundation models in other domains has inspired the development of "depth foundation models"—deep neural networks trained on large datasets with strong zero-shot generalization capabilities. This paper surveys the evolution of deep learning architectures and paradigms for depth estimation across the monocular, stereo, multi-view, and monocular video settings. We explore the potential of these models to address existing challenges and provide a comprehensive overview of large-scale datasets that can facilitate their development. By identifying key architectures and training strategies, we aim to highlight the path towards robust depth foundation models, offering insights into their future research and applications.

**Index Terms**—Depth Estimation, Foundation Models, 3D Computer Vision



## 1 INTRODUCTION

Depth estimation stands as a cornerstone in the field of 3D computer vision. This task has been a focal point for researchers due to its critical role in various applications such as 3D reconstruction, 3D generative models, robotics, autonomous driving, and AR/VR technologies. However, algorithms often struggle to achieve high-quality and consistent depth recovery akin to human perception, which leverages prior knowledge of the scene and the world. Traditional depth recovery methods typically rely on active sensing hardware, including commercially available LiDAR, time-of-flight (ToF) sensors, and ultrasonic probes. These sensors estimate depth by measuring the time taken for photons or sound waves to travel back and forth. Despite their accuracy, the high cost of these sensors limits their widespread application. Additionally, active sensors often suffer from low resolution and significant noise interference. For instance, the LiDAR sensor on an iPhone can only achieve a reconstruction resolution within a limited range and struggles with high precision for very close or distant objects. Moreover, these sensors are sensitive to environmental lighting conditions, making them less effective in outdoor, high-light scenarios.

Recently, there has been a growing interest in vision-based depth estimation methods that avoid active depth-sensing hardware, instead relying on readily available cameras commonly found in everyday devices. Compared to active sensor-based approaches, vision methods are cost-

effective, offer an unrestricted depth range, are less affected by environmental conditions, and provide high resolution. For example, a standard iPhone camera can easily capture 4K resolution RGB information. However, existing vision-based depth estimation algorithms still face numerous challenges. Monocular depth estimation, in particular, is highly ill-posed, and standard deep learning algorithms struggle to achieve high-precision results. To introduce constraints and reduce ill-posedness, researchers often explore depth estimation algorithms with multiple camera inputs or scenarios with more extensive scene observations, such as stereo, multi-view, or video-based depth estimation. Yet, these methods are often trained on small-scale synthetic data, leading to instability across spatial and temporal domains, poor generalization across different scenes and input types, and difficulties in overcoming the domain gap between synthetic and real-world data.

With the validation and rise of scaling laws in natural language processing, text-based image generation, and video generation models, the concept of foundation models has emerged. Foundation models are deep neural networks trained on vast amounts of data, exhibiting emergent zero-shot generalization capabilities in other domains. To achieve such capabilities, researchers focus on the scale and variation of input training data, leveraging large-scale models from other domains and ingeniously constructing self-supervision architectures. We define scalable depth estimation models and architectures capable of absorbing massive data as "depth foundation models". For depth estimation tasks, including monocular, stereo, multi-view, and monocular video depth estimation, corresponding foundation models have the potential to address the aforementioned generalization issues and provide key solutions to long-standing Computer Vision tasks.

- The first two authors contributed equally to this work.
- Corresponding author: Hujun Bao.
- Ruizhen Hu is affiliated with Shenzhen University, China. Haoyu Guo is affiliated with Shanghai AI Lab, China. All the other authors are affiliated with Zhejiang University, China.

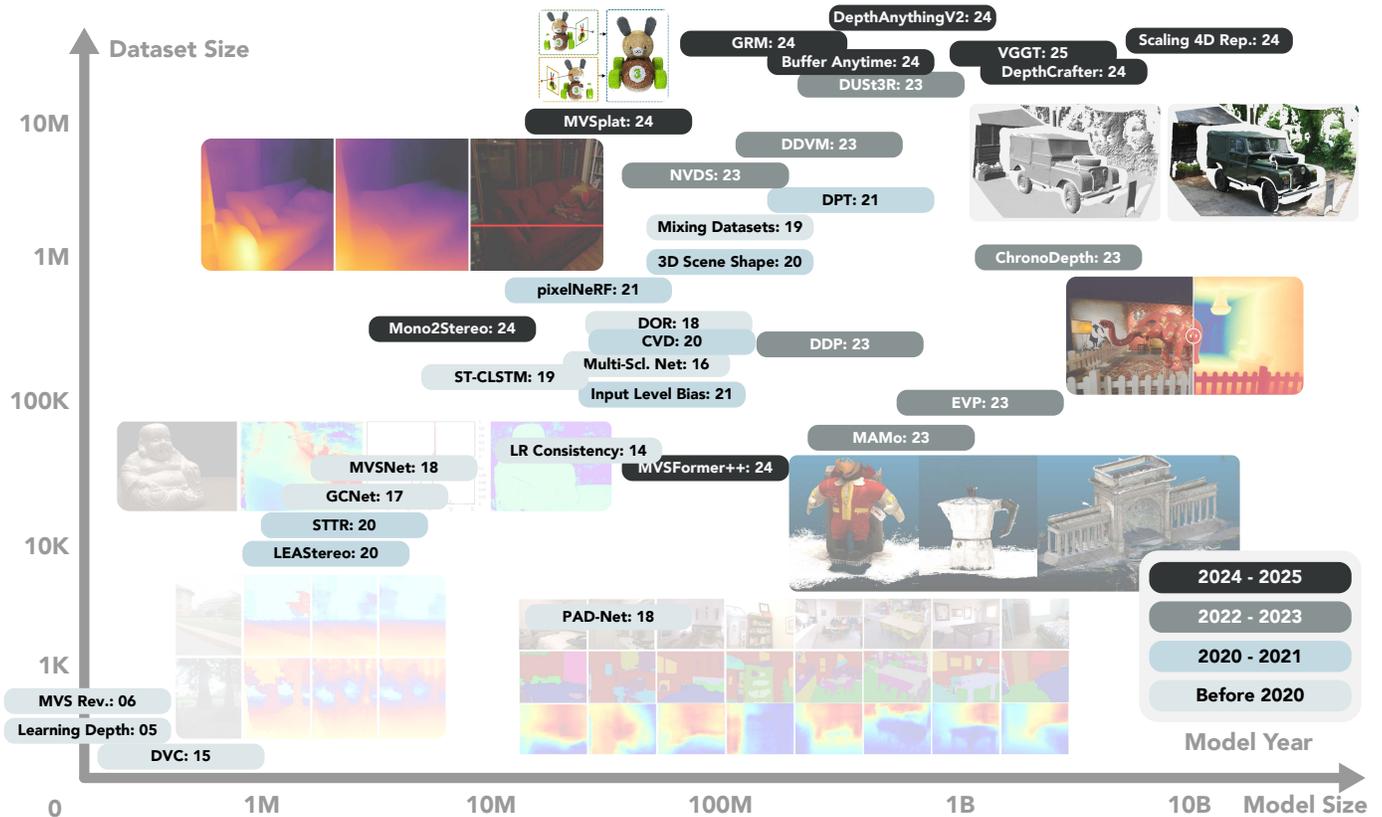


Fig. 1: **Scaling trends in model capacity and data volume for depth estimation.** Each point represents a published method, positioned by its approximate model size (bottom axis, logarithmic scale of parameter count) and dataset size (left axis, logarithmic scale of training images), and colored by publication year. Early models (lightest shading, before 2020) relied on sub-million-parameter networks trained on only thousands of images, yielding limited generalization. Over 2020–2021 (mid-tone), methods grew to tens of millions of parameters and hundred-thousand-image datasets. The most recent models (darkest shading, 2022–2025) scale further into the billions of parameters and multi-million to ten-million-image with strong generalization ability and demonstrating the potential evolution toward depth foundation models.

This paper aims to survey the evolution towards depth foundation models and paradigms for depth estimation across the monocular, stereo, multi-view, and monocular video settings.

- We explore the development of deep learning model architectures and learning paradigms for each task and identify key paradigms with foundational capability or potential.
- To aid the development of such depth foundation models, we also provide comprehensive surveys on large-scale datasets in each respective subfield.
- We also list the current key challenges faced by the foundational architectures in each task to provide insight into future works.

## 2 SURVEY SCOPE

This paper primarily concentrates on depth estimation methods that leverage deep learning, with a particular emphasis on foundation models that utilize large-scale architectures and extensive datasets. We begin by defining depth foundation models and then outline the depth estimation tasks that will be addressed in the following sections.

### 2.1 Definition of Depth Foundation Models

We provide a brief overview of the development of foundation models in language model to facilitate the understanding of depth foundation models. The field of language model have experienced explosive growth with the establishment of foundation models in recent years. This progress stems from the ability of these models to learn universal language and patterns from massive datasets, enabling them to generalize powerfully across various downstream tasks.

Convolutional neural networks and long short-term memory networks [51] plays as the main role at the early stage of language models, with limited network and data scales. The concept of Word Embeddings [52] and the introduction of the self-attention Mechanism [53], allowed the model to process all words in a sequence simultaneously, vastly improving parallel computation efficiency and the ability to capture long-range dependencies. The original Transformer model had a relatively small number of parameters, but its architecture laid the groundwork for subsequent large-scale models. BERTs [54] and GPTs [55] can be considered as the beginning of foundation models in large language models (LLMs). Proposed by Google, BERT is a bidirectional pre-trained model based on the Transformer architecture, enabling better understanding the polysemy of

Domain	Tasks				Dataset Name	Dynamic	Scenes#	Frames#	Resolution	Anno
	Mono	Stereo	MV	Video						
Real-World	✓				A2D2 [1]	Dynamic	3	394K	1920 × 1208	Sparse
	✓	✓			Arogoverse2 [2], [3]	Dynamic	1000	2.14M	2048 × 1550	Sparse
	✓	✓			CityScapes [4]	Dynamic	2975	89K	2048 × 1024	Stereo
	✓				DDAD [5]	Dynamic	200	98K	1936 × 1216	Sparse
	✓	✓			DIML [6]	Static	200+	2M	1334 × 756	Sparse
	✓				DIODE [7]	Static	30	27K	1024 × 768	Sparse
	✓				DSEC [8]	Dynamic	53	26K	1440 × 1080	Sparse
	✓				HM3D [9]	Static	1000	+	+	Dense
	✓				iBims-1 [10]	Static	20	54	640 × 480	Sparse
	✓				Lyft [11]	Dynamic	366	158K	1224 × 1024	Sparse
	✓				Mapillary PSD [12]	Dynamic	50K	750K	1920 × 1080	Sparse
	✓				NuScenes [13]	Dynamic	1K	40K	1600 × 900	Sparse
	✓				NYU [14]	Static	464	435K	640 × 480	Dense
	✓				Pandaset [15]	Dynamic	103	8K	1920 × 1080	Sparse
	✓				Replica [16]	Static	18	36K	1200 × 680	Dense
	✓				Taskonomy [17]	Static	600	4.5M	512 × 512	Dense
	✓				KITTI [18]	Dynamic	22	41K	1242 × 375	Sparse
	✓		✓		DrivingStereo [19]	Dynamic	184	180K	1762 × 800	Sparse
	✓		✓		InStereo2K [20]	Static	50	2K	1080 × 860	Dense
	✓		✓	✓	Argoverse [21]	Dynamic	113	6,624	2056 × 2464	Sparse
	✓		✓	✓	ETH3D [22]	Static	25	1K	6233 × 4146	Dense
	✓			✓	Waymo [23]	Dynamic	1,150	160K	1920 × 1280	Sparse
	✓			✓	UASOL [24]	Static	676	160.9K	2280 × 1282	Dense
	✓			✓	DTU [25]	Static	124	42.5K	1600 × 1200	Dense
	✓			✓	BlendedMVS [26]	Static	113	17k	2048 × 1536	Dense
	✓			✓	Tanks and Temples [27]	Static	21	147K	3840 × 2160	Sparse
	✓			✓	WildRGBD [28]	Static	23K	6M	480 × 640	Dense
	✓			✓	MVImgNet [29]	Static	220K	6.8M	1080 × 1920	Dense
	✓			✓	ARKitScenes [30]	Static	5,047	450M	256 × 192	Dense
	✓			✓	ARKitScenes-HighRes [30]	Static	5,047	450M	1920 × 1440	Dense
	✓			✓	Matterport3D [31]	Static	90	194.4K	1280 × 1024	Dense
	✓			✓	ScanNet [32]	Static	1,513	2.5M	640 × 480	Dense
	✓			✓	ScanNet++ [33]	Static	1,858	3.7M+	1920 × 1440	Dense
Synthetic	✓		✓		Sintel [34]	Dynamic	10	1K	1024 × 436	Dense
		✓			SceneFlow [35]	Dynamic	9	40K	960 × 540	Dense
		✓			CREStereo [36]	Static	0	103K	1920 × 1080	Dense
		✓			FallingThings [37]	Dynamic	3	62K	960 × 540	Dense
		✓			FSD [38]	Dynamic	12	1M	1280 × 720	Dense
		✓		✓	UnrealStereo4K [39]	Static	8	7,720	3840 × 2160	Dense
		✓		✓	Spring [40]	Dynamic	47	6K	1920 × 1080	Dense
		✓		✓	TartanAir [41]	Static	163	1M	640 × 480	Dense
		✓		✓	VirtualKITTI2 [42]	Dynamic	5	25K	1242 × 375	Dense
				✓	MatrixCity [43]	Static	3K	519K	1000 × 1000	Dense
	✓			✓	MVS-Synth [44]	Dynamic	120	12K	1920 × 1080	Dense
	✓			✓	3D Ken Burns [45]	Static	32	536K	512 × 512	Dense
	✓			✓	Dynamic Replica [46]	Dynamic	484	145K	1280 × 720	Dense
	✓			✓	OmniObject3D [47]	Static	6K	600K	800 × 800	Dense
	✓	✓		✓	IRS [48]	Static	70	100K	960 × 540	Dense
				✓	PointOdyssey [49]	Dynamic	131	200K	540 × 960	Dense
				✓	BEDLAM [50]	Dynamic	10,450	380K	1280 × 720	Dense

TABLE 1: **Datasets.** The term  indicates that the scenes in these datasets are object-centric. Similarly,  and  refer to indoor and outdoor scenes, respectively. The depth annotations can be classified as dense or sparse, depending on whether most pixels have corresponding depth values. †: HM3D is a digital twin dataset created from real-world data, with the number of frames and resolutions being user-defined.

words in a sentence. Bert is trained on Toronto BookCorpus (800 million words) and English Wikipedia (2.5 billion words). BERT-Base has 110 million parameters, and BERT-Large has 340 million parameters. Proposed by OpenAI, GPT is a unidirectional generative pre-trained model based on the Transformer architecture. GPT models learn language patterns by predicting the next word, excelling in text generation tasks. The GPT-3 is trained on a dataset which is larger than 45 TB, along with 175 billion parameters.

The development of depth estimation models is illustrated in Fig. 1. Considering the foundation models scale in the areas of language models, we define a depth foundation model as one that is trained on a large-scale dataset (over 10 million images) and employs models with a substantial number of parameters (over 1 billion). Additionally, depth foundation models should exhibit strong generalizability across multiple data domains.

## 2.2 Depth Estimation Tasks

This survey covers several tasks, including monocular depth estimation, stereo depth estimation, multi-view depth estimation, and monocular video depth estimation using foundation models. Let  $\mathbf{I} = \{I_{k,t}, k = 1, \dots, \mathcal{K}, t = 1, \dots, \mathcal{T}\}$  represent a collection of RGB images, where  $\mathcal{K}$  denotes the number of cameras and  $\mathcal{T}$  is the number of timestamps for the frames. In the case of monocular depth estimation, the input consists of a single image  $I_{1,1}$ . For stereo depth estimation, the input comprises a pair of images  $\{I_{1,1}, I_{2,1}\}$ . In multi-view depth estimation, the input is a set of images captured at the same timestamp but varying in spatial locations, represented as  $\{I_{k,1}, k = 1, \dots, \mathcal{K}\}$ . For monocular video depth estimation, the input consists of a sequence of images captured by a monocular camera at different timestamps, represented as  $\{I_{1,t}, t = 1, \dots, \mathcal{T}\}$ . The scope of our survey excludes the task of multi-view video depth estimation, which can be represented as the most general form of inputs:  $\{I_{k,t}, k = 1, \dots, \mathcal{K}, t = 1, \dots, \mathcal{T}\}$ . This is due to the fact that foundation models for this task have not yet been thoroughly explored.

For each task, we begin by reviewing the background and evolution of deep learning models specific to the task. We then delve into the development of foundation models. Prominent examples of foundation models include transformer-based models and diffusion models. Furthermore, we also discuss the large-scale datasets used for training these foundation models, encompassing both synthetic and real-world datasets, which enable the models to generalize effectively across diverse scenes. Finally, we address valuable problems faced by existing depth foundation models.

## 3 OVERVIEW OF DEPTH ESTIMATION

In this section, we provide an overlook of paradigms and datasets used in depth models.

**Paradigms.** In monocular image depth estimation, models have progressed from direct depth regression to affine-invariant depth, depth classification, and canonical camera depth. This progression facilitates more accurate depth map

predictions using just a single input image. In stereo image depth estimation, by utilizing the principles of stereo geometry, models can concentrate on matching corresponding pixels in image pairs. This has driven the evolution of paradigms from cost-volume methods to the attention mechanism, iterative optimizers, and ultimately to scalable training approaches. In multi-view image depth estimation, similar to stereo paradigms, the incorporation of multi-view information has facilitated an evolution from patch-match stereo to cost volume methods, followed by a coarse-to-fine strategy and the implementation of token attention mechanisms. In monocular image depth estimation, by incorporating an additional dimension of timestamps, models leverage temporal correlation and test-time optimization paradigms to establish the relationship between temporal and spatial information. Additionally, scaling up training enhances depth estimation performance and guides models toward becoming depth foundation models.

**Datasets.** We summarize the datasets commonly used in depth estimation tasks in Tab. 1. These datasets can be categorized into two classes: real-world captured and synthesized. Each dataset may be applicable to multiple tasks, and we specify the tasks associated with each dataset in the table. Furthermore, we present details about the scenes within the datasets, including whether they provide metric information, if they include dynamic scenes, the number of scenes and frames, the image resolution, and whether the annotations are dense, based on the proportion of pixels that have corresponding depth values.

## 4 MONOCULAR IMAGE DEPTH ESTIMATION

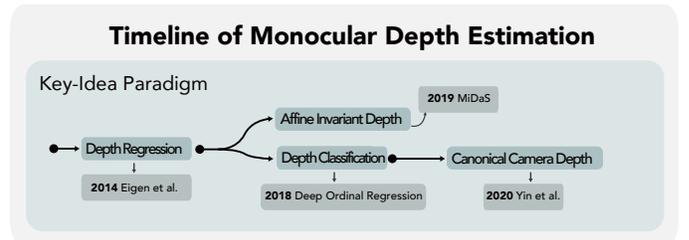


Fig. 2: Overview of the key-idea paradigm evolution of monocular image depth estimation. From early direct regression and classification methods, through affine-invariant and canonical camera depth estimation, monocular depth estimation models have shown increasingly stronger generalization capabilities, paving the way for the emergence of depth foundation models.

Monocular depth estimation aims to predict per-pixel depth distances for a scene from a single RGB image, establishing a geometric mapping from 2D images to 3D scenes. The core challenges lie in scale ambiguity and the lack of geometric information inherent in monocular vision. Compared to active depth sensors, monocular depth estimation requires no additional hardware and relies solely on image content to infer scene structure, making it valuable for various applications such as 3D reconstruction, video editing, autonomous driving, and augmented reality. In recent years, with the advancement of deep learning techniques,

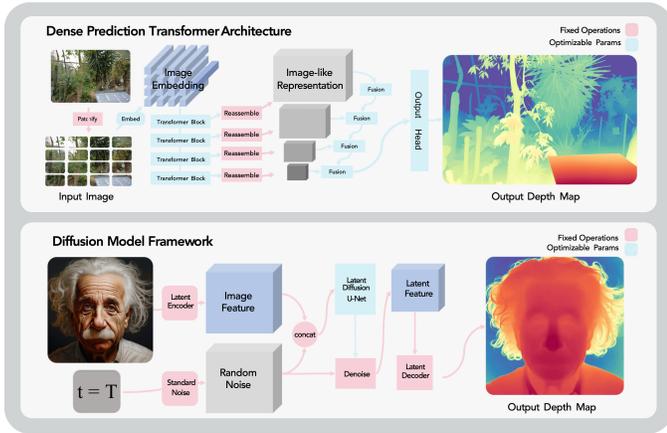


Fig. 3: **Representative pipelines of monocular image depth estimation.** The pink component denotes operations without learnable parameters or with fixed parameters, while the blue component indicates operations with optimizable parameters. Vision Transformer-based approaches, leveraging their lightweight architectures, enable real-time monocular depth estimation. However, due to the presence of convolutional operations in their architectures, they may lose detailed features. Diffusion Model-based methods treat RGB images as conditional inputs, effectively preserving fine-grained details. Nevertheless, their denoising processes impose computational costs, making it challenging to achieve real-time performance.

monocular depth estimation has gradually shifted from traditional geometric methods to data-driven end-to-end learning paradigms, evolving toward foundation models with strong generalization capabilities, high accuracy, and robustness.

**Evolution of model architectures.** The evolution of monocular depth estimation has been driven by progressive innovations across network architectures, depth representations, and learning paradigms. Early approaches predominantly employed handcrafted image processing pipelines [56], [57], [58], [59], with seminal works like CNN [60] establishing the deep learning foundation. Subsequent refinements incorporated U-Net [61] and ResNet [62] residual blocks to enhance spatial continuity. The advent of dense prediction transformers [63] marked a paradigm shift: DPT [63] introduced global attention mechanisms through patch-wise sequence modeling, addressing CNN-inherent locality constraints. The latest frontier involves diffusion models [64], [65], exemplified by Marigold [66], which leverage conditional denoising frameworks to transfer generative priors into depth estimation, significantly improving geometric consistency and cross-domain generalization.

**Evolution of method paradigms.** Parallel advancements emerged in depth representation and learning strategies. Early methods [60], [61], [62] focused on absolute depth prediction but faced scale ambiguity challenges, prompting innovations like scale-invariant loss [67], [68] for relative depth estimation. Another advancement came through bins classification approaches [69], [70], [71], [72], [73], [74], which discretize the continuous depth space into multiple

bins. This reformulates depth estimation as a per-pixel classification problem, effectively handling non-uniform depth distributions and capturing uncertainty in predictions, particularly at depth discontinuities. Recent breakthroughs [75], [76], [77], [78], [79] integrate camera parameters to resolve metric ambiguity while maintaining scale awareness. The data paradigm evolved through large-scale multi-view stereo (MVS) datasets like MegaDepth [80], enabling unprecedented generalization. Supervision methods expanded beyond fully labeled training: self-supervised approaches exploited photometric consistency in stereo sequences [81], while Depth Anything [82] demonstrated pseudo-label distillation’s effectiveness for knowledge transfer. Multi-task frameworks [83], [84], [85], [86] further enhanced robustness through joint depth-flow-pose estimation with geometric constraints, complemented by geometry-aware losses like Virtual Normal that explicitly enforce surface regularity. This multifaceted progression underscores how architectural innovation, representation learning, and supervision paradigms collectively advance monocular depth estimation toward human-level scene understanding.

**Evolution towards foundation models.** Models for monocular depth estimation are becoming larger and more data-intensive, evolving towards foundation models. In terms of *Model Size*, advanced foundation models primarily utilize two key methods: Vision Transformer (ViT) architectures [87] and diffusion-based generative models [88]. Dense Prediction Transformers (DPT) are the main architecture for modern depth estimation [63], featuring 343 million parameters. Unlike older fully convolutional networks, DPT’s ViT backbone keeps high-resolution features and a global view throughout, leading to more detailed and consistent depth estimates. Pre-Trained Diffusion models are used as strong depth estimators to improve how well monocular depth estimation works, thanks to their existing visual knowledge, which ranges from 200 million to 1 billion parameters. For example, Marigold [66] slightly modifies text-to-image latent diffusion models (like Stable Diffusion v2 [89]) to predict depth from images. In terms of *Data scale*, powerful hardware and better cameras have led to many high-quality depth datasets, typically ranging from thousands to millions of images; further details can be found in Tab. 1. The current trend is to train models on large datasets to make them more adaptable. When training for large-scale depth estimation, models typically use loss functions that ignore scale and shift differences in the data [67]. This means models learn to predict relative depth, making relative depth estimation a standard task in computer vision. However, metric depth estimation, which provides absolute distance measurements, is crucial for real-world applications. Currently, there are two main ways to get metric depth: 1) Methods like those in [76], [77], [78], [79], [90], [91] combine camera intrinsic estimations with relative depth predictions to get metric-scaled outputs through geometric calculations. 2) Directly learning metric depth: Approaches such as those in [92], [93], [94] train models directly on data that includes scale information, allowing them to predict metric depth without extra steps.

**Valuable problems.** Recent methods have made notable progress, yet critical challenges persist across four primary

dimensions. *Depth accuracy.* Depth accuracy enhancement remains the foremost pursuit, aiming to develop a “visual LiDAR” system where RGB cameras rival dedicated depth sensors. Current methods exhibit 4-5% relative error on standard benchmarks [78], escalating to 20-50% in challenging scenarios [95], while hardware sensors consistently achieve sub-1% accuracy. *Absolute scale recovery.* Though improved through works like Metric3D [77], it demonstrates fragility under complex illumination and texture-deficient conditions, necessitating more robust geometric priors. The *Data Efficiency Bottleneck* manifests through compounded limitations: pseudo-label noise restricts supervision quality while synthetic-to-real domain gaps constrain model generalization, demanding innovative low-cost high-precision annotation paradigms. *Multi-task generalization.* This presents an open research frontier, as current approaches struggle to unify depth estimation with complementary tasks like semantic segmentation and surface normal prediction within a single foundational model architecture. These interconnected challenges collectively underscore the need for fundamental breakthroughs in geometric understanding, data utilization, and cross-task knowledge integration to bridge the performance gap between learning-based methods and physical sensing systems.

## 5 STEREO IMAGE DEPTH ESTIMATION

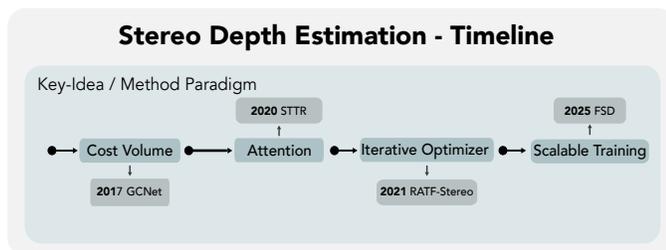


Fig. 4: **Overview of the key-idea paradigm evolution of stereo image depth estimation.** The paradigms have transitioned from cost-volume methods to attention mechanisms and iterative optimization techniques to effectively match the features of stereo images. The incorporation of monocular and diffusion priors facilitates large-scale training, paving the way for foundation models.

Stereo Depth Estimation aims to estimate the per-pixel depth for a scene given the relative pose of a pair of stereo cameras and a pair of RGB observations as input. Compared to monocular depth estimation, stereo depth estimation can utilize disparity priors to estimate depth through epipolar feature matching and triangulation, which theoretically yields better results than monocular depth estimation. Additionally, since stereo depth estimation incorporates the known relative pose of the cameras, the estimated depth inherently includes scale information, addressing the scale ambiguity present in monocular depth estimation. The core challenge of stereo depth estimation lies in accurately establishing the correspondence between pixels in the left and right images. The traditional matching process can be easily affected by factors such as changes in lighting, occlusions, repetitive textures, and weak textures.

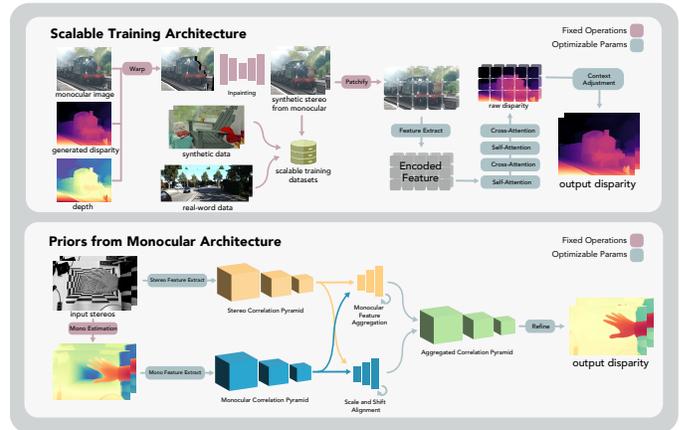


Fig. 5: **Representative pipelines of stereo image depth estimation.** The first architecture is for scalable training, which leverages all available datasets along with pseudo stereo pairs synthesized from a monocular dataset, to train a foundation model. The second architecture is migrating knowledge from a monocular foundation model to the stereo model, making it possible to achieve a stereo foundation model from relatively small-scale training datasets.

With the rapid development of deep learning technology, stereo depth estimation has gradually become popular for using neural networks to replace traditional epipolar feature matching methods. The neural network architecture paradigms mainly include the CNN-based cost-volume paradigm, the Transformer-based attention paradigm, and the RNN-based iterative optimization paradigm.

**Evolution of model architectures.** The CNN-based models are one of the most classic and widely used technical approaches in stereo depth estimation [39], [96], [97], [98], [99], [100], [101], [102], [103], [104]. Early methods, such as PSMNet [97], used simple convolution layers to extract features, while subsequent works adopted more complex architectures like ResNet [105], DenseNet [106], and NAS search techniques [101] to extract more robust features. Since 2020, with the success of Transformers [53] in the field of language models and the advancements of Vision Transformers [87] in computer vision, the attention paradigm was proposed [107], [108], [109], [110], [111], [112], [113]. To reduce computational complexity, [108] introduces local attention or sparse attention mechanisms. Inspired by RNNs [114] and the application of the RAFT network architecture in optical flow estimation, [115], [116], [117] construct a correlation pyramid and use an RNN model to iteratively optimize the disparity map. [117] introduces learnable update strategies that dynamically adjust the update direction through attention mechanisms. [38], [118], [119], [120] and [121] utilize monocular foundation models as priors and diffusion models as inpainting tools for stereo generation, respectively, advancing stereo depth estimation through large-scale training.

**Evolution of model paradigms.** The CNN-based Cost Volume paradigm [39], [96], [97], [98], [99], [100], [101], [102], [103], [104] constructs a cost volume using 2D or 3D CNNs to represent the matching cost between the

left and right images, and then perform Cost Aggregation as a post-processing step to ultimately regress the disparity map. The attention paradigm [107], [108], [109], [110], [111], [112], [113], [117] models the epipolar matching problem as a sequence-to-sequence task, utilizing Self-Attention and Cross-Attention mechanisms to capture long-range dependencies between pixels. Within a single image, the self-attention mechanism models the relationships between pixels to capture contextual information, while the Cross-Attention mechanism matches corresponding pixels between the left and right images. Specifically, each pixel in the left image interacts with all pixels in the right image to compute attention weights. The iterative optimization paradigm [115], [116], [117] dynamically updates the Cost Volume or feature maps based on the current disparity map to capture more accurate matching information. It can balance performance and speed through early stopping, making it suitable for time-sensitive applications. Additionally, this paradigm exhibits strong robustness to initial errors. With significant advancements in simulation technology, generative techniques, and monocular depth estimation methods in recent years, foundation models [38], [41], [119], [120], [121], [122], [123] in the field of stereo depth estimation are beginning to emerge.

**Evolution towards foundation models.** Foundation models are emerging as the new paradigm for stereo depth estimation, leading to an increase in data intensity. However, in contrast to foundation models used for monocular depth estimation, the *Model Size* for stereo depth estimation does not see a significant increase, ranging from 3.5 million to 11 million parameters [96], [101], [107], [115]. Instead, advances in monocular depth estimation foundation models allow stereo tasks to leverage monocular priors to enhance depth estimation. There are currently two main approaches to utilizing these priors. One approach [38], [118] involves injecting features from monocular depth models into the cost volume. The other approach [119], [120] applies stereo metrics to scale monocular depth estimates and then uses refinement networks to achieve better depth estimation results. There are two methods to enable *Large Data Scale Training* for stereo depth estimation. 1) creating high-fidelity virtual scenes using advanced simulation and rendering technologies, such as FSD [38] (1 million image pairs) and TartanAir [41] (1 million image pairs) and Stereo anything [122] (30 million image pairs) generate virtual stereo pairs by using estimated scene depth from monocular images [67], [82], allowing pixels to be warped to pseudo-stereo viewpoints and inpainting to fill in gaps. Recent advancements in diffusion models have led to stereoGen [121] (35 thousand image pairs) using stable diffusion as an inpainting tool.

**Valuable Problems.** Research on foundation models for stereo depth estimation is still in its early stages, and we believe there are several key issues worth exploring in the future: *Limited data*. Generated stereo data from monocular faces challenges such as insufficient accuracy in monocular depth estimation and difficulties in filling in warp holes. Additionally, there remains a domain gap between synthetic

data and the real world, and the diversity of the synthetic datasets is still not rich enough. *Lack of end-to-end training paradigm*. Current methods that leverage monocular priors treat monocular foundation models as cues, lacking end-to-end training of large model parameters. The model parameter counts are relatively small, and there is a lack of foundation model designs tailored for stereo tasks. *Limit utilization of available datasets*. There is an insufficient application of cross-domain datasets. Aside from Stereo Anything [122], most approaches typically utilize only one or two datasets, failing to fully leverage the existing stereo datasets. *Under-utilization of the diffusion architecture and priors*. A diffusion foundation model for stereo depth estimation, similar to the Marigold model [66], has not been explored.

## 6 MULTI-VIEW IMAGE DEPTH ESTIMATION

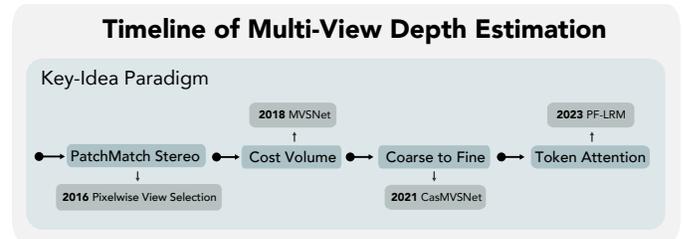


Fig. 6: Overview of the key-idea paradigm evolution of multi-view image depth estimation. From early heuristic matching, through 2D CNNs and 3D CNNs, to advanced frameworks utilizing transformers and diffusion models, multi-view depth estimation models have shown increasing robustness and accuracy.

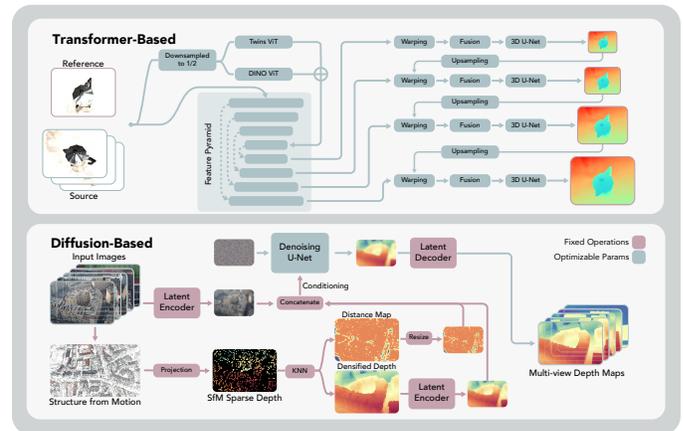


Fig. 7: Representative pipelines of multi-view image depth estimation. As shown in the upper part of the figure, the Transformer-based architecture utilizes a Transformer model to extract global features and employs a 3D-geometry Transformer to facilitate cross-view information interaction (the Cost Volume method is depicted in the figure). In contrast, the Diffusion-based model in the lower part of the figure leverages the SfM (Structure from Motion) point clouds from multi-view images as prior information and generates depth maps through a diffusion model.

**Evolution of model architectures.** The model architectures for multi-view depth estimation have evolved from early heuristic matching to state-of-the-art Transformer and diffusion-based models. Early approaches [124], [125] typically assume that camera parameters were already known, obtained via SfM [126], [127], [128], SLAM [129], [130], or robotic arm control [131], [132], and rely on traditional dense pixel matching techniques [133] for depth estimation. This heuristic matching lays the groundwork for future developments. With the advent of deep learning, researchers have begun employing 2D CNNs [134], [135] to process the cost volume, which requires low computational cost and enables real-time scene reconstruction. At the same time, alternative approaches utilize 3D CNNs [44], [117], [136], [137], [138], [139], [140], [141], [142], [143] to achieve more accurate depth estimation. Recent trends have started to explore the use of Transformers [144], [145], [146], [147], [148], [149], [150], [151], [152], [153], [154] or diffusion models [155], which leverage global self-attention mechanisms or diffusion-based strategies to capture richer feature representations and enhance inference quality.

**Evolution of method paradigms.** The method paradigms for multi-view depth estimation have evolved from early PatchMatch-based matching strategies to modern token attention mechanisms. Traditional methods [124], [125] rely on the PatchMatch algorithm [133] for pixel-level matching—a simple yet influential approach that sets the stage for subsequent innovations. With the integration of deep learning, constructing a cost volume by back-projecting multi-view features using camera parameters becomes the mainstream [44], [134], [136], [137], [138], allowing for more accurate depth regression by computing feature correlations. Given the high memory consumption of cost volumes, modern approaches [117], [139], [140], [141], [142], [143] adopt a coarse-to-fine strategy: starting with a low resolution and a large number of depth hypotheses, and progressively refining the cost volume with predictions from previous stages to achieve a balance between high resolution and high accuracy. Looking ahead, emerging research is investigating the incorporation of token attention [147], [149], [150], [151], [152], [153], [154], [156], [157], [158], [159], [160], [161], [162] mechanisms in the depth estimation process, aiming to better capture both local and global contexts and offering promising new directions for multi-view depth estimation.

**Evolution towards foundation models.** In terms of *Model Size*, recent advances have been driven by the emergence of Transformer-based foundation models. MVSTR [144] is the first to introduce the Transformer architecture to multi-view depth estimation, employing a global-context Transformer and a 3D-geometry Transformer for intra-view global feature extraction and inter-view information interaction. MVFormer [145] proposes to use a pre-trained Vision Transformer (ViT) to enhance Multi-View Stereo (MVS) tasks, leveraging priors learned from large-scale datasets. MVFormer++ [146] further employs a pre-trained DINOv2 model with 1.1 billion parameters and introduces distinct attention mechanisms tailored for the feature encoder and cost volume regularization. Both PF-LRM [163] and DUST3R [156] utilize ViT to directly predict Point Maps, with DUST3R employing the CroCo model featuring a com-

plete encoder-decoder Transformer architecture. In terms of *Data Scale*, the evolution shows a clear trend from small-scale datasets to large-scale training. Early methods like MVSNet [138] used only 27K images from the DTU dataset, while MVFormer++ [146] expanded to 40K images from the DTU and BlendedMVS datasets. The breakthrough came with sparse-view methods: PF-LRM [163] utilizes 1M objects from Objaverse and MVImgNet datasets, DUST3R [156] is trained on 17M image pairs from eight datasets, including Habitat, MegaDepth, and Waymo, GRM [147] uses 40M objects from Objaverse, and VGGT [149] employs 30M images from multiple datasets, including Co3Dv2, BlendMVS, and MegaDepth. These models establish a paradigm shift from traditional optimization-based methods to data-driven architectures capable of unified feature learning and cross-view reasoning.

**Valuable Problems.** *Sparse view reconstruction.* In real-world applications, capturing a scene with dense and complete views may be feasible due to constraints on reachability. Therefore, the ability to reconstruct complete scene geometry from partial observations by leveraging prior knowledge represents an important direction for future research. *Find-grained depth estimation.* Current feed-forward methods [144], [145], [146], [147] can efficiently predict multi-view image depth in a single forward pass. However, accurately capturing fine-grained geometry remains challenging, as it demands high-precision geometric prediction capabilities from neural networks. *Depth estimation of objects with complex materials.* Reflective or transparent scenes pose significant challenges for geometry estimation due to their complex optical properties. Incorporating learned priors into depth estimation presents a promising approach for accurately capturing the geometry of complex materials.

## 7 MONOCULAR VIDEO DEPTH ESTIMATION

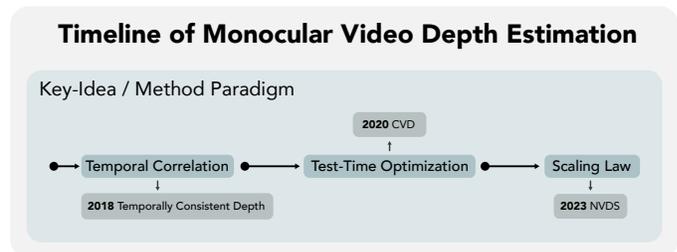


Fig. 8: Overview of the key-idea paradigm evolution of monocular video depth estimation. From RNN-based temporal modeling (2019) to CNNs with test-time optimization (2020, CVD), transformer-based scaling (2023, NVDS), and video diffusion for enhanced stability and generalization.

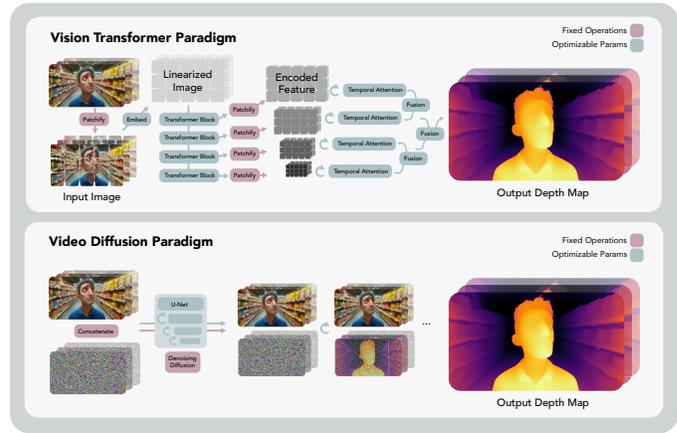
Video depth estimation aims to estimate per-frame depth from a given monocular video while ensuring temporal consistency across the sequence. Compared to monocular depth estimation, its primary challenge lies in maintaining consistency over time. In contrast to multi-view stereo (MVS) [146], [147] and other multi-view depth estimation methods [144], [145], it needs to handle dynamic scenes, which imposes further challenges. Since video depth estimation requires both temporally consistent and accurate depth

predictions, and its technical paradigm integrates elements from both monocular and multi-view depth estimation, we consider it the ultimate problem that a depth foundation model should address.

**Evolution of model architectures.** With the advancement of deep learning technologies and network architectures, various network structures have been applied in recent years to improve the temporal consistency of depth estimation. With the rise and widespread adoption of RNNs in language models, in 2019, [164] proposed using the LSTM mechanism to integrate temporal information, thereby enhancing the stability of video depth estimation. Meanwhile, in image tasks, CNN architectures (including U-Net and ResNet) have shown excellent performance. In 2020, [165] introduced test-time optimization to single-frame depth estimation methods based on CNN architectures, amplifying the capabilities of traditional CNNs through optimization and alignment. Subsequently, attention mechanisms gained significant attention. [151], [152], [153], [166] applied attention mechanisms to integrate temporal features, maintaining memory through attention mechanisms rather than LSTMs, achieving more accurate and consistent depth estimation results. Later, with the remarkable success of diffusion models in the field of image generation and the impressive results of video diffusion models, in 2024, [167], [168] proposed using video diffusion models to enhance consistency, achieving unprecedented stability and predictive performance.

**Evolution of method paradigms.** In recent years, innovative technical paradigms have been proposed to enhance the stability of video depth estimation. In 2019, [164] introduced a memory mechanism to enable networks to integrate multi-frame information, implicitly encoding content from other frames to assist in depth estimation for the current frame. Subsequently, several works explored similar memory paradigms, such as MAMo [166]. In 2020, [165] proposed using test-time optimization to post-process predicted depth results. Since this approach better leverages prior techniques like SfM and SLAM while also performing bundle adjustment (BA) on depth, it often yields superior results compared to purely generalized methods. Later, with the rise of scaling laws, interest in generalized methods was reignited, and attention shifted to the generation and utilization of large-scale training data. In 2023, [169] introduced a network-based post-processing method for monocular depth estimation results, finding a middle ground between direct video depth prediction and test-time training (TTT). It also introduced the first representative large-scale Video Depth in the Wild (VDW) dataset. Video depth estimation has since begun to evolve toward the development of foundation models.

**Evolution towards foundation models.** In terms of *Model Size*, recent work has focused on using scalable network architectures to model video depth. Buffer Anytime [171] proposed using a Temporal ViT architecture with 343 million parameters to integrate temporal information and generate a large amount of pseudo ground truth (GT) based on single-image priors. Specifically, the input images are first patchified and converted into tokens recognizable by ViT. Then, self-attention is applied to these tokens. To fuse temporal information, cross-attention is performed on the



**Fig. 9: Representative pipelines of monocular video depth estimation.** The Vision Transformer (ViT) [63] paradigm processes patchified inputs through self-attention and temporal attention mechanisms to integrate spatial and temporal depth cues [170]. The Video Diffusion paradigm leverages denoising diffusion models, using concatenated depth and image features, optionally enriched with CLIP embeddings, to generate consistent video depth estimates [168]. These scalable architectures enhance generalization and enable zero-shot depth estimation across diverse datasets.

single-frame attention results, reprocessing the features to obtain tokens integrated with temporal information. Finally, the tokens are decoded and un-patchified to produce the final depth estimation results. Another technical paradigm [82], [167], [168] proposed using video diffusion models with 200 million to 1 billion parameters to aid in recovering video depth. Specifically, these methods use the input video as a condition in the video diffusion denoising process. For the noisy depth during denoising, it is concatenated with the corresponding frame’s RGB color and, optionally, with CLIP embeddings as input to the denoising UNet. In terms of *Data scale*, the insufficient volume of video depth data has always been a critical issue. In 2023, NVDS [169] constructed the first large-scale, diverse video stereo depth dataset (14.2K videos, 2.24M frames) by using a video stereo matching method to generate pseudo labels. In 2024, methods like [82], [167], [168] addressed the lack of video data by leveraging video priors from video generation models. Similar to NVDS [169], DepthCrafter [168] annotated 200K videos using a video stereo matching method. Depth Any Video [82] created 40K videos, totaling 6M frames, using a game engine. Video Depth Anything [170] proposed using an image teacher to provide pseudo labels to compensate for the shortage of video depth ground truth (GT) labels.

**Valuable problems.** *Geometric inconsistency.* When camera motion is present in a monocular video, the ability to estimate consistent depth for the same statistical scene observed across multiple frames is crucial. To achieve this, jointly modeling camera motion and depth estimation is a promising approach and future direction for enhancing geometric consistency. *Temporal inconsistency.* In the presence of dynamic objects, the lack of multi-view observation in monocular video makes it particularly challenging to estimate their geometry. The method needs to learn strong

priors to predict temporally consistent depth for dynamic objects across multiple frames. *Monocular video depth training data*. Learning strong dynamic priors needs extensive and diverse training data. However, due to the presence of dynamic objects, collecting real-world ground-truth monocular video training data with accurate ground-truth depth often requires additional depth sensors, limiting its scale-up capabilities. Scaling up monocular video training data, either using real-world Internet unlabeled video with a self-supervised training strategy or simulating and rendering realistic synthetic data, is a valuable direction to be explored.

## 8 APPLICATIONS

### 8.1 3D Reconstruction

For multi-view 3D reconstruction task, once the depth for each view is estimated, they can be directly fused to obtain the reconstructed scene geometry, which is usually represented as a point cloud or triangle mesh. MVSNNet [138] and CasMVSNNet [140] adopt fusible [172], which first converts depth map of each view into a point cloud, then filters out points with poor consistency by projecting them into other views and checking, and finally fuses the point clouds from all views to obtain the reconstruction result of the entire scene. Simplerecon [173] and Murre [155] use TSDF fusion [174], which converts depth maps into sparse truncated SDF grids, averages them across multiple views, and finally extracts the surface using marching cubes to obtain a triangle mesh. Direct fusion methods are relatively efficient but require high-quality depth maps. NeuRIS [175] and MonoSDF [176] use the relative depth predicted from a monocular depth estimation method as a supervision signal, constraining the SDF field through a specially designed depth loss, thereby enhancing the quality of the 3D reconstruction.

### 8.2 Novel View Synthesis

In recent years, NeRF-based [177] and 3D Gaussian Splatting-based [178] methods have made significant progress in the task of novel view synthesis. High quality depth estimations can serve as strong priors for these models, enhancing their performance and accelerating convergence. Specifically, the depth map can be utilized as sampling guidance for NeRF-based methods [176], [179], [180] and as a coarse initialization for 3D Gaussian Splatting-based methods [181], [182], [183]. Additionally, the dense depth map can also be employed to learn a dense SDF field during the training process, allowing for the alignment of the geometry of NeRF or 3D Gaussians with the SDF field. This alignment can improve the geometry of the reconstructed results and facilitate the synthesis of high-quality novel views, particularly for perspectives that are distant from the training views. Some existing works [184], [185], [186], [187], [188] depend on lidar point clouds or RGB-D images as geometric priors. However, both lidar point clouds and RGB-D images can be costly to acquire, requiring additional sensors. High-quality depth estimations can serve as a substitute for these methods, unleashing the potential of novel view synthesis techniques.

### 8.3 Video World Models

As the popularity of diffusion models grows, video diffusion models have been proposed to generalize the image synthesis pipeline to video generation. With the success of SORA [189] and other video foundation models [190], [191], there exist several attempts [192], [193], [194] at exploring video models' capabilities as world models. Having a foundational model for depth estimation, preferably on videos, would significantly bridge the gap between image-only generation models' understanding of 3D space and motion. As the ability to predict depth would indicate, the model at least possesses the capability of distinguishing the size and placement pattern of everyday objects, scenes, and people. Having depth cues for video generation models could potentially serve as a breaking point for further boosting current world models' ability to understand everyday scenes and might even stimulate the underlying generalization ability even further, leading to more spatially and temporally consistent generation and future prediction results.

### 8.4 Robotics and Autonomous Driving

In robotics and autonomous driving, accurate and reliable depth perception plays a pivotal role in tasks such as navigation, obstacle detection, and collision avoidance. Traditional solutions often rely on LiDAR or stereo camera systems, both of which come with increased hardware costs and complexity. Depth foundation models learned from large-scale datasets have the potential to deliver high-quality depth estimates from a single camera, making them particularly attractive for cost-sensitive real-world applications. Recent methods [94] demonstrate that monocular depth estimation can be integrated into robotic perception pipelines, serving as a complement or even a substitute for more expensive sensors. For instance, some works [76], [77] adopt monocular depth estimation to enhance SLAM or visual odometry frameworks, showing improvements in localization and mapping under challenging lighting or weather conditions. Similarly, in autonomous driving, depth estimates can facilitate large-scale 3D reconstructions for building realistic simulation environments, supporting algorithm development and testing [184]. Furthermore, by leveraging depth priors learned from diverse scenes, these models exhibit promising generalization capabilities, potentially enabling robust domain adaptation across varied environments—from urban streets to off-road terrains—thus paving the way for more versatile and scalable robotic and self-driving solutions.

## 9 FUTURE WORK

As discussed in the previous sections, there exist several fundamental problems to be solved before we reach a general-purpose depth foundation model, namely, data and consistency.

**Data.** For all of the discussed depth estimation tasks, including monocular image, stereo image, multi-view image, and monocular video depth estimation, the lack of accurate, large-scale, high-quality, and high-variability data is currently the main concern for constructing and training a depth foundation model. Due to the unique nature of

the depth estimation task, current approaches to acquiring data usually fall into two categories: depth sensor or synthetic rendering. For depth sensors, the main approach is to utilize LiDARs or ultrasonic devices. However, the acquired ground truth depth maps are usually incomplete or noisy due to the sensitive nature of the depth sensing devices. For synthetic data generation, there exist several attempts at curating high-quality hand-crafted, large static or dynamic 3D scenes by artists. However, these data are naturally limited to a small scale due to the amount of work required. Future works should focus on either utilizing self-supervision techniques to better transfer the knowledge of vast image and video data to the task of depth estimation, or developing a better approach for simulation and generation, providing artist-quality synthetic rendering and depth pairs to boost generalization ability.

**Consistency.** This includes both spatial and temporal consistency. For the task of monocular image depth estimation, current methods typically fall short when merging depth estimation results together from different timestamps and viewports of the same scene. For video depth estimation, although a vast amount of work has investigated the issue of temporal consistency throughout the target video, they still fail to produce accurate and consistent results when given multiple viewports of the same 3D scene or trying to unproject and merge the prediction results [94]. Notably, multi-view video reconstruction methods [195], [196], [197] have proved the existence of the dynamic and multi-view inductive bias of the 4D world by providing accurate reconstruction from only image-based optimization objects. Future work should focus on exploring the intrinsic 3D or dynamic inductive bias present in the dynamic 3D world, further mitigating the problem of spatial and temporal inconsistency.

## 10 CONCLUSION

Since 2022, the advancements in foundation models within the natural language processing domain, along with the emergence of scaling laws, have led to a significant increase in the development of foundation models in the field of computer vision. In recent years, numerous foundation models have been introduced for depth estimation tasks, and new models continue to emerge at a rapid pace, making it challenging for practitioners to stay updated with the latest developments.

In this timely paper, we present a comprehensive survey of foundation models for depth estimation tasks, covering their background, development, and the latest advancements. We also address the valuable problems faced by existing depth foundation models and their downstream applications. We aim for this paper to serve as a valuable guide for practitioners and researchers interested in depth estimation foundation models.

Finally, there remain numerous challenges and opportunities in the realm of depth foundation models. We believe that as foundation models continue to evolve and depth estimation tasks advance, we will witness an increasing number of sophisticated and practical applications in the future.

## 11 DECLARATION

**Funding and acknowledgments.** This work was partially supported by NSFC (No. 62172364, No. U24B20154, No. 62402427), Zhejiang Provincial Natural Science Foundation of China (No. LR25F020003), and Information Technology Center and State Key Lab of CAD&CG, Zhejiang University.

**Competing interests.** This survey offers an analysis of recent vision-based depth estimation research and its trend towards depth foundation models, and does not introduce new datasets or materials, nor involve any competing interests.

**Author contributions.** Zhen Xu and Hongyu Zhou were responsible for the overall writing of the manuscript. Sida Peng, Yiyi Liao, Yue Wang, Ruizhen Hu, Xiaowei Zhou, and Hujun Bao provided critical supervision and guidance throughout the project, shaping its framework, refining technical discussions, and ensuring clarity and coherence, providing valuable oversight and feedback during drafting and revision. The remaining co-authors supported the work by evaluating key publications and charting the evolution timeline of depth estimation architectures across the monocular, stereo, multi-view, and monocular video depth estimation tasks to deliver a thorough survey.

## REFERENCES

- [1] J. Geyer, Y. Kassahun, M. Mahmudi, X. Ricou, R. Durgesh, A. S. Chung, L. Hauswald, V. H. Pham, M. Mühlegg, S. Dorn *et al.*, "A2d2: Audi autonomous driving dataset," *arXiv preprint arXiv:2004.06320*, 2020.
- [2] B. Wilson, W. Qi, T. Agarwal, J. Lambert, J. Singh, S. Khandelwal, B. Pan, R. Kumar, A. Hartnett, J. K. Pontes, D. Ramanan, P. Carr, and J. Hays, "Argoverse 2: Next generation datasets for self-driving perception and forecasting," in *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks (NeurIPS Datasets and Benchmarks 2021)*, 2021.
- [3] J. Lambert and J. Hays, "Trust, but verify: Cross-modality fusion for hd map change detection," in *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks (NeurIPS Datasets and Benchmarks 2021)*, 2021.
- [4] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.
- [5] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3d packing for self-supervised monocular depth estimation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [6] J. Cho, D. Min, Y. Kim, and K. Sohn, "Diml/cvl rgb-d dataset: 2m rgb-d images of natural indoor and outdoor scenes," *arXiv preprint arXiv:2110.11590*, 2021.
- [7] I. Vasiljevic, N. Kolkin, S. Zhang, R. Luo, H. Wang, F. Z. Dai, A. F. Daniele, M. Mostajabi, S. Basart, M. R. Walter *et al.*, "Diode: A dense indoor and outdoor depth dataset," *arXiv preprint arXiv:1908.00463*, 2019.
- [8] M. Gehrig, W. Aarents, D. Gehrig, and D. Scaramuzza, "Dsec: A stereo event camera dataset for driving scenarios," *IEEE Robotics and Automation Letters*, 2021.
- [9] S. K. Ramakrishnan, A. Gokaslan, E. Wijmans, O. Maksymets, A. Clegg, J. M. Turner, E. Undersander, W. Galuba, A. Westbury, A. X. Chang, M. Savva, Y. Zhao, and D. Batra, "Habitat-matterport 3d dataset (HM3d): 1000 large-scale 3d environments for embodied AI," in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. [Online]. Available: <https://arxiv.org/abs/2109.08238>
- [10] T. Koch, L. Liebel, F. Fraundorfer, and M. Korner, "Evaluation of cnn-based single-image depth estimation methods," in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2018, pp. 0–0.

- [11] J. Houston, G. Zuidhof, L. Bergamini, Y. Ye, L. Chen, A. Jain, S. Omari, V. Iglovikov, and P. Ondruska, "One thousand and one hours: Self-driving motion prediction dataset," in *Conference on Robot Learning*. PMLR, 2021, pp. 409–418.
- [12] M. L. Antequera, P. Gargallo, M. Hofinger, S. R. Buló, Y. Kuang, and P. Kotschieder, "Mapiillary planet-scale depth dataset," in *European Conference on Computer Vision*. Springer, 2020, pp. 589–604.
- [13] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nusenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 621–11 631.
- [14] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from rgb-d images," in *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7–13, 2012, Proceedings, Part V 12*. Springer, 2012, pp. 746–760.
- [15] P. Xiao, Z. Shao, S. Hao, Z. Zhang, X. Chai, J. Jiao, Z. Li, J. Wu, K. Sun, K. Jiang *et al.*, "Pandaset: Advanced sensor suite dataset for autonomous driving," in *2021 IEEE international intelligent transportation systems conference (ITSC)*. IEEE, 2021, pp. 3095–3101.
- [16] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma, A. Clarkson, M. Yan, B. Budge, Y. Yan, X. Pan, J. Yon, Y. Zou, K. Leon, N. Carter, J. Briales, T. Gillingham, E. Mueggler, L. Pesqueira, M. Savva, D. Batra, H. M. Strasdat, R. D. Nardi, M. Goesele, S. Lovegrove, and R. Newcombe, "The Replica dataset: A digital replica of indoor spaces," *arXiv preprint arXiv:1906.05797*, 2019.
- [17] A. R. Zamir, A. Sax, W. Shen, L. J. Guibas, J. Malik, and S. Savarese, "Taskonomy: Disentangling task transfer learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3712–3722.
- [18] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 3354–3361.
- [19] G. Yang, X. Song, C. Huang, Z. Deng, J. Shi, and B. Zhou, "Drivingstereo: A large-scale dataset for stereo matching in autonomous driving scenarios," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 899–908.
- [20] W. Bao, W. Wang, Y. Xu, Y. Guo, S. Hong, and X. Zhang, "Instereo2k: a large real dataset for stereo matching in indoor scenes," *Science China Information Sciences*, vol. 63, pp. 1–11, 2020.
- [21] M.-F. Chang, J. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan *et al.*, "Argoverse: 3d tracking and forecasting with rich maps," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 8748–8757.
- [22] T. Schops, J. L. Schonberger, S. Galliani, T. Sattler, K. Schindler, M. Pollefeys, and A. Geiger, "A multi-view stereo benchmark with high-resolution images and multi-camera videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3260–3269.
- [23] P. Sun, H. Kretschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine *et al.*, "Scalability in perception for autonomous driving: Waymo open dataset," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2446–2454.
- [24] Z. Bauer, F. Gomez-Donoso, E. Cruz, S. Orts-Escolano, and M. Cazorla, "Uasol, a large-scale high-resolution outdoor stereo dataset," *Scientific Data*, vol. 6, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:201654192>
- [25] H. Aanæs, R. R. Jensen, G. Vogiatzis, E. Tola, and A. B. Dahl, "Large-scale data for multiple-view stereopsis," *International Journal of Computer Vision*, vol. 120, pp. 153–168, 2016.
- [26] Y. Yao, Z. Luo, S. Li, J. Zhang, Y. Ren, L. Zhou, T. Fang, and L. Quan, "Blendedmvs: A large-scale dataset for generalized multi-view stereo networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 1790–1799.
- [27] A. Knapitsch, J. Park, Q.-Y. Zhou, and V. Koltun, "Tanks and temples: Benchmarking large-scale scene reconstruction," *ACM Transactions on Graphics (ToG)*, vol. 36, no. 4, pp. 1–13, 2017.
- [28] H. Xia, Y. Fu, S. Liu, and X. Wang, "Rgbd objects in the wild: scaling real-world 3d object learning from rgb-d videos," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 22 378–22 389.
- [29] X. Yu, M. Xu, Y. Zhang, H. Liu, C. Ye, Y. Wu, Z. Yan, C. Zhu, Z. Xiong, T. Liang, G. Chen, S. Cui, and X. Han, "Mvimgnet: A large-scale dataset of multi-view images," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9150–9161, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:257482733>
- [30] G. Baruch, Z. Chen, A. Dehghan, T. Dimry, Y. Feigin, P. Fu, T. Gebauer, B. Joffe, D. Kurz, A. Schwartz *et al.*, "Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data," *arXiv preprint arXiv:2111.08897*, 2021.
- [31] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3d: Learning from rgb-d data in indoor environments," *arXiv preprint arXiv:1709.06158*, 2017.
- [32] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "ScanNet: Richly-annotated 3d reconstructions of indoor scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5828–5839.
- [33] C. Yeshwanth, Y.-C. Liu, M. Nießner, and A. Dai, "ScanNet++: A high-fidelity dataset of 3d indoor scenes," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 12–22.
- [34] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, "A naturalistic open source movie for optical flow evaluation," in *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7–13, 2012, Proceedings, Part VI 12*. Springer, 2012, pp. 611–625.
- [35] N. Mayer, E. Ilg, P. Hausser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox, "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4040–4048.
- [36] J. Li, P. Wang, P. Xiong, T. Cai, Z. Yan, L. Yang, J. Liu, H. Fan, and S. Liu, "Practical stereo matching via cascaded recurrent network with adaptive correlation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 263–16 272.
- [37] J. Tremblay, T. To, and S. Birchfield, "Falling things: A synthetic dataset for 3d object detection and pose estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 2038–2041.
- [38] B. Wen, M. Trepte, J. Aribido, J. Kautz, O. Gallo, and S. Birchfield, "Foundationstereo: Zero-shot stereo matching," *arXiv preprint arXiv:2501.09898*, 2025.
- [39] F. Tosi, Y. Liao, C. Schmitt, and A. Geiger, "Smd-nets: Stereo mixture density networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 8942–8952.
- [40] L. Mehl, J. Schmalfluss, A. Jahedi, Y. Nalivayko, and A. Bruhn, "Spring: A high-resolution high-detail dataset and benchmark for scene flow, optical flow and stereo," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4981–4991.
- [41] W. Wang, D. Zhu, X. Wang, Y. Hu, Y. Qiu, C. Wang, Y. Hu, A. Kapoor, and S. Scherer, "Tartanair: A dataset to push the limits of visual slam," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 4909–4916.
- [42] Y. Cabon, N. Murray, and M. Humenberger, "Virtual kitti 2," 2020.
- [43] Y. Li, L. Jiang, L. Xu, Y. Xiangli, Z. Wang, D. Lin, and B. Dai, "Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond," *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 3182–3192, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:263135139>
- [44] P.-H. Huang, K. Matzen, J. Kopf, N. Ahuja, and J.-B. Huang, "Deepmvs: Learning multi-view stereopsis," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2821–2830.
- [45] S. Niklaus, L. Mai, J. Yang, and F. Liu, "3d ken burns effect from a single image," *ACM Transactions on Graphics (ToG)*, vol. 38, no. 6, pp. 1–15, 2019.
- [46] N. Karaev, I. Rocco, B. Graham, N. Neverova, A. Vedaldi, and C. Rupprecht, "Dynamicstereo: Consistent dynamic depth from

- stereo videos," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13 229–13 239.
- [47] T. Wu, J. Zhang, X. Fu, Y. Wang, J. Ren, L. Pan, W. Wu, L. Yang, J. Wang, C. Qian *et al.*, "Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 803–814.
- [48] Q. Wang, S. Zheng, Q. Yan, F. Deng, K. Zhao, and X. Chu, "Irs: A large naturalistic indoor robotics stereo dataset to train deep models for disparity and surface normal estimation," in *2021 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2021, pp. 1–6.
- [49] Y. Zheng, A. W. Harley, B. Shen, G. Wetzstein, and L. J. Guibas, "Pointodysey: A large-scale synthetic dataset for long-term point tracking," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 19 855–19 865.
- [50] M. J. Black, P. Patel, J. Tesch, and J. Yang, "Bedlam: A synthetic dataset of bodies exhibiting detailed lifelike animated motion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 8726–8737.
- [51] A. Graves and A. Graves, "Long short-term memory," *Supervised sequence labelling with recurrent neural networks*, pp. 37–45, 2012.
- [52] F. Almeida and G. Xexéo, "Word embeddings: A survey," *arXiv preprint arXiv:1901.09069*, 2019.
- [53] A. Waswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *NIPS*, 2017.
- [54] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 2019, pp. 4171–4186.
- [55] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [56] A. Saxena, S. Chung, and A. Ng, "Learning depth from single monocular images," *Advances in neural information processing systems*, vol. 18, 2005.
- [57] D. Hoiem, A. A. Efros, and M. Hebert, "Recovering surface layout from an image," *International Journal of Computer Vision*, vol. 75, pp. 151–172, 2007.
- [58] C. Liu, J. Yuen, A. Torralba, J. Sivic, and W. T. Freeman, "Sift flow: Dense correspondence across different scenes," in *Computer Vision—ECCV 2008: 10th European Conference on Computer Vision, Marseille, France, October 12–18, 2008, Proceedings, Part III 10*. Springer, 2008, pp. 28–42.
- [59] A. Saxena, M. Sun, and A. Y. Ng, "Make3d: Learning 3d scene structure from a single still image," *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 5, pp. 824–840, 2008.
- [60] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," *Advances in neural information processing systems*, vol. 27, 2014.
- [61] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2650–2658.
- [62] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, "Deeper depth prediction with fully convolutional residual networks," in *2016 Fourth international conference on 3D vision (3DV)*. IEEE, 2016, pp. 239–248.
- [63] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 12 179–12 188.
- [64] Y. Ji, Z. Chen, E. Xie, L. Hong, X. Liu, Z. Liu, T. Lu, Z. Li, and P. Luo, "Ddp: Diffusion model for dense visual prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 21 741–21 752.
- [65] S. Saxena, C. Herrmann, J. Hur, A. Kar, M. Norouzi, D. Sun, and D. J. Fleet, "The surprising effectiveness of diffusion models for optical flow and monocular depth estimation," *Advances in Neural Information Processing Systems*, vol. 36, pp. 39 443–39 469, 2023.
- [66] B. Ke, A. Obukhov, S. Huang, N. Metzger, R. C. Daudt, and K. Schindler, "Repurposing diffusion-based image generators for monocular depth estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9492–9502.
- [67] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun, "Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 3, pp. 1623–1637, 2020.
- [68] W. Yin, X. Wang, C. Shen, Y. Liu, Z. Tian, S. Xu, C. Sun, and D. Renyin, "Diversedepth: Affine-invariant depth prediction using diverse data," *arXiv preprint arXiv:2002.00569*, 2020.
- [69] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, "Deep ordinal regression network for monocular depth estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2002–2011.
- [70] J.-H. Lee and C.-S. Kim, "Monocular depth estimation using relative depth maps," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 9729–9738.
- [71] S. F. Bhat, I. Alhashim, and P. Wonka, "Adabins: Depth estimation using adaptive bins," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 4009–4018.
- [72] —, "Localbins: Improving depth estimation by learning local distributions," in *European Conference on Computer Vision*. Springer, 2022, pp. 480–496.
- [73] S. Zhang, L. Yang, M. B. Mi, X. Zheng, and A. Yao, "Improving deep regression with ordinal entropy," *arXiv preprint arXiv:2301.08915*, 2023.
- [74] S. Shao, Z. Pei, X. Wu, Z. Liu, W. Chen, and Z. Li, "Iebins: Iterative elastic bins for monocular depth estimation," *Advances in Neural Information Processing Systems*, vol. 36, pp. 53 025–53 037, 2023.
- [75] W. Yin, J. Zhang, O. Wang, S. Niklaus, L. Mai, S. Chen, and C. Shen, "Learning to recover 3d scene shape from a single image," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 204–213.
- [76] W. Yin, C. Zhang, H. Chen, Z. Cai, G. Yu, K. Wang, X. Chen, and C. Shen, "Metric3d: Towards zero-shot metric 3d prediction from a single image," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 9043–9053.
- [77] M. Hu, W. Yin, C. Zhang, Z. Cai, X. Long, H. Chen, K. Wang, G. Yu, C. Shen, and S. Shen, "Metric3d v2: A versatile monocular geometric foundation model for zero-shot metric depth and surface normal estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [78] L. Piccinelli, Y.-H. Yang, C. Sakaridis, M. Segu, S. Li, L. Van Gool, and F. Yu, "Unidepth: Universal monocular metric depth estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 10 106–10 116.
- [79] A. Bochkovskii, A. Delaunoy, H. Germain, M. Santos, Y. Zhou, S. R. Richter, and V. Koltun, "Depth pro: Sharp monocular metric depth in less than a second," *arXiv preprint arXiv:2410.02073*, 2024.
- [80] Z. Li and N. Snavely, "Megadepth: Learning single-view depth prediction from internet photos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2041–2050.
- [81] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, "Digging into self-supervised monocular depth estimation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 3828–3838.
- [82] L. Yang, B. Kang, Z. Huang, X. Xu, J. Feng, and H. Zhao, "Depth anything: Unleashing the power of large-scale unlabeled data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 10 371–10 381.
- [83] X. Qi, R. Liao, Z. Liu, R. Urtasun, and J. Jia, "Geonet: Geometric neural network for joint depth and surface normal estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 283–291.
- [84] D. Xu, W. Ouyang, X. Wang, and N. Sebe, "Pad-net: Multi-tasks guided prediction-and-distillation network for simultaneous depth estimation and scene parsing," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 675–684.
- [85] X. Fu, W. Yin, M. Hu, K. Wang, Y. Ma, P. Tan, S. Shen, D. Lin, and X. Long, "Geowizard: Unleashing the diffusion priors for 3d geometry estimation from a single image," in *European Conference on Computer Vision*. Springer, 2024, pp. 241–258.
- [86] A. Eftekhar, A. Sax, J. Malik, and A. Zamir, "Omniata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 786–10 796.

- [87] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations*, 2020.
- [88] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [89] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10684–10695.
- [90] R. Wang, S. Xu, C. Dai, J. Xiang, Y. Deng, X. Tong, and J. Yang, “Moge: Unlocking accurate monocular geometry estimation for open-domain images with optimal training supervision,” 2024. [Online]. Available: <https://arxiv.org/abs/2410.19115>
- [91] R. Wang, S. Xu, Y. Dong, Y. Deng, J. Xiang, Z. Lv, G. Sun, X. Tong, and J. Yang, “Moge-2: Accurate monocular geometry with metric scale and sharp details,” 2025. [Online]. Available: <https://arxiv.org/abs/2507.02546>
- [92] S. F. Bhat, R. Birkel, D. Wofk, P. Wonka, and M. Müller, “Zoedepth: Zero-shot transfer by combining relative and metric depth,” 2023. [Online]. Available: <https://arxiv.org/abs/2302.12288>
- [93] M. Viola, K. Qu, N. Metzger, B. Ke, A. Becker, K. Schindler, and A. Obukhov, “Marigold-dc: Zero-shot monocular depth completion with guided diffusion,” 2024.
- [94] H. Lin, S. Peng, J. Chen, S. Peng, J. Sun, M. Liu, H. Bao, J. Feng, X. Zhou, and B. Kang, “Prompting depth anything for 4k resolution accurate metric depth estimation,” in *CVPR*, 2025.
- [95] H. Yang, D. Huang, W. Yin, C. Shen, H. Liu, X. He, B. Lin, W. Ouyang, and T. He, “Depth any video with scalable synthetic data,” *arXiv preprint arXiv:2410.10815*, 2024.
- [96] A. Kendall, H. Martirosyan, S. Dasgupta, P. Henry, R. Kennedy, A. Bachrach, and A. Bry, “End-to-end learning of geometry and context for deep stereo regression,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 66–75.
- [97] J.-R. Chang and Y.-S. Chen, “Pyramid stereo matching network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5410–5418.
- [98] X. Guo, K. Yang, W. Yang, X. Wang, and H. Li, “Group-wise correlation stereo network,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 3273–3282.
- [99] F. Zhang, V. Prisacariu, R. Yang, and P. H. Torr, “Ga-net: Guided aggregation net for end-to-end stereo matching,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 185–194.
- [100] Z. Shen, Y. Dai, X. Song, Z. Rao, D. Zhou, and L. Zhang, “Pcwnet: Pyramid combination and warping cost volume for stereo matching,” in *European conference on computer vision*. Springer, 2022, pp. 280–297.
- [101] X. Cheng, Y. Zhong, M. Harandi, Y. Dai, X. Chang, H. Li, T. Drummond, and Z. Ge, “Hierarchical neural architecture search for deep stereo matching,” *Advances in neural information processing systems*, vol. 33, pp. 22158–22169, 2020.
- [102] A. Badki, A. Troccoli, K. Kim, J. Kautz, P. Sen, and O. Gallo, “Bi3d: Stereo depth estimation via binary classifications,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1600–1608.
- [103] H. Xu and J. Zhang, “Aanet: Adaptive aggregation network for efficient stereo matching,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 1959–1968.
- [104] M. Yang, F. Wu, and W. Li, “Waveletstereo: Learning wavelet coefficients of disparity map in stereo matching,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 12885–12894.
- [105] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [106] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- [107] Z. Li, X. Liu, N. Drenkow, A. Ding, F. X. Creighton, R. H. Taylor, and M. Unberath, “Revisiting stereo depth estimation from a sequence-to-sequence perspective with transformers,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 6197–6206.
- [108] W. Guo, Z. Li, Y. Yang, Z. Wang, R. H. Taylor, M. Unberath, A. Yuille, and Y. Li, “Context-enhanced stereo transformer,” in *European Conference on Computer Vision*. Springer, 2022, pp. 263–279.
- [109] Z. Liu, Y. Li, and M. Okutomi, “Global occlusion-aware transformer for robust stereo matching,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 3535–3544.
- [110] H. Xu, J. Zhang, J. Cai, H. Rezatofighi, F. Yu, D. Tao, and A. Geiger, “Unifying flow, stereo and depth estimation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [111] Q. Su and S. Ji, “Chitransformer: Towards reliable stereo from cues,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1939–1949.
- [112] Weinzaepfel, Philippe and Leroy, Vincent and Lucas, Thomas and Bréjier, Romain and Cabon, Yohann and Arora, Vaibhav and Antsfeld, Leonid and Chidlovskii, Boris and Csurka, Gabriela and Revaud Jérôme, “CroCo: Self-Supervised Pre-training for 3D Vision Tasks by Cross-View Completion,” in *NeurIPS*, 2022.
- [113] J. Lou, W. Liu, Z. Chen, F. Liu, and J. Cheng, “Elfnet: Evidential local-global fusion for stereo matching,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 17784–17793.
- [114] L. R. Medsker, L. Jain *et al.*, “Recurrent neural networks,” *Design and Applications*, vol. 5, no. 64–67, p. 2, 2001.
- [115] L. Lipson, Z. Teed, and J. Deng, “Raft-stereo: Multilevel recurrent field transforms for stereo matching,” in *2021 International Conference on 3D Vision (3DV)*. IEEE, 2021, pp. 218–227.
- [116] G. Xu, X. Wang, X. Ding, and X. Yang, “Iterative geometry encoding volume for stereo matching,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 21919–21928.
- [117] G. Xu, X. Wang, Z. Zhang, J. Cheng, C. Liao, and X. Yang, “Ige++: iterative multi-range geometry encoding volumes for stereo matching,” *arXiv preprint arXiv:2409.00638*, 2024.
- [118] L. Bartolomei, F. Tosi, M. Poggi, and S. Mattocchia, “Stereo anywhere: Robust zero-shot deep stereo matching even where either stereo or mono fail,” *arXiv preprint arXiv:2412.04472*, 2024.
- [119] J. Cheng, L. Liu, G. Xu, X. Wang, Z. Zhang, Y. Deng, J. Zang, Y. Chen, Z. Cai, and X. Yang, “Monster: Marry monodepth to stereo unleashes power,” *arXiv preprint arXiv:2501.08643*, 2025.
- [120] H. Jiang, Z. Lou, L. Ding, R. Xu, M. Tan, W. Jiang, and R. Huang, “Defom-stereo: Depth foundation model based stereo matching,” *arXiv preprint arXiv:2501.09466*, 2025.
- [121] X. Wang, H. Yang, G. Xu, J. Cheng, M. Lin, Y. Deng, J. Zang, Y. Chen, and X. Yang, “StereoGen: High-quality stereo image generation from a single image,” *arXiv preprint arXiv:2501.08654*, 2025.
- [122] X. Guo, C. Zhang, Y. Zhang, D. Nie, R. Wang, W. Zheng, M. Poggi, and L. Chen, “Stereo anything: Unifying stereo matching with large-scale mixed data,” *arXiv preprint arXiv:2411.14053*, 2024.
- [123] J. Watson, O. M. Aodha, D. Turmukhambetov, G. J. Brostow, and M. Firman, “Learning stereo from single images,” in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 2020, pp. 722–740.
- [124] M. Bleyer, C. Rhemann, and C. Rother, “Patchmatch stereo-stereo matching with slanted support windows,” in *Bmvc*, vol. 11, no. 2011, 2011, pp. 1–11.
- [125] J. L. Schönberger, E. Zheng, J.-M. Frahm, and M. Pollefeys, “Pixelwise view selection for unstructured multi-view stereo,” in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part III 14*. Springer, 2016, pp. 501–518.
- [126] J. L. Schönberger and J.-M. Frahm, “Structure-from-motion revisited,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4104–4113.
- [127] B. Triggs, P. F. McLauchlan, R. I. Hartley, and A. W. Fitzgibbon, “Bundle adjustment—a modern synthesis,” in *Vision Algorithms: Theory and Practice: International Workshop on Vision Algorithms Corfu, Greece, September 21–22, 1999 Proceedings*. Springer, 2000, pp. 298–372.
- [128] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [129] M. Montemerlo, S. Thrun, D. Koller, B. Wegbreit *et al.*, “Fastslam: A factored solution to the simultaneous localization and mapping problem,” *Aaai/iaai*, vol. 593598, pp. 593–598, 2002.

- [130] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, "Monoslam: Real-time single camera slam," *IEEE transactions on pattern analysis and machine intelligence*, vol. 29, no. 6, pp. 1052–1067, 2007.
- [131] R. Y. Tsai, R. K. Lenz *et al.*, "A new technique for fully autonomous and efficient 3 d robotics hand/eye calibration," *IEEE Transactions on robotics and automation*, vol. 5, no. 3, pp. 345–358, 1989.
- [132] K. Daniilidis, "Hand-eye calibration using dual quaternions," *The International Journal of Robotics Research*, vol. 18, no. 3, pp. 286–298, 1999.
- [133] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman, "Patchmatch: A randomized correspondence algorithm for structural image editing," *ACM Trans. Graph.*, vol. 28, no. 3, p. 24, 2009.
- [134] K. Wang and S. Shen, "Mvdepthnet: Real-time multiview depth estimation neural network," in *2018 International conference on 3d vision (3DV)*. IEEE, 2018, pp. 248–257.
- [135] Z. Yang, Z. Ren, Q. Shan, and Q. Huang, "Mvs2d: Efficient multi-view stereo via attention-driven 2d convolutions," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 8574–8584.
- [136] A. Kar, C. Häne, and J. Malik, "Learning a multi-view stereo machine," *Advances in neural information processing systems*, vol. 30, 2017.
- [137] M. Ji, J. Gall, H. Zheng, Y. Liu, and L. Fang, "Surfacenet: An end-to-end 3d neural network for multiview stereopsis," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2307–2315.
- [138] Y. Yao, Z. Luo, S. Li, T. Fang, and L. Quan, "Mvsnet: Depth inference for unstructured multi-view stereo," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 767–783.
- [139] S. Cheng, Z. Xu, S. Zhu, Z. Li, L. E. Li, R. Ramamoorthi, and H. Su, "Deep stereo using adaptive thin volume representation with uncertainty awareness," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2524–2534.
- [140] X. Gu, Z. Fan, S. Zhu, Z. Dai, F. Tan, and P. Tan, "Cascade cost volume for high-resolution multi-view stereo and stereo matching," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2495–2504.
- [141] J. Yang, W. Mao, J. M. Alvarez, and M. Liu, "Cost volume pyramid based depth inference for multi-view stereo," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 4877–4886.
- [142] H. Yi, Z. Wei, M. Ding, R. Zhang, Y. Chen, G. Wang, and Y.-W. Tai, "Pyramid multi-view stereo net with self-adaptive view aggregation," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*. Springer, 2020, pp. 766–782.
- [143] F. Wang, S. Galliani, C. Vogel, and M. Pollefeys, "Itermvs: Iterative probability estimation for efficient multi-view stereo," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 8606–8615.
- [144] X. Wang, Z. Zhu, F. Qin, Y. Ye, G. Huang, X. Chi, Y. He, and X. Wang, "Mvster: Epipolar transformer for efficient multi-view stereo," 2022.
- [145] C. Cao, X. Ren, and Y. Fu, "Mvsformer: Multi-view stereo by learning robust image features and temperature-based depth," *Transactions of Machine Learning Research*, 2023.
- [146] X. R. Chenjie Cao and Y. Fu, "Mvsformer++: Revealing the devil in transformer's details for multi-view stereo," in *International Conference on Learning Representations (ICLR)*, 2024.
- [147] Y. Xu, Z. Shi, W. Yifan, S. Peng, C. Yang, Y. Shen, and W. Gordon, "Grm: Large gaussian reconstruction model for efficient 3d reconstruction and generation," *arxiv: 2403.14621*, 2024.
- [148] J. Wang, N. Karaev, C. Rupprecht, and D. Novotny, "Vggsfm: Visual geometry grounded deep structure from motion," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 21 686–21 697.
- [149] J. Wang, M. Chen, N. Karaev, A. Vedaldi, C. Rupprecht, and D. Novotny, "Vggt: Visual geometry grounded transformer," *arXiv preprint arXiv:2503.11651*, 2025.
- [150] V. Leroy, Y. Cabon, and J. Revaud, "Grounding image matching in 3d with mast3r," 2024.
- [151] J. Lu, T. Huang, P. Li, Z. Dou, C. Lin, Z. Cui, Z. Dong, S.-K. Yeung, W. Wang, and Y. Liu, "Align3r: Aligned monocular depth estimation for dynamic videos," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 22 820–22 830.
- [152] Q. Wang, Y. Zhang, A. Holynski, A. A. Efros, and A. Kanazawa, "Continuous 3d perception model with persistent state," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 10 510–10 522.
- [153] H. Wang and L. Agapito, "3d reconstruction with spatial memory," *arXiv preprint arXiv:2408.16061*, 2024.
- [154] B. Duisterhof, L. Zust, P. Weinzaepfel, V. Leroy, Y. Cabon, and J. Revaud, "Mast3r-sfm: a fully-integrated solution for unconstrained structure-from-motion," *arXiv preprint arXiv:2409.19152*, 2024.
- [155] H. Guo, H. Zhu, S. Peng, H. Lin, Y. Yan, T. Xie, W. Wang, X. Zhou, and H. Bao, "Multi-view reconstruction via sfm-guided monocular depth estimation," in *CVPR*, 2025.
- [156] S. Wang, V. Leroy, Y. Cabon, B. Chidlovskii, and J. Revaud, "Dust3r: Geometric 3d vision made easy," in *CVPR*, 2024.
- [157] S. Zhang, J. Wang, Y. Xu, N. Xue, C. Rupprecht, X. Zhou, Y. Shen, and G. Wetzstein, "Flare: Feed-forward geometry, appearance and camera estimation from uncalibrated sparse views," *arXiv preprint arXiv:2502.12138*, 2025.
- [158] J. Yang, A. Sax, K. J. Liang, M. Henaff, H. Tang, A. Cao, J. Chai, F. Meier, and M. Feiszli, "Fast3r: Towards 3d reconstruction of 1000+ images in one forward pass," *arXiv preprint arXiv:2501.13928*, 2025.
- [159] W. Jang, P. Weinzaepfel, V. Leroy, L. Agapito, and J. Revaud, "Pow3r: Empowering unconstrained 3d reconstruction with camera and scene priors," *arXiv preprint arXiv:2503.17316*, 2025.
- [160] Z. Li, R. Tucker, F. Cole, Q. Wang, L. Jin, V. Ye, A. Kanazawa, A. Holynski, and N. Snavely, "Megasam: Accurate, fast and robust structure and motion from casual dynamic videos," in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 10 486–10 496.
- [161] J. Shriram, A. Trevithick, L. Liu, and R. Ramamoorthi, "Realmdreamer: Text-driven 3d scene generation with inpainting and depth diffusion," *International Conference on 3D Vision (3DV)*, 2025.
- [162] J. Zhang, C. Herrmann, J. Hur, V. Jampani, T. Darrell, F. Cole, D. Sun, and M.-H. Yang, "Monst3r: A simple approach for estimating geometry in the presence of motion," *arXiv preprint arxiv:2410.03825*, 2024.
- [163] P. Wang, H. Tan, S. Bi, Y. Xu, F. Luan, K. Sunkavalli, W. Wang, Z. Xu, and K. Zhang, "Pf-lrm: Pose-free large reconstruction model for joint pose and shape prediction," *arXiv preprint arXiv:2311.12024*, 2023.
- [164] H. Zhang, C. Shen, Y. Li, Y. Cao, Y. Liu, and Y. Yan, "Exploiting temporal consistency for real-time video depth estimation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1725–1734.
- [165] X. Luo, J.-B. Huang, R. Szeliski, K. Matzen, and J. Kopf, "Consistent video depth estimation," *ACM Transactions on Graphics (ToG)*, vol. 39, no. 4, pp. 71–1, 2020.
- [166] R. Yasarla, H. Cai, J. Jeong, Y. Shi, R. Garrepalli, and F. Porikli, "Mamo: Leveraging memory and attention for monocular video depth estimation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 8754–8764.
- [167] J. Shao, Y. Yang, H. Zhou, Y. Zhang, Y. Shen, M. Poggi, and Y. Liao, "Learning temporally consistent video depth from video diffusion priors," *arXiv preprint arXiv:2406.01493*, 2024.
- [168] W. Hu, X. Gao, X. Li, S. Zhao, X. Cun, Y. Zhang, L. Quan, and Y. Shan, "Depthrafter: Generating consistent long depth sequences for open-world videos," *arXiv preprint arXiv:2409.02095*, 2024.
- [169] Y. Wang, M. Shi, J. Li, Z. Huang, Z. Cao, J. Zhang, K. Xian, and G. Lin, "Neural video depth stabilizer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 9466–9476.
- [170] S. Chen, H. Guo, S. Zhu, F. Zhang, Z. Huang, J. Feng, and B. Kang, "Video depth anything: Consistent depth estimation for super-long videos," *arXiv preprint arXiv:2501.12375*, 2025.
- [171] Z. Kuang, T. Zhang, K. Zhang, H. Tan, S. Bi, Y. Hu, Z. Xu, M. Hasan, G. Wetzstein, and F. Luan, "Buffer anytime: Zero-shot video depth and normal from image priors," *arXiv preprint arXiv:2411.17249*, 2024.
- [172] S. Galliani, K. Lasinger, and K. Schindler, "Massively parallel multiview stereopsis by surface normal diffusion," in *Proceedings*

- of the *IEEE International Conference on Computer Vision*, 2015, pp. 873–881.
- [173] D. Sörmann, A. Dünser, F. Fraundorfer, and D. Schmalstieg, “Simplerecon: 3d reconstruction without 3d convolutions,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12 153–12 162.
- [174] B. Curless and M. Levoy, “A volumetric method for building complex models from range images,” *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pp. 303–312, 1996.
- [175] J. Wang, A. Kar, and S. Fidler, “Neuris: Neural reconstruction of indoor scenes using normal priors,” in *European Conference on Computer Vision*. Springer, 2022, pp. 647–664.
- [176] Z. Yu, S. Peng, M. Niemeyer, T. Sattler, and A. Geiger, “Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction,” in *Advances in Neural Information Processing Systems*, vol. 35, 2022, pp. 3646–3658.
- [177] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [178] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, “3d gaussian splatting for real-time radiance field rendering,” *ACM Transactions on Graphics*, vol. 42, no. 4, July 2023. [Online]. Available: <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>
- [179] J. Guo, N. Deng, X. Li, Y. Bai, B. Shi, C. Wang, C. Ding, D. Wang, and Y. Li, “Streetsurf: Extending multi-view implicit surface reconstruction to street views,” *arXiv preprint arXiv:2306.04988*, 2023.
- [180] S. Miao, J. Huang, D. Bai, W. Qiu, B. Liu, A. Geiger, and Y. Liao, “Efficient depth-guided urban view synthesis,” in *European Conference on Computer Vision*. Springer, 2024, pp. 90–107.
- [181] Z. Fan, K. Wen, W. Cong, K. Wang, J. Zhang, X. Ding, D. Xu, B. Ivanovic, M. Pavone, G. Pavlakos *et al.*, “Instantsplat: Sparse-view sfm-free gaussian splatting in seconds,” *arXiv preprint arXiv:2403.20309*, 2024.
- [182] H. Zhou, L. Lin, J. Wang, Y. Lu, D. Bai, B. Liu, Y. Wang, A. Geiger, and Y. Liao, “Hugsim: A real-time, photo-realistic and closed-loop simulator for autonomous driving,” *arXiv preprint arXiv:2412.01718*, 2024.
- [183] S. Miao, J. Huang, D. Bai, X. Yan, H. Zhou, Y. Wang, B. Liu, A. Geiger, and Y. Liao, “Evolspat: Efficient volume-based gaussian splatting for urban view synthesis,” *arXiv preprint arXiv:2503.20168*, 2025.
- [184] Y. Yan, H. Lin, C. Zhou, W. Wang, H. Sun, K. Zhan, X. Lang, X. Zhou, and S. Peng, “Street gaussians: Modeling dynamic urban scenes with gaussian splatting,” in *ECCV*, 2024.
- [185] E. Sandström, K. Tateno, M. Oechsle, M. Niemeyer, L. Van Gool, M. R. Oswald, and F. Tombari, “Splat-slam: Globally optimized rgb-only slam with 3d gaussians,” *arXiv preprint arXiv:2405.16544*, 2024.
- [186] Y. Pan, X. Zhong, L. Jin, L. Wiesmann, M. Popović, J. Behley, and C. Stachniss, “Pings: Gaussian splatting meets distance fields within a point-based implicit neural map,” *arXiv preprint arXiv:2502.05752*, 2025.
- [187] H. Matsuki, R. Murai, P. H. Kelly, and A. J. Davison, “Gaussian splatting slam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 18 039–18 048.
- [188] Y. Mao, X. Yu, Z. Zhang, K. Wang, Y. Wang, R. Xiong, and Y. Liao, “Ngel-slam: Neural implicit representation-based global consistent low-latency slam system,” in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 6952–6958.
- [189] OpenAI, “Video generation models as world simulators,” <https://openai.com/index/video-generation-models-as-world-simulators/>, 2024, [Accessed 15-03-2025].
- [190] W. Hong, M. Ding, W. Zheng, X. Liu, and J. Tang, “Cogvideo: Large-scale pretraining for text-to-video generation via transformers,” *arXiv preprint arXiv:2205.15868*, 2022.
- [191] Z. Yang, J. Teng, W. Zheng, M. Ding, S. Huang, J. Xu, Y. Yang, W. Hong, X. Zhang, G. Feng *et al.*, “Cogvideox: Text-to-video diffusion models with an expert transformer,” *arXiv preprint arXiv:2408.06072*, 2024.
- [192] M. Rigtter, T. Gupta, A. Hilmkil, and C. Ma, “Avid: Adapting video diffusion models to world models,” *arXiv preprint arXiv:2410.12822*, 2024.
- [193] Q. Zhang, S. Zhai, M. A. Bautista, K. Miao, A. Toshev, J. Susskind, and J. Gu, “World-consistent video diffusion with explicit 3d modeling,” *arXiv preprint arXiv:2412.01821*, 2024.
- [194] B. Kang, Y. Yue, R. Lu, Z. Lin, Y. Zhao, K. Wang, G. Huang, and J. Feng, “How far is video generation from world model: A physical law perspective,” *arXiv preprint arXiv:2411.02385*, 2024.
- [195] Z. Xu, Y. Xu, Z. Yu, S. Peng, J. Sun, H. Bao, and X. Zhou, “Representing long volumetric video with temporal gaussian hierarchy,” *ACM Transactions on Graphics*, vol. 43, no. 6, November 2024. [Online]. Available: <https://zju3dv.github.io/longvolcap>
- [196] Z. Xu, S. Peng, H. Lin, G. He, J. Sun, Y. Shen, H. Bao, and X. Zhou, “4k4d: Real-time 4d view synthesis at 4k resolution,” in *CVPR*, 2024.
- [197] Z. Xu, T. Xie, S. Peng, H. Lin, Q. Shuai, Z. Yu, G. He, J. Sun, H. Bao, and X. Zhou, “Easyvolcap: Accelerating neural volumetric video research,” *SIGGRAPH Asia 2023 Technical Communications*, 2023.