

NTIRE 2025 Challenge on HR Depth from Images of Specular and Transparent Surfaces

Pierluigi Zama Ramirez Fabio Tosi Luigi Di Stefano Radu Timofte
 Alex Costanzino Matteo Poggi Samuele Salti Stefano Mattoccia Zhe Zhang
 Yang Yang Wu Chen Anlong Ming Mingshuai Zhao Mengying Yu Shida Gao
 Xiangfeng Wang Feng Xue Jun Shi Yong Yang Yong A Yixiang Jin
 Dingzhe Li Aryan Shukla Liam Frija-Altarc Matthew Toews Hui Geng
 Tianjiao Wan Zijian Gao Qisheng Xu Kele Xu Zijian Zang
 Jameer Babu Pinjari Kuldeep Purohit Mykola Lavreniuk Jing Cao Shenyi Li
 Kui Jiang Junjun Jiang Yong Huang

Abstract

This paper reports on the NTIRE 2025 challenge on HR Depth From images of Specular and Transparent surfaces, held in conjunction with the New Trends in Image Restoration and Enhancement (NTIRE) workshop at CVPR 2025. This challenge aims to advance the research on depth estimation, specifically to address two of the main open issues in the field: high-resolution and non-Lambertian surfaces. The challenge proposes two tracks on stereo and single-image depth estimation, attracting about 177 registered participants. In the final testing stage, 4 and 4 participating teams submitted their models and fact sheets for the two tracks.

1. Introduction

Reversing the image formation process to model the 3D structure of the world represents one of the quintessential tasks studied by computer vision. For this purpose, estimating depth from images often lays the foundation of this process, as well as the entry point to higher-level applications such as augmented reality, autonomous or assisted driving, robotics, and more. Recovering this information from images represents a cheaper and more viable alternative to the use of *active* depth sensors – such as Radars, LiDARs, Time-of-Flight (ToF), and others – which are known for

their higher cost and multiple limitations preventing their unconstrained deployment in any environment. Furthermore, the disruptive advent of deep learning in computer vision has made the former strategy more and more preferable over active sensors, also thanks to the recent development of the first *foundational* models for depth estimation and, in general, 3D vision. Although this brought a rapid evolution of the depth estimation models observed in the last decade, this task remains far from being solved in the presence of some particularly challenging conditions. Among the many, we argue that two matters of interest are common to the different approaches devoted to estimating depth from images – and even to active sensors.

The first is a longstanding challenge, common to any computer vision task: spatial resolution. Indeed, although nowadays we have color cameras capable of capturing frames up to dozens of Megapixels (Mpx) whereas active sensors fall far behind, processing high-resolution images poses several challenges, in terms of computational requirements as well as data and training methodologies to deploy deep models capable of exploiting such rich information.

The second consists of the ambiguity that may occur in some images: this may assume different forms, such as the lack of texture or the absence of perspective cues, and affect different kinds of depth estimation strategies. Specifically, the presence of *non-Lambertian* surfaces represents a challenge to any of these approaches, as well as to active sensors – since materials featuring this property often violate the basic assumptions upon which active sensors are built, e.g., with LiDARs beams being refracted or surpassing transparent surfaces. This makes ground-truth depth annotations, often sourced through active sensors, very hard to collect and, therefore, very rare in the training data leveraged by most state-of-the-art image-based depth estimation

*Pierluigi Zama Ramirez (pierluigi.zama@unibo.it), Alex Costanzino, Fabio Tosi, Matteo Poggi, Samuele Salti, Stefano Mattoccia, Luigi Di Stefano and Radu Timofte are the NTIRE 2025 HR Depth from Images of Specular and Transparent Surfaces challenge organizers. The other authors participated in the challenge. Sections B and C contains the authors' team names and affiliations. The NTIRE website: <https://cvlai.net/ntire/2025/>

models, making these latter failing to estimate the distance of a transparent surface in favor of the distance of objects behind it, or the surface of a mirror in place of the depth of the reflected objects. Although these latter examples might not represent real failure cases, since the definition of depth itself becomes ambiguous in such circumstances, we argue it is for some popular applications, for instance when properly perceiving the real depth for transparent objects may be crucial to accomplishing a higher level task, like grasping some glassy objects or navigating in an indoor environment where glass doors may be common.

This NTIRE 2025 Challenge on HR Depth from Images of Specular and Transparent Surfaces aims at pushing the development of state-of-the-art methodologies answering to the aforementioned challenges. Following the previous, successful editions [76, 77], we build our challenge over the Booster dataset [127, 130], a benchmark peculiarly encompassing both high-resolution and non-Lambertian surfaces, thanks to its 12Mpx images and the abundant presence of transparent and reflective objects. Following our tradition, the challenge is organized into two tracks: one devoted to *Stereo* approaches, where depth is measured through triangulation from the *disparity* estimated between pixels into two rectified frames, and the other focusing on single-image frameworks (*Mono*). The challenge attracted up to 177 registered participants. Among them, 5 and 4 teams, respectively, for the monocular and stereo tracks, submitted their models and fact sheets during the final phase. Some adopt the most recent foundational models in the field as off-the-shelf solutions, whereas others develop their own custom frameworks. The outcome of this edition of our challenge is reported and discussed in detail in Section 4.

2. Related Work

Deep Stereo Matching. Ten years ago already, the community started facing stereo depth estimation with deep neural networks [131], becoming the standard approach to this task over the years [72, 74]. At first, two main families of end-to-end models were developed, respectively, 2D [60, 65, 70, 73, 84, 91, 97, 99, 119, 125] and 3D [7, 12, 13, 18, 31, 45, 46, 89, 104, 110, 118, 134] architectures. In the 2020s [103], the advent of new paradigms to deal with dense matching tasks, such as the use of optimization-based frameworks [98] or transformers [29, 57, 64] ignited the development of two new lines of research. The former in particular, starting with RAFT-Stereo [98], rapidly become the most popular approach [39, 52, 105, 115, 117, 133]. This translated into a steady saturation of the most popular benchmarks, from KITTI 2012 [25] and 2015 [67], to ETH3D [87] and Middlebury 2014 [86]. Other approaches try to refine disparity maps for high-resolution predictions [2, 102]. Lately, the first foundational models for stereo depth estimation have appeared [3, 11, 38, 113],

achieving a consistent step forward in terms of zero-shot generalization and robustness in handling non-Lambertian surfaces. Indeed, as we will notice in the remainder of this paper, some of these solutions were successfully deployed on the Booster dataset [130] as well.

Monocular Depth Estimation. Deep learning allowed facing highly ill-posed tasks, such as estimating depth out of a single image [8, 19, 50, 75, 109], thanks to the increasing availability of large-scale, annotated datasets [8, 19, 50, 75, 109] or to the emergence of self-supervised paradigms [26–28, 30, 37, 40, 71, 100, 101, 112, 128, 135, 136] replacing the need for explicit depth annotation with principles of multi-view geometry, for instance by casting the depth estimation process into an image reconstruction problem during training thanks to the availability of either paired stereo images or monocular videos. A major trend emerging in the twenties consists of the development of affine-invariant depth estimation models [78, 80], capable of generalizing beyond the single-dataset domain. MiDaS [80] took this direction first, training a deep network on a mixture of multiple datasets to achieve cross-domain generalization, then followed by DPT [78], and others focusing on recovering the real shapes from the deformed point cloud obtained from monocular depth maps [124] or restoring high-frequency details [56, 68] at a higher resolution. Affine-invariant models recently converged into the first foundational models for single-image depth estimation, such as the Depth Anything series [121, 122], or the newest diffusion-derived frameworks such as Marigold [44], GeoWizard [22], Lotus [33] and others, then extended to deal with video depth estimation [34, 43, 88].

Lately, the ability of single-image depth estimation to effectively handle transparent and reflective surfaces has gained relevance, also thanks to the advent of benchmarks dedicated to this purpose [127]. On this track, some approaches developed an annotation pipeline to obtain reliable pseud-labels for non-Lambertian objects, by using pre-trained monocular depth estimation models jointly with material segmentation masks [17] or diffusion models [106], whereas others employed depth completion approaches [14, 85] to fill the holes in the depth maps occurring in correspondence of transparent surfaces. Furthermore, some of the latest foundational models such as Depth Anything v2 [122] expose surprising effectiveness at perceiving transparent or mirroring surfaces, as shown in the remainder.

Competitions/Challenges on Depth Estimation. The depth estimation task, both from stereo and monocular images, as been the object of several challenges taking place in the previous years, or even concurrently with ours. Among them, the Robust Vision Challenge (ROB) [132] covering both, the Dense Depth for Autonomous Driving challenge (DDAD) [24], the Fast and Accurate Single-Image Depth Estimation on Mobile Devices Challenge (MAI) [35], the

Argoverse Stereo Challenge [49] and the Monocular Depth Estimation Challenge (MDEC) [69, 92–94]. Finally, we recall the previous editions of this challenge [76, 77, 126], part of the NTIRE workshop at CVPR 2023 and 2024 and the TRICKY workshop at ECCV 2024.

NTIRE 2025 Challenges. This challenge is one of the NTIRE 2025¹ Workshop associated challenges on: ambient lighting normalization [108], reflection removal in the wild [120], shadow removal [107], event-based image deblurring [95], image denoising [96], XGC quality assessment [62], UGC video enhancement [83], night photography rendering [20], image super-resolution (x4) [9], real-world face restoration [10], efficient super-resolution [81], HR depth estimation [129], efficient burst HDR and restoration [51], cross-domain few-shot object detection [23], short-form UGC video quality assessment and enhancement [54, 55], text to image generation model quality assessment [32], day and night raindrop removal for dual-focused images [53], video quality assessment for video conferencing [36], low light image enhancement [63], light field super-resolution [111], restore any image model (RAIM) in the wild [58], raw restoration and super-resolution [15] and raw reconstruction from RGB on smartphones [16].

3. NTIRE Challenge on HR Depth from Images of Specular and Transparent Surfaces

We organize the NTIRE 2025 Challenge on HR Depth from Images of Specular and Transparent Surfaces to further push the community toward developing newer, state-of-the-art solutions to properly deal with high-resolution images and non-Lambertian surfaces – such as mirrors and glasses. We now outline the main characteristics of our challenge.

Tracks. This challenge is composed of two tracks: *Stereo*, devoted to methods estimating the disparity between pairs of rectified images, and *Mono*, which instead allows estimating depth from a single input image only.

- **Track 1: Stereo.** This track demands the participants to obtain high-quality, high-resolution dense disparity maps from 12 Mpx stereo frames. The resolution itself represents one of the main challenges, as it is prohibitive for most state-of-the-art existing stereo networks. Furthermore, the presence of non-Lambertian objects violating the common assumptions made in stereo matching makes this track even harder.
- **Track 2: Mono.** In parallel, this track requires to estimate depth out of a single 12Mpx frame. In this case, the inherent ill-posed nature of the problem represents one of the main challenges. Additionally, the presence of objects belonging to the long-tail of training data used for

this task – such as transparent objects and mirrors – further makes it more complex.

Datasets. We build our challenge around the Booster dataset [127, 130], composed of 419 high-resolution balanced and unbalanced stereo pairs, captured in 64 different scenes and respectively distributed into 228 and 191 pairs for training and testing purposes – with 38 and 26 for the two sets respectively. A newer version of the dataset [127] extended it with a second testing split, tailored for evaluating monocular depth estimation approaches over 187 single frames, captured in 21 different environments.

As in the previous editions [76, 77], we use the original 228 training stereo pair as a shared *training split*, common to both tracks. We then select two distinct *validation splits*, by selecting frames with different illuminations from 3 scenes of the stereo and monocular testing splits – i.e., *Microwave*, *Mirror1*, *Pots* for the Stereo track, and *Desk*, *Mirror3*, *Sanitaries* for the Mono track respectively, for a total of 15 validation samples for each track from total 26 and 28 available from the selected scenes. The remaining images in the two original testing splits become the official stereo and mono *testing splits* for this challenge, for a total of 169 and 159 samples.

Evaluation Protocol. For each track, Stereo and Mono respectively, we adopt the official metrics reported on the Booster benchmark [127, 130]. For Stereo track, we compute the percentage of pixels with disparity errors larger than a threshold τ (bad- τ , with $\tau \in [2, 4, 6, 8]$), as well as the Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) in pixel. For Mono track, we compute the percentage of pixels having the maximum between the prediction/ground-truth and ground-truth/prediction ratios lower than a threshold ($\delta < i$, with i being 1.05, 1.15, and 1.25) and the absolute error relative to the ground truth value (Abs Rel.), as well as the mean absolute error (MAE), and Root Mean Squared Error (RMSE). Following the latest edition [77], we compute metrics on three different sets of pixels as in [17]: *ToM* regions – i.e., those belonging to non-Lambertian surfaces – *All* pixels and *Others* – i.e., the difference between *All* and *ToM* sets. To rank submissions and determine the winner, we use bad-2 and $\delta < 1.05$ – respectively for Stereo and Mono tracks – averaged over all pixels, highlighted in **red** in the tables. Specifically, we define two rankings based on performance on *ToM* and *All* regions, respectively². Finally, as most state-of-the-art monocular networks estimate depth up to an unknown pair of scale and shift factors, before computing metrics we recover metric depth from predicted maps \hat{d} as $\alpha\hat{d} + \beta$, with α, β being scale and shift factors. Following [80], α, β are estimated with Least Square Estimation (LSE) regression over the ground truth depth map d :

¹<https://www.cvlai.net/ntire/2025/>

²we will highlight that the two coincide on the Mono track

Team	ToM							All						Other						
	Rank	bad-2	bad-4	bad-6	bad-8	MAE	RMSE	Rank	bad-2	bad-4	bad-6	bad-8	MAE	RMSE	bad-2	bad-4	bad-6	bad-8	MAE	RMSE
SRC-B	#1	42.47	25.39	18.91	14.21	4.66	7.55	#3	25.50	12.27	7.81	5.42	2.34	5.46	22.54	9.19	5.21	3.31	1.81	4.43
Robot01-vRobotit	#2	45.20	28.30	20.47	16.37	5.14	9.39	#4	30.02	13.34	8.59	6.25	2.86	7.23	28.75	11.22	6.82	4.67	2.54	6.61
NJUST-KMG	#3	50.25	27.97	19.97	16.28	5.79	8.77	#2	22.64	9.76	6.29	4.77	2.44	5.68	18.75	6.36	3.71	2.62	1.82	4.45
weouibaguette	#4	52.30	32.57	26.49	23.63	18.54	25.40	#1	19.27	8.74	6.44	5.43	5.06	12.22	13.64	3.51	1.76	1.04	1.27	3.48
CREStereo [baseline]	#5	59.64	47.26	40.27	35.41	24.69	42.28	#5	35.75	23.51	18.98	16.42	12.13	28.46	8.34	13.93	7.91	4.64	2.95	7.59

Table 1. **Stereo Track: Evaluation on the Challenge Test Set.** Predictions evaluated at full resolution (4112×3008) on All pixels and pixels belonging to ToM (Transparent or Mirror) or Other materials. In **gold**, **silver**, and **bronze**, we show first, second, and third-rank approaches, respectively. We rank methods on two metrics, $\delta < 1.05$ computed on either ToM or All pixels.

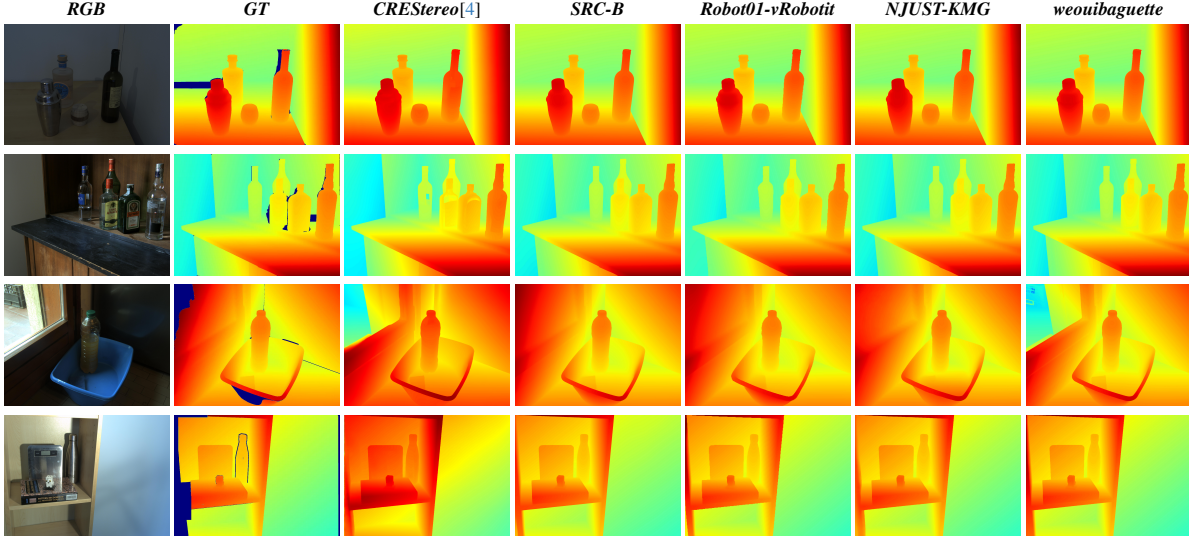


Figure 1. **Qualitative results – Stereo track.** From left to right: RGB reference image, ground-truth disparity, predictions by CREStereo [52], SRC-B, Robot01-vRobotit, NJUST-KMG, and weouibaguette.

$$(\alpha, \beta) = \arg \min_{\alpha, \beta} \sum_p (\alpha \hat{d}(p) + \beta - d(p))^2 \quad (1)$$

where p are the pixel locations having both predictions and ground truth depths available.

4. Challenge Results

For each track, four teams participated in the final evaluation phase, with their outcomes detailed in Sections 4.1 and 4.2. A brief explanation of each approach for both stereo and mono tracks is provided in Section 5.1 and Section 5.2, while the team composition is detailed in Sections B and C.

4.1. Track 1: Stereo

Table 1 reports the results for this first track. At the bottom, we report the baseline – i.e., CREStereo [52]. From left to right, we report bad- τ metrics, MAE, and RMSE metrics for *ToM*, *All*, and *Other* pixels respectively. On the right of the team’s name, we report their overall rank, computed according to bad-2 errors – the most restrictive metric – on *ToM* and *All* regions.

All submitted methods outperformed the baseline on *ToM* and *All* pixels, with **SRC-B** achieving the lowest error rates on *ToM* pixels and **weouibaguette** *All* pixels, while the baseline still performs the best on *Other* pixels.

Interestingly, unlike previous iterations of the challenge, while there is a somewhat clear trend on *ToM* pixels, such as low bad-2 errors correlating with low MAE and RMSE scores, the results *Other* and *All* are consistently mixed up, making it hard to identify a jack-of-all-trades model. Fig. 1 depicts some qualitative results from the stereo testing set. We can appreciate how the submitted methods tend to deal better with some specific challenges, such as the bottles in row 2, which CREStereo falters to infer without discontinuities, and the book in column 4, where the submitted methods produce much smoother results.

4.2. Track 2: Mono

Table 2 shows the results for the second track. At the very bottom, we report the results achieved by the baseline method – i.e., the ZoeDepth [4] model using the weights provided by the authors. From left to right, we report deltas, Abs Rel., MAE, and RMSE metrics for *ToM*, *All*, and *Other* pixels respectively. We report two different rankings, according to the performance observed on $\delta < 1.05$ – the most restrictive metric – computed over *ToM* and *All* pixels.

Unlike the stereo track, all of the submitted methods consistently outperformed the ZoeDepth baseline, with **Lavre-niuk** achieving the best accuracy values on *ToM*, *All* and *Other* pixels. Indeed, conversely to what was observed in

Team	ToM							All							Other						
	Rank	$\delta < 1.05$	$\delta < 1.15$	$\delta < 1.25$	Abs Rel.	MAE	RMSE	Rank	$\delta < 1.05$	$\delta < 1.15$	$\delta < 1.25$	Abs Rel.	MAE	RMSE	$\delta < 1.05$	$\delta < 1.15$	$\delta < 1.25$	Abs Rel.	MAE	RMSE	
Lavreniuk	#1	87.67	99.15	99.84	2.54	2.71	3.54	#1	84.19	97.73	99.50	3.63	3.12	6.02	82.13	96.93	99.27	4.11	3.39	6.94	
colab	#2	86.64	99.60	99.83	2.59	2.58	3.76	#2	82.93	98.53	99.46	3.53	2.97	6.11	80.22	97.88	99.27	4.03	3.33	6.95	
PreRdw	#3	84.58	99.42	99.86	2.70	2.79	3.64	#3	80.47	96.73	98.45	3.99	3.41	6.59	79.10	95.95	98.15	4.41	3.65	7.46	
IPCV	#4	62.66	95.43	99.19	4.60	4.70	6.01	#4	65.61	91.54	97.34	6.37	5.30	9.70	62.70	90.47	97.10	7.00	5.63	10.67	
ZoeDepth [Baseline]	#5	45.21	82.27	93.06	8.04	8.71	9.57	#5	61.31	87.97	94.38	7.60	6.38	10.88	60.23	87.43	93.71	8.34	6.31	12.18	

Table 2. **Mono Track: Evaluation on the Challenge Test Set.** Predictions evaluated at full resolution (4112×3008) on All pixels and pixels belonging to ToM (Transparent or Mirror) or Other materials. In **gold**, **silver**, and **bronze**, we show first, second, and third-rank approaches, respectively. We rank methods on two metrics, $\delta < 1.05$ computed on either ToM or All pixels.

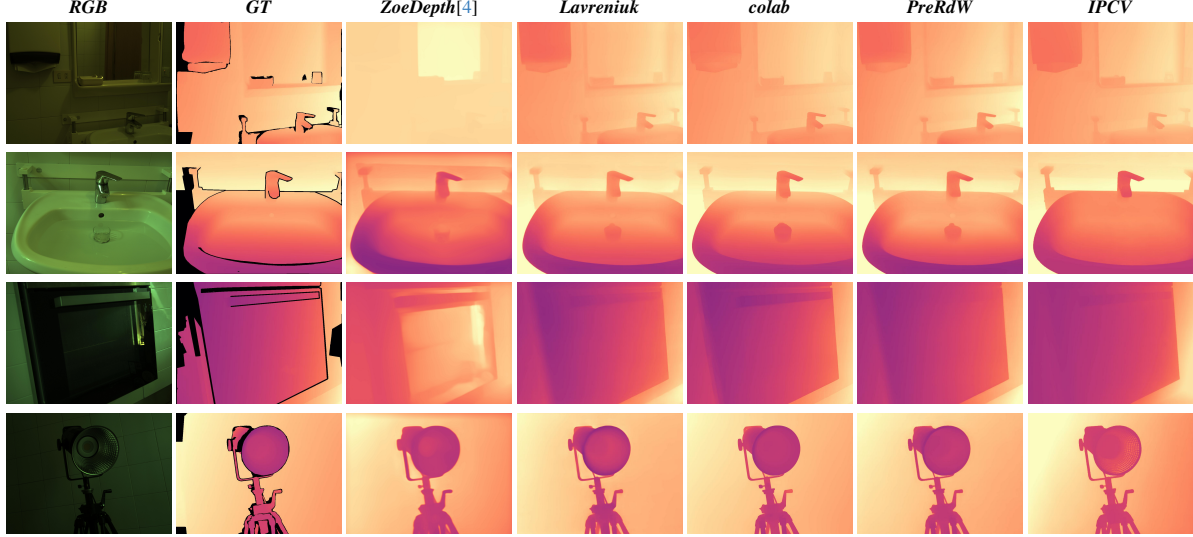


Figure 2. **Qualitative results – Mono track.** From left to right: RGB reference image, ground-truth disparity, predictions by ZoeDepth [4], Lavreniuk, colab, PreRdw, and IPCV.

the stereo track, **Lavreniuk** represents the most versatile method being the top performer on all the pixel categories. For what concerns *ToM* regions, the top #3 methods are able to push the strictest accuracy metric – $\delta < 1.05$ – beyond 85%, with a remarkable 15% improvement with respect to last year, as well as to reduce the Abs Rel. below 3%. The improvements are reflected on *All* and *Other* pixels as well. Despite the minor gain with respect to the baseline, compared to what was observed on *ToM* regions, the improvement is yet consistent.

Fig. 2 shows some qualitative examples from the mono testing set. Similarly to last edition, any of the submitted models can properly handle *ToM* regions, such as for the oven in the third row, while still struggling on mirrors or water surfaces, as in first and second rows.

5. Challenge Methods

5.1. Track 1: Stereo

5.1.1 Baseline - CREStereo [52]

For the first track, we set the state-of-the-art CREStereo architecture [52] as our baseline. This model consists of a hierarchical network employing a recurrent refinement process, designed to update the predicted disparity map in coarse-to-fine manner. This process is implemented based

on an adaptive group correlation layer (AGCL), where an alternate 2D-1D local search strategy with deformable windows is employed for robust matching even in the presence of imperfect rectification. The AGCL module computes correlations between pixels in local search windows, in contrast to what the all-pairs correlation module from RAFT-Stereo [61] does, reducing the computational requirements. To obtain the final predictions, we process images at quarter resolution using the original weights released by the authors, then we upsample predicted disparity maps to the original resolution through bilinear interpolation.

5.1.2 Team 1 - NJUST-KMG

The NJUST-KMG team (CodaLab: *chenyin*) adapts DEFOM-Stereo [38] by integrating Depth Anything V2. Their approach improves CNN encoders with a depth foundation model, introduces a scale update module before the delta update module, and leverages reflective data with multi-scale sampling for training. The feature encoder fuses DPT and CNN feature maps at 1/4 resolution for matching, while the context encoder combines multi-level DPT and CNN features.

Disparity estimation is initialized with Depth Anything V2’s depth, followed by scale correction and detail refinement using pyramid lookup. The system is pre-trained

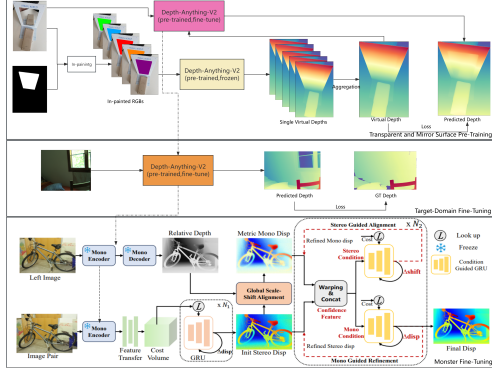


Figure 3. **Network Architecture – Team Robot01-vRobotit.**

on KITTI [67], Middlebury [86], and ETH3D [87], then fine-tuned on Booster [130] and CREStereo [52] datasets. Quarter-resolution inference ensures efficient processing of high-resolution inputs.

5.1.3 Team 2 - Robot01-vRobotit

The Robot01-vRobotit team (CodaLab: *Bupt-chenwu*) introduces an improved MonSter [11] based on Depth Anything V2. Their approach leverages the complementary strengths of monocular depth estimation and stereo matching in a dual-branch architecture, where monocular depth provides global structural information while stereo matching refines pixel-level geometric details.

The training process is divided into two stages: first, the Depth Anything V2-Large model is fine-tuned to improve its accuracy in TOM regions of the Booster dataset; second, the MonSter network is fine-tuned using the improved monocular model from the first stage. During training, extensive data augmentation is applied, including resizing, random cropping, saturation adjustment, color jittering, and spatial scaling. The fine-tuning is performed on NVIDIA RTX 4090 GPUs for 1000 epochs, with a learning rate of 1×10^{-4} and batch size of 12.

5.1.4 Team 3 - SRC-B [Stereo]

The SRC-B team (CodaLab: *pixinsight*) presents “Multi-Scale-Mono-Stereo,” a method that integrates monocular depth network features to improve depth estimation in high-resolution images with non-Lambertian surfaces. The approach adopts a data augmentation strategy similar to AS-Grasp [90], leveraging Blender to generate additional stereo training samples from the AI2THOR [48] dataset. To effectively use the pre-trained monocular model, they adopt the stereo branch structure from MonSter [11], incorporating the pretrained ViT encoder of Depth Anything V2 with frozen parameters. A feature transfer network is introduced to downsample and transform the ViT-extracted features into a multi-scale feature pyramid, which is then

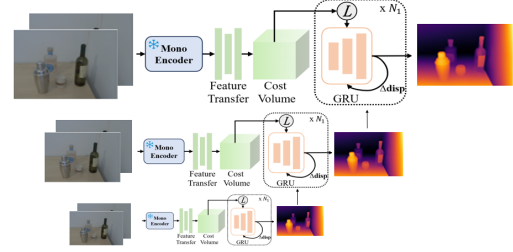


Figure 4. **Network Architecture – Team SRC-B [Stereo].**

concatenated with features extracted by IGEV [116]. To enhance performance on high-resolution images, the team integrates a stacked cascaded architecture during training, enabling the network to adaptively propagate depth information across different scales. The network is implemented in PyTorch and trained on NVIDIA RTX 3090 GPUs, with the Depth Anything V2 module’s weights kept fixed while fine-tuning the stereo module for an additional 100,000 steps.

5.1.5 Team 4 - weouibaguette

The weouibaguette team employs FoundationStereo (FS) [113], finding that the original network outperforms all custom-trained models they developed. Using the Booster training set, they evaluated FS under various parameter configurations and observed, similar to last year’s NTIRE winner (MiMcAlgo), that inference yielded better overall metrics on down-sampled images.

The team experimented with different down-sampling factors (0.5, 0.3, 0.25, and 0.1) and tested inference with and without hierarchical processing. The output of FS was resized using linear interpolation. Their results indicate that optimal performance was achieved with a resizing factor of 0.2 and hierarchical inference enabled.

5.2. Track 2: Mono

5.2.1 Baseline - ZoeDepth [4]

For this second track, we set the ZoeDepth model as our baseline, a state-of-the-art framework for single-image depth estimation. It builds over the DPT backbone [79], enhanced through a metric bins module implemented for learning metric depth rather than an affine-invariant output. As for the Stereo track, we obtain the predicted depth maps by using the original weights made available by the authors.

5.2.2 Team 1 - Lavreniuk

The Lavreniuk team’s “DeepBlend” method uses the Depth Anything v2 (DAv2) model as the primary depth estimator, with additional modifications to improve accuracy in challenging scenarios. For training, images are categorized into two groups: transparent objects and mirrors. For both categories, pseudo-labeling depth masks are created using a refined blending technique that fuses the original image with

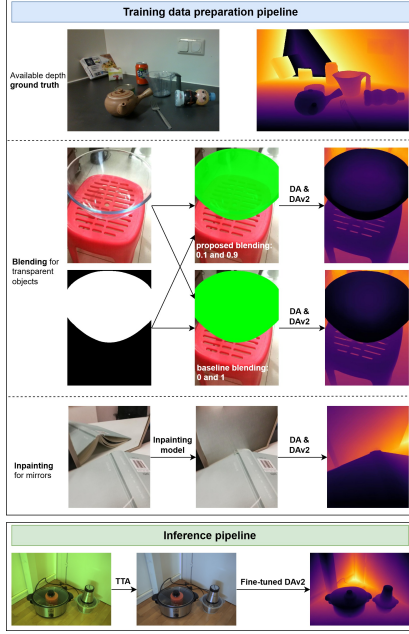


Figure 5. Network Architecture – Team Lavreniuk.

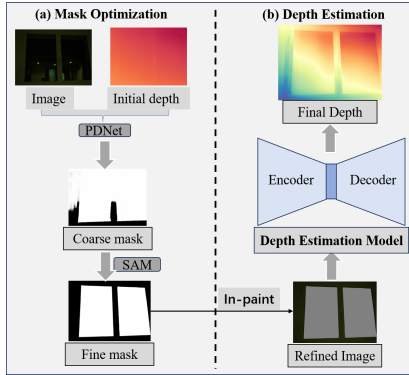


Figure 6. Network Architecture – Team colab.

a transparent object mask, improving upon [17]. For mirror surfaces, an additional restoration step is applied, where inpainted images generated using a fast Fourier convolution-based model serve as an auxiliary input for depth estimation. Experiments with various depth estimation models led to selecting Depth Anything [121] and Depth Anything v2 for pseudo-labeling. Interestingly, averaging their outputs yields better results than using DAv2 alone, which performs better for transparent and mirror surfaces. As blending and inpainting remove the transparent or mirror surfaces from images, DA effectively complements DAv2 there. The final DAv2 model is fine-tuned on a carefully selected subset of datasets, including TransCG [21], ClearGrasp [85], MIDepth [59], Hammer [41], HouseCat6D [42], MSD [123], Trans10K [114], and Booster [130]. During training, extensive data augmentations are applied to enhance robustness to varying lighting conditions. At inference, test-time augmentations (flipping and color jittering) further refine the final depth predictions.

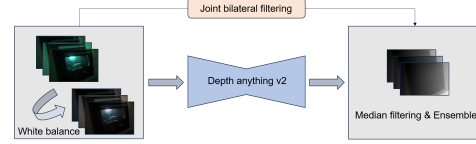


Figure 7. Network Architecture – Team PreRdw.

5.2.3 Team 2 - colab

The colab team (CodaLab: *what*) presents “DepthInpaint,” an efficient depth estimation optimization framework for ToM objects combining two-stage ToM mask optimization and in-paint mechanism. Through this two-stage approach, it overcomes the physical limitations of traditional monocular depth perception and significantly improves the accuracy of depth estimation for ToM surfaces. The method first generated a rough ToM surfaces mask based on the PDNet model [66], which combines RGB image and depth information for mirror segmentation. Then, DepthInpaint refines the rough mask with Segment Anything Model [47] (SAM) based on the image and initial mask. Following this, a physics-guided image inpaint strategy is applied to the masked regions, eliminating artifacts from specular highlights and medium refraction. Finally, Depth Pro [5] is used to generate depth metric from inpaint images. The key innovations of the approach are: 1) a second-stage mask optimization strategy for ToM surfaces, and 2) a mask-guided image-in-paint mechanism.

5.2.4 Team 3 - IPCV

The IPCV team (CodaLab: *JameerBabu*) employs Marigold [44] a monocular depth estimation model that leverages the visual knowledge embedded in diffusion-based image generators. Specifically, Marigold is built upon the architecture of Stable Diffusion [1], a latent diffusion model. The core component of this architecture is a denoising U-Net, which operates within the latent space of the model. To adapt Stable Diffusion for depth estimation, Marigold employs a fine-tuning protocol that focuses on this denoising U-Net while preserving the integrity of the latent space. This fine-tuning is performed using synthetic RGB-D datasets, such as Hypersim [82] and Virtual KITTI [6], and can be completed within a few days on a single GPU. For implementation, the team directly used inference with the pretrained weights from Marigold repository.

5.2.5 Team 4 - PreRdw

The PreRdw team (CodaLab: *jingc*) presents “Reflective Depth Wizard”, a method that takes advantage of Depth Anything V2’s inherent generalization capability. The team fine-tuned the model for accurately estimating depth on challenging reflective and transparent materials, initializing training using the pre-trained Hypersim model as their

foundation. To prepare the training data, they converted Booster depth annotations into usable depth maps and computed scale (s) and shift (t) parameters to properly align the model’s predictions with the ground truth depths during training. All input images were resized to 518×714 pixels for processing. During inference, the team implemented several enhancements: (i) Color calibration: Applied the Gray World algorithm to normalize the input image color distributions; (ii) Depth refinement: Processed raw depth predictions with a bilateral filter to smooth surfaces while preserving edge details; (iii) Ensemble optimization: Combined results across varying lighting conditions and applied median filtering to reduce noise in final depth maps.

Acknowledgments

This work was partially supported by the Humboldt Foundation. We thank the NTIRE 2025 sponsors: ByteDance, Meituan, Kuaishou, and University of Wurzburg (Computer Vision Lab).

A. NTIRE 2025 Organizers

Title:

NTIRE 2025 Challenge on HR Depth from Images of Specular and Transparent Surfaces

Members:

Pierluigi Zama Ramirez¹ (pierluigi.zama@unibo.it), Alex Costanzino¹, Fabio Tosi¹, Matteo Poggi¹, Samuele Salti¹, Stefano Mattoccia¹, Luigi Di Stefano¹, Radu Timofte²

Affiliations:

¹ University of Bologna, Italy

² Computer Vision Lab, University of Würzburg, Germany

B. Track 1: Teams and Affiliations

NJUST-KMG

Members:

Zhe Zhang¹ (zhe.zhang@njust.edu.cn), Yang Yang (yyang@njust.edu.cn)¹

Affiliations:

¹ Nanjing University of Science and Technology, China

Robot01-vRobotit

Members:

Wu Chen¹ (chenw@bupt.edu.cn), Anlong Ming¹ (mal@bupt.edu.cn), Mingshuai Zhao¹ (mingshuai.z@bupt.edu.cn), Mengying Yu¹ (yumengying@bupt.edu.cn), Shida Gao¹ (gaostar2024@bupt.edu.cn), Xiangfeng Wang¹ (xiangfeng_w@foxmail.com), Feng Xue² (feng.xue@unitn.it)

Affiliations:

¹ Beijing University of Posts and Telecommunications,

China

² University of Trento, Italy

Samsung R&D Institute China-Beijing (SRC-B)

Members:

Jun Shi¹ (jun7.shi@samsung.com), Yong Yang¹, Yong A¹, Yixiang Jin¹, Dingzhe Li¹

Affiliations:

¹ Samsung R&D Institute China-Beijing (SRC-B)

weouibaguette

Members:

Aryan Shukla¹ (aryan.shukla.1@ens.etsmtl.ca), Liam Frija-Altarac¹ (liam.frija-altarac.1@ens.etsmtl.ca), Matthew Toews¹

Affiliations:

¹ École de technologie supérieure (ÉTS), Montréal, Canada

C. Track 2: Teams and Affiliations

colab

Members:

Hui Geng¹ (gengh666666@163.com), Tianjiao Wan¹, Zijian Gao¹, Qisheng Xu¹, Kele Xu¹, Zijian Zang²

Affiliations:

¹ National University of Defense Technology, Changsha, China

² Fudan University

IPCVC

Members:

Jameer Babu Pinjari¹ (jameer.jb@gmail.com), Kuldeep Purohit¹ (kuldeep_purohit3@gmail.com)

Affiliations:

¹ Independent Researchers

Lavreniuk

Members:

Mykola Lavreniuk¹ (nick_93@ukr.net)

Affiliations:

¹ Space Research Institute NASU-SSAU, Kyiv, Ukraine

PreRdw

Members:

Jing Cao¹ (caojing@stu.hit.edu.cn), Shenyi Li¹ (504774430@qq.com), Kui Jiang¹ (jiangkui@hit.edu.cn), Junjun Jiang¹ (jiangjunjun@hit.edu.cn), Yong Huang¹ (huangyong@hit.edu.cn)

Affiliations:

¹ Harbin Institute of Technology, China

References

- [1] Stable diffusion v1.5 model card, 2022.
- [2] Filippo Aleotti, Fabio Tosi, Pierluigi Zama Ramirez, Matteo Poggi, Samuele Salti, Stefano Mattoccia, and Luigi Di Stefano. Neural disparity refinement for arbitrary resolution stereo. In *International Conference on 3D Vision*, 2021. 3DV.
- [3] Luca Bartolomei, Fabio Tosi, Matteo Poggi, and Stefano Mattoccia. Stereo anywhere: Robust zero-shot deep stereo matching even where either stereo or mono fail. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025.
- [4] Shariq Farooq Bhat, Reiner Birkel, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth, 2023.
- [5] Aleksei Bochkovskii, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R Richter, and Vladlen Koltun. Depth pro: Sharp monocular metric depth in less than a second. *arXiv preprint arXiv:2410.02073*, 2024.
- [6] Johann Cabon, Naila Murray, and Martin Humenberger. Virtual kitti 2. *arXiv preprint arXiv:2001.10773*, 2020.
- [7] Jia-Ren Chang and Yong-Sheng Chen. Pyramid stereo matching network. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5410–5418, 2018.
- [8] Weifeng Chen, Zhao Fu, Dawei Yang, and Jia Deng. Single-image depth perception in the wild. In *Proc. NeurIPS*, 2016.
- [9] Zheng Chen, Kai Liu, Jue Gong, Jingkai Wang, Lei Sun, Zongwei Wu, Radu Timofte, Yulun Zhang, et al. NTIRE 2025 challenge on image super-resolution (x4): Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [10] Zheng Chen, Jingkai Wang, Kai Liu, Jue Gong, Lei Sun, Zongwei Wu, Radu Timofte, Yulun Zhang, et al. NTIRE 2025 challenge on real-world face restoration: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [11] Junda Cheng, Longliang Liu, Gangwei Xu, Xianqi Wang, Zhaoxing Zhang, Yong Deng, Jinliang Zang, Yurui Chen, Zhipeng Cai, and Xin Yang. Monster: Marry monodepth to stereo unleashes power. *arXiv preprint arXiv:2501.08643*, 2025.
- [12] Xinjing Cheng, Peng Wang, and Ruigang Yang. Learning depth with convolutional spatial propagation network. *IEEE transactions on pattern analysis and machine intelligence*, 42(10):2361–2379, 2019.
- [13] Xuelian Cheng, Yiran Zhong, Mehrtash Harandi, Yuchao Dai, Xiaojun Chang, Hongdong Li, Tom Drummond, and Zongyuan Ge. Hierarchical neural architecture search for deep stereo matching. *Advances in Neural Information Processing Systems*, 33, 2020.
- [14] Jaehoon Choi, Dongki Jung, Yonghan Lee, Deokhwa Kim, Dinesh Manocha, and Donghwan Lee. Selfdeco: Self-supervised monocular depth completion in challenging indoor environments. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 467–474. IEEE, 2021.
- [15] Marcos Conde, Radu Timofte, et al. NTIRE 2025 challenge on raw image restoration and super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [16] Marcos Conde, Radu Timofte, et al. Raw image reconstruction from RGB on smartphones. NTIRE 2025 challenge report. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [17] Alex Costanzino, Pierluigi Zama Ramirez, Matteo Poggi, Fabio Tosi, Stefano Mattoccia, and Luigi Di Stefano. Learning depth estimation for transparent and mirror surfaces. In *The IEEE International Conference on Computer Vision*, 2023. ICCV.
- [18] Shivam Duggal, Shenlong Wang, Wei-Chiu Ma, Rui Hu, and Raquel Urtasun. Deeppruner: Learning efficient stereo matching via differentiable patchmatch. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4384–4393, 2019.
- [19] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Proc. NeurIPS*, 2014.
- [20] Egor Ershov, Sergey Korchagin, Alexei Khalin, Artyom Panshin, Arseniy Terekhin, Ekaterina Zaychenkova, Georgiy Lobarev, Vsevolod Plokhotnyuk, Denis Abramov, Elisey Zhdanov, Sofia Dorogova, Yasin Mamedov, Nikola Banic, Georgii Perevozchikov, Radu Timofte, et al. NTIRE 2025 challenge on night photography rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [21] Hongjie Fang, Hao-Shu Fang, Sheng Xu, and Cewu Lu. Transcg: A large-scale real-world dataset for transparent object depth completion and a grasping baseline. *IEEE Robotics and Automation Letters*, 7(3):7383–7390, 2022.
- [22] Xiao Fu, Wei Yin, Mu Hu, Kaixuan Wang, Yuexin Ma, Ping Tan, Shaojie Shen, Dahua Lin, and Xiaoxiao Long. Geowizard: Unleashing the diffusion priors for 3d geometry estimation from a single image. In *ECCV*, 2024.
- [23] Yuqian Fu, Xingyu Qiu, Bin Ren Yanwei Fu, Radu Timofte, Nicu Sebe, Ming-Hsuan Yang, Luc Van Gool, et al. NTIRE 2025 challenge on cross-domain few-shot object detection: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [24] Adrien Gaidon, Greg Shakhnarovich, Rares Ambrus, Victor Guizilini, Igor Vasiljevic, Matthew Walter, Sudeep Pillai, and Nick Kolkin. Dense depth for autonomous driving (DDAD) challenge (<https://sites.google.com/view/mono3d-workshop>), 2021.
- [25] Andreas Geiger, Martin Roser, and Raquel Urtasun. Efficient large-scale stereo matching. In *Asian conference on computer vision*, pages 25–38. Springer, 2010.
- [26] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proc. CVPR*, 2017.
- [27] Clément Godard, Oisín Mac Aodha, Michael Firman, and

- Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *Proc. ICCV*, 2019.
- [28] Juan Luis GonzalezBello and Munchurl Kim. Forget about the lidar: Self-supervised depth estimators with med probability volumes. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 12626–12637. Curran Associates, Inc., 2020.
- [29] Weiyu Guo, Zhaoshuo Li, Yongkui Yang, Zheng Wang, Russell H Taylor, Mathias Unberath, Alan Yuille, and Yingwei Li. Context-enhanced stereo transformer. In *European Conference on Computer Vision*, pages 263–279. Springer, 2022.
- [30] Xiaoyang Guo, Hongsheng Li, Shuai Yi, Jimmy Ren, and Xiaogang Wang. Learning monocular depth by distilling cross-domain stereo networks. In *Proc. ECCV*, 2018.
- [31] Xiaoyang Guo, Kai Yang, Wukui Yang, Xiaogang Wang, and Hongsheng Li. Group-wise correlation stereo network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3273–3282, 2019.
- [32] Shuhao Han, Haotian Fan, Fangyuan Kong, Wenjie Liao, Chunle Guo, Chongyi Li, Radu Timofte, et al. NTIRE 2025 challenge on text to image generation model quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [33] Jing He, Haodong Li, Wei Yin, Yixun Liang, Leheng Li, Kaiqiang Zhou, Hongbo Zhang, Bingbing Liu, and Ying-Cong Chen. Lotus: Diffusion-based visual foundation model for high-quality dense prediction. *arXiv preprint arXiv:2409.18124*, 2024.
- [34] Wenbo Hu, Xiangjun Gao, Xiaoyu Li, Sijie Zhao, Xiaodong Cun, Yong Zhang, Long Quan, and Ying Shan. Depthrafter: Generating consistent long depth sequences for open-world videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025.
- [35] Andrey Ignatov, Grigory Malivenko, David Plowman, Samarth Shukla, and Radu Timofte. Fast and accurate single-image depth estimation on mobile devices, mobile ai 2021 challenge: Report. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 2545–2557, June 2021.
- [36] Varun Jain, Zongwei Wu, Quan Zou, Louis Florentin, Henrik Turbell, Sandeep Siddhartha, Radu Timofte, et al. NTIRE 2025 challenge on video quality enhancement for video conferencing: Datasets, methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [37] Huaizu Jiang, Gustav Larsson, Michael Maire Greg Shakhnarovich, and Erik Learned-Miller. Self-supervised relative depth learning for urban scene understanding. In *Proc. ECCV*, 2018.
- [38] Hualie Jiang, Zhiqiang Lou, Laiyan Ding, Rui Xu, Minglang Tan, Wenjie Jiang, and Rui Huang. Defomstereo: Depth foundation model based stereo matching. *arXiv preprint arXiv:2501.09466*, 2025.
- [39] Junpeng Jing, Jiankun Li, Pengfei Xiong, Jiangyu Liu, Shuaicheng Liu, Yichen Guo, Xin Deng, Mai Xu, Lai Jiang, and Leonid Sigal. Uncertainty guided adaptive warping for robust and efficient stereo matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3318–3327, October 2023.
- [40] Adrian Johnston and Gustavo Carneiro. Self-supervised monocular trained depth estimation using self-attention and discrete disparity volume. In *Proc. CVPR*, 2020.
- [41] HyunJun Jung, Patrick Ruhkamp, Guangyao Zhai, Nikolas Brasch, Yitong Li, Yannick Verdie, Jifei Song, Yiren Zhou, Anil Armagan, Slobodan Ilic, et al. On the importance of accurate geometry data for dense 3d vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 780–791, 2023.
- [42] HyunJun Jung, Shun-Cheng Wu, Patrick Ruhkamp, Guangyao Zhai, Hannah Schieber, Giulia Rizzoli, Pengyuan Wang, Hongcheng Zhao, Lorenzo Garattoni, Sven Meier, et al. Housecat6d-a large-scale multi-modal category level 6d object perception dataset with household objects in realistic scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22498–22508, 2024.
- [43] Bingxin Ke, Dominik Narnhofer, Shengyu Huang, Lei Ke, Torben Peters, Katerina Fragkiadaki, Anton Obukhov, and Konrad Schindler. Video depth without video models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025.
- [44] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing diffusion-based image generators for monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [45] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [46] Sameh Khamis, Sean Fanello, Christoph Rhemann, Adarsh Kowdle, Julien Valentin, and Shahram Izadi. Stereonet: Guided hierarchical refinement for real-time edge-aware depth prediction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 573–590, 2018.
- [47] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023.
- [48] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli Vander-Bilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.
- [49] Henrik Kretzschmar, Alex Liniger, Jose M. Alvarez, Yan Wang, Vincent Casser, Fisher Yu, Marco Pavone, Bo Li, Andreas Geiger, Peter Ondruska, Li Erran Li, Dragomir Angelov, John Leonard, and Luc Van Gool. Argoverse stereo competition (<https://cvpr2022.wad.vision/>), 2021, 2022.
- [50] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Fed-

- erico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *2016 Fourth international conference on 3D vision (3DV)*, pages 239–248. IEEE, 2016.
- [51] Sangmin Lee, Eunpil Park, Angel Canelo, Hyunhee Park, Youngjo Kim, Hyungju Chun, Xin Jin, Chongyi Li, Chun-Le Guo, Radu Timofte, et al. NTIRE 2025 challenge on efficient burst hdr and restoration: Datasets, methods, and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [52] Jiankun Li, Peisen Wang, Pengfei Xiong, Tao Cai, Ziwei Yan, Lei Yang, Jiangyu Liu, Haoqiang Fan, and Shuaicheng Liu. Practical stereo matching via cascaded recurrent network with adaptive correlation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16263–16272, 2022.
- [53] Xin Li, Yeying Jin, Xin Jin, Zongwei Wu, Bingchen Li, Yufei Wang, Wenhan Yang, Yu Li, Zhibo Chen, Bihan Wen, Robby Tan, Radu Timofte, et al. NTIRE 2025 challenge on day and night raindrop removal for dual-focused images: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [54] Xin Li, Xijun Wang, Bingchen Li, Kun Yuan, Yizhen Shao, Suhang Yao, Ming Sun, Chao Zhou, Radu Timofte, and Zhibo Chen. NTIRE 2025 challenge on short-form ugc video quality assessment and enhancement: Kwaisr dataset and study. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [55] Xin Li, Kun Yuan, Bingchen Li, Fengbin Guan, Yizhen Shao, Zihao Yu, Xijun Wang, Yiting Lu, Wei Luo, Suhang Yao, Ming Sun, Chao Zhou, Zhibo Chen, Radu Timofte, et al. NTIRE 2025 challenge on short-form ugc video quality assessment and enhancement: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [56] Zhenyu Li, Shariq Farooq Bhat, and Peter Wonka. Patchfusion: An end-to-end tile-based framework for high-resolution monocular metric depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [57] Zhaoshuo Li, Xingtong Liu, Nathan Drenkow, Andy Ding, Francis X Creighton, Russell H Taylor, and Mathias Unberath. Revisiting stereo depth estimation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6197–6206, 2021.
- [58] Jie Liang, Radu Timofte, Qiaosi Yi, Zhengqiang Zhang, Shuaizheng Liu, Lingchen Sun, Rongyuan Wu, Xindong Zhang, Hui Zeng, Lei Zhang, et al. NTIRE 2025 the 2nd restore any image model (RAIM) in the wild challenge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [59] Yuan Liang, Zitian Zhang, Chuhua Xian, and Shengfeng He. Delving into multi-illumination monocular depth estimation: A new dataset and method. *IEEE Transactions on Multimedia*, 2024.
- [60] Zhengfa Liang, Yiliu Feng, Yulan Guo, Hengzhu Liu, Wei Chen, Linbo Qiao, Li Zhou, and Jianfeng Zhang. Learning for disparity estimation through feature constancy. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [61] Lahav Lipson, Zachary Teed, and Jia Deng. RAFT-Stereo: Multilevel Recurrent Field Transforms for Stereo Matching. *arXiv preprint arXiv:2109.07547*, 2021.
- [62] Xiaohong Liu, Xiongkuo Min, Qiang Hu, Xiaoyun Zhang, Jie Guo, et al. NTIRE 2025 XGC quality assessment challenge: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [63] Xiaoning Liu, Zongwei Wu, Florin-Alexandru Vasluianu, Hailong Yan, Bin Ren, Yulun Zhang, Shuhang Gu, Le Zhang, Ce Zhu, Radu Timofte, et al. NTIRE 2025 challenge on low light image enhancement: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [64] Jieming Lou, Weide Liu, Zhuo Chen, Fayao Liu, and Jun Cheng. Elfnet: Evidential local-global fusion for stereo matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17784–17793, 2023.
- [65] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [66] Haiyang Mei, Bo Dong, Wen Dong, Pieter Peers, Xin Yang, Qiang Zhang, and Xiaopeng Wei. Depth-aware mirror segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3044–3053, 2021.
- [67] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [68] S Mahdi H Miangoleh, Sebastian Dille, Long Mai, Sylvain Paris, and Yagiz Aksoy. Boosting monocular depth estimation models to high-resolution via content-adaptive multi-resolution merging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9685–9694, 2021.
- [69] Anton Obukhov, Matteo Poggi, Fabio Tosi, Ripudaman Singh Arora, Jaime Spencer, Chris Russell, Simon Hadfield, Richard Bowden, Shuaihang Wang, Zhenxin Ma, Weijie Chen, Baobei Xu, Fengyu Sun, Di Xie, Jiang Zhu, Mykola Lavreniuk, Haining Guan, Qun Wu, Yupei Zeng, Chao Lu, Huanran Wang, Guangyuan Zhou, Hao-tian Zhang, Jianxiong Wang, Qiang Rao, Chunjie Wang, Xiao Liu, Zhiqiang Lou, Hualie Jiang, Yihao Chen, Rui Xu, Minglang Tan, Zihan Qin, Yifan Mao, Jiayang Liu, Jialei Xu, Yifan Yang, Wenbo Zhao, Junjun Jiang, Xi-anming Liu, Mingshuai Zhao, Anlong Ming, Wu Chen, Feng Xue, Mengying Yu, Shida Gao, Xiangfeng Wang, Gbenga Omotara, Ramy Farag, Jacket Demby’s, Seyed Mohammad Ali Tousi, Guilherme N. DeSouza, Tuan-Anh Yang, Minh-Quang Nguyen, Thien-Phuc Tran, Albert Luginov,

- and Muhammad Shahzad. The fourth monocular depth estimation challenge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2025.
- [70] Jiahao Pang, Wenxiu Sun, Jimmy SJ. Ren, Chengxi Yang, and Qiong Yan. Cascade residual learning: A two-stage convolutional neural network for stereo matching. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [71] Matteo Poggi, Filippo Aleotti, Fabio Tosi, and Stefano Mattoccia. On the uncertainty of self-supervised monocular depth estimation. In *Proc. CVPR*, 2020.
- [72] Matteo Poggi, Seungryong Kim, Fabio Tosi, Sunok Kim, Filippo Aleotti, Dongbo Min, Kwanghoon Sohn, and Stefano Mattoccia. On the confidence of stereo matching in a deep-learning era: a quantitative evaluation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [73] Matteo Poggi and Fabio Tosi. Federated online adaptation for deep stereo. In *CVPR*, 2024.
- [74] Matteo Poggi, Fabio Tosi, Konstantinos Batsos, Philippos Mordohai, and Stefano Mattoccia. On the synergies between machine learning and binocular stereo for depth estimation from images: a survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [75] Michael Ramamonjisoa, Yuming Du, and Vincent Lepetit. Predicting sharp and accurate occlusion boundaries in monocular depth estimation using displacement fields. In *Proc. CVPR*, 2020.
- [76] Pierluigi Zama Ramirez, Fabio Tosi, Luigi Di Stefano, Radu Timofte, Alex Costanzino, Matteo Poggi, Samuele Salti, Stefano Mattoccia, Jun Shi, Dafeng Zhang, et al. Ntire 2023 challenge on hr depth from images of specular and transparent surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1384–1395, 2023.
- [77] Pierluigi Zama Ramirez, Fabio Tosi, Luigi Di Stefano, Radu Timofte, Alex Costanzino, Matteo Poggi, Samuele Salti, Stefano Mattoccia, Yangyang Zhang, Cailin Wu, Zhuangda He, Shuangshuang Yin, Jiaxu Dong, Yangchenxu Liu, Hao Jiang, Jun Shi, Yong A, Yixiang Jin, Dingzhe Li, Bingxin Ke, Anton Obukhov, Tinafu Wang, Nando Metzger, Shengyu Huang, Konrad Schindler, Yachuan Huang, Jiaqi Li, Junrui Zhang, Yiran Wang, Zihao Huang, Tianqi Liu, Zhiguo Cao, Pengzhi Li, Jui-Lin Wang, Wenjie Zhu, Hui Geng, Yuxin Zhang, Long Lan, Kele Xu, Tao Sun, Qisheng Xu, Sourav Saini, Aashray Gupta, Sahaj K. Mistry, Aryan Shukla, Vinit Jakhetiya, Sunil Jaiswal, Yuejin Sun, Zhuofan Zheng, Yi Ning, Jen-Hao Cheng, Hou-I Liu, Hsiang-Wei Huang, Cheng-Yen Yang, Zhongyu Jiang, Yi-Hao Peng, Aishi Huang, and Jenq-Neng Hwang. Ntire 2024 challenge on hr depth from images of specular and transparent surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 6499–6512, June 2024.
- [78] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. *ICCV*, 2021.
- [79] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12179–12188, October 2021.
- [80] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(3), 2022.
- [81] Bin Ren, Hang Guo, Lei Sun, Zongwei Wu, Radu Timofte, Yawei Li, et al. The tenth NTIRE 2025 efficient super-resolution challenge report. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [82] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10912–10922, 2021.
- [83] Nickolay Safonov, Alexey Bryntsev, Andrey Moskalenko, Dmitry Kulikov, Dmitriy Vatolin, Radu Timofte, et al. NTIRE 2025 challenge on UGC video enhancement: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [84] Tonmoy Saikia, Yassine Marrakchi, Arber Zela, Frank Hutter, and Thomas Brox. Autodispnet: Improving disparity estimation with automl. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1812–1823, 2019.
- [85] Shreeyak Sajjan, Matthew Moore, Mike Pan, Ganesh Nagaraja, Johnny Lee, Andy Zeng, and Shuran Song. Clear grasp: 3d shape estimation of transparent objects for manipulation. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3634–3642. IEEE, 2020.
- [86] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nešić, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *German conference on pattern recognition*, pages 31–42. Springer, 2014.
- [87] Thomas Schops, Johannes L Schonberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3260–3269. IEEE, 2017.
- [88] Jiahao Shao, Yuanbo Yang, Hongyu Zhou, Youmin Zhang, Yujun Shen, Vitor Guizilini, Yue Wang, Matteo Poggi, and Yiyi Liao. Learning temporally consistent video depth from video diffusion priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025.
- [89] Zhelun Shen, Yuchao Dai, and Zhibo Rao. Cfnet: Cascade and fused cost volume for robust stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13906–13915, June 2021.
- [90] Jun Shi, Yixiang Jin, Dingzhe Li, Haoyu Niu, Zhezhu Jin, He Wang, et al. Asgrasp: Generalizable transparent object reconstruction and grasping from rgb-d active stereo cam-

- era. *arXiv preprint arXiv:2405.05648*, 2024.
- [91] Xiao Song, Xu Zhao, Hanwen Hu, and Liangji Fang. Edgestereo: A context integrated residual pyramid network for stereo matching. In *ACCV*, 2018.
 - [92] Jaime Spencer, C. Stella Qian, Chris Russell, Simon Hadfield, Erich Graf, Wendy Adams, Andrew J. Schofield, James H. Elder, Richard Bowden, Heng Cong, Stefano Mattoccia, Matteo Poggi, Zeeshan Khan Suri, Yang Tang, Fabio Tosi, Hao Wang, Youmin Zhang, Yusheng Zhang, and Chaoqiang Zhao. The monocular depth estimation challenge. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) Workshops*, pages 623–632, January 2023.
 - [93] Jaime Spencer, C. Stella Qian, Michaela Trescakova, Chris Russell, Simon Hadfield, Erich Graf, Wendy Adams, Andrew J. Schofield, James Elder, Richard Bowden, Ali Anwar, Hao Chen, Xiaozhi Chen, Kai Cheng, Yuchao Dai, Huynh Thai Hoa, Sadat Hossain, Jianmian Huang, Mohan Jing, Bo Li, Chao Li, Baojun Li, Zhiwen Liu, Stefano Mattoccia, Siegfried Mercelis, Myungwoo Nam, Matteo Poggi, Xiaohua Qi, Jiahui Ren, Yang Tang, Fabio Tosi, Linh Trinh, S M Nadim Uddin, Khan Muhammad Umair, Kaixuan Wang, Yufei Wang, Yixing Wang, Mochu Xiang, Guangkai Xu, Wei Yin, Jun Yu, Qi Zhang, and Chaoqiang Zhao. The second monocular depth estimation challenge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2023.
 - [94] Jaime Spencer, Fabio Tosi, Matteo Poggi, Ripudaman Singh Arora, Chris Russell, Simon Hadfield, Richard Bowden, GuangYuan Zhou, ZhengXin Li, Qiang Rao, YiPing Bao, Xiao Liu, Dohyeong Kim, Jinseong Kim, Myunghyun Kim, Mykola Lavreniuk, Rui Li, Qing Mao, Jiang Wu, Yu Zhu, Jinqiu Sun, Yanning Zhang, Suraj Patni, Aradhye Agarwal, Chetan Arora, Pihai Sun, Kui Jiang, Gang Wu, Jian Liu, Xianming Liu, Junjun Jiang, Xidan Zhang, Jianing Wei, Fangjun Wang, Zhiming Tan, Jiabao Wang, Albert Luginov, Muhammad Shahzad, Seyed Hosseini, Aleksander Trajcevski, and James H. Elder. The third monocular depth estimation challenge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2024.
 - [95] Lei Sun, Andrea Alfano, Peiqi Duan, Shaolin Su, Kaiwei Wang, Boxin Shi, Radu Timofte, Danda Pani Paudel, Luc Van Gool, et al. NTIRE 2025 challenge on event-based image deblurring: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
 - [96] Lei Sun, Hang Guo, Bin Ren, Luc Van Gool, Radu Timofte, Yawei Li, et al. The tenth ntire 2025 image denoising challenge report. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
 - [97] Vladimir Tankovich, Christian Hane, Yinda Zhang, Adarsh Kowdle, Sean Fanello, and Sofien Bouaziz. Hitnet: Hierarchical iterative tile refinement network for real-time stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14362–14372, June 2021.
 - [98] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *European conference on computer vision*, pages 402–419. Springer, 2020.
 - [99] Alessio Tonioni, Fabio Tosi, Matteo Poggi, Stefano Mattoccia, and Luigi Di Stefano. Real-time self-adaptive deep stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 195–204, 2019.
 - [100] Fabio Tosi, Filippo Aleotti, Matteo Poggi, and Stefano Mattoccia. Learning monocular depth estimation infusing traditional stereo knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
 - [101] Fabio Tosi, Filippo Aleotti, Pierluigi Zama Ramirez, Matteo Poggi, Samuele Salti, Luigi Di Stefano, and Stefano Mattoccia. Distilled semantics for comprehensive scene understanding from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
 - [102] Fabio Tosi, Filippo Aleotti, Pierluigi Zama Ramirez, Matteo Poggi, Samuele Salti, Stefano Mattoccia, and Luigi Di Stefano. Neural Disparity Refinement. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 46(12):8900–8917, 2024.
 - [103] Fabio Tosi, Luca Bartolomei, and Matteo Poggi. A survey on deep stereo matching in the twenties. *International Journal of Computer Vision*, pages 1–32, 2025.
 - [104] Fabio Tosi, Yiyi Liao, Carolin Schmitt, and Andreas Geiger. Smd-nets: Stereo mixture density networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
 - [105] Fabio Tosi, Alessio Tonioni, Daniele De Gregorio, and Matteo Poggi. Nerf-supervised deep stereo. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 855–866, June 2023.
 - [106] Fabio Tosi, Pierluigi Zama Ramirez, and Matteo Poggi. Diffusion models for monocular depth estimation: Overcoming challenging conditions. In *European Conference on Computer Vision (ECCV)*, 2024.
 - [107] Florin-Alexandru Vasluianu, Tim Seizinger, Zhuyun Zhou, Cailian Chen, Zongwei Wu, Radu Timofte, et al. NTIRE 2025 image shadow removal challenge report. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
 - [108] Florin-Alexandru Vasluianu, Tim Seizinger, Zhuyun Zhou, Zongwei Wu, Radu Timofte, et al. NTIRE 2025 ambient lighting normalization challenge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
 - [109] Lijun Wang, Jianming Zhang, Yifan Wang, Huchuan Lu, and Xiang Ruan. CLIFFNet for monocular depth estimation with hierarchical embedding loss. In *Proc. ECCV*, 2020.
 - [110] Yan Wang, Zihang Lai, Gao Huang, Brian H Wang, Laurens Van Der Maaten, Mark Campbell, and Kilian Q Weinberger. Anytime stereo image depth estimation on mobile devices. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 5893–5900, 2019.
 - [111] Yingqian Wang, Zhengyu Liang, Fengyuan Zhang, Lvli Tian, Longguang Wang, Juncheng Li, Juegang Yang, Radu Timofte, Yulan Guo, et al. NTIRE 2025 challenge on light

- field image super-resolution: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [112] Jamie Watson, Michael Firman, Gabriel J. Brostow, and Daniyar Turmukhambetov. Self-supervised monocular depth hints. In *Proc. ICCV*, 2019.
- [113] Bowen Wen, Matthew Trepte, Joseph Aribido, Jan Kautz, Orazio Gallo, and Stan Birchfield. Foundationstereo: Zero-shot stereo matching. *arXiv*, 2025.
- [114] Enze Xie, Wenjia Wang, Wenhai Wang, Mingyu Ding, Chunhua Shen, and Ping Luo. Segmenting transparent objects in the wild. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*, pages 696–711. Springer, 2020.
- [115] Gangwei Xu, Xianqi Wang, Xiaohuan Ding, and Xin Yang. Iterative geometry encoding volume for stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21919–21928, 2023.
- [116] Gangwei Xu, Xianqi Wang, Xiaohuan Ding, and Xin Yang. Iterative geometry encoding volume for stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21919–21928, 2023.
- [117] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, Fisher Yu, Dacheng Tao, and Andreas Geiger. Unifying flow, stereo and depth estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [118] Gengshan Yang, Joshua Manela, Michael Happpold, and Deva Ramanan. Hierarchical deep stereo matching on high-resolution images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5515–5524, 2019.
- [119] Guorun Yang, Hengshuang Zhao, Jianping Shi, Zhidong Deng, and Jiaya Jia. Segstereo: Exploiting semantic information for disparity estimation. In *ECCV*, pages 636–651, 2018.
- [120] Kangning Yang, Jie Cai, Ling Ouyang, Florin-Alexandru Vasluianu, Radu Timofte, Jiaming Ding, Huiming Sun, Lan Fu, Jinlong Li, Chiu Man Ho, Zibo Meng, et al. NTIRE 2025 challenge on single image reflection removal in the wild: Datasets, methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [121] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024.
- [122] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *Advances in Neural Information Processing Systems*, 37:21875–21911, 2024.
- [123] Xin Yang, Haiyang Mei, Ke Xu, Xiaopeng Wei, Baocai Yin, and Rynson WH Lau. Where is my mirror? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8809–8818, 2019.
- [124] Wei Yin, Jianming Zhang, Oliver Wang, Simon Niklaus, Long Mai, Simon Chen, and Chunhua Shen. Learning to recover 3d scene shape from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 204–213, 2021.
- [125] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for match density estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6044–6053, 2019.
- [126] Pierluigi Zama Ramirez, Alex Costanzino, Fabio Tosi, Matteo Poggi, Luigi Di Stefano, Jean-Baptiste Weibel, Dominik Bauer, Doris Antensteiner, Markus Vincze, Jiaqi Li, Yachuan Huang, Junrui Zhang, Yiran Wang, Jinghong Zheng, Liao Shen, Zhiguo Cao, Ziyang Song, Zerong Wang, Ruijie Zhu, Hao Zhang, Rui Li, Jiang Wu, Xian Li, Yu Zhu, Jinqiu Sun, Yanning Zhang, Pihai Sun, Yuanqi Yao, Wenbo Zhao, Kui Jiang, Junjun Jiang, Mykola Lavreniuk, and Jui-Lin Wang. Tricky 2024 challenge on monocular depth from images of specular and transparent surfaces. In *European Conference on Computer Vision Workshops*, 2024. ECCVW.
- [127] Pierluigi Zama Ramirez, Alex Costanzino, Fabio Tosi, Matteo Poggi, Samuele Salti, Luigi Di Stefano, and Stefano Mattoccia. Booster: a benchmark for depth from images of specular and transparent surfaces. *arXiv preprint arXiv:2301.08245*, 2023.
- [128] Pierluigi Zama Ramirez, Matteo Poggi, Fabio Tosi, Stefano Mattoccia, and Luigi Di Stefano. Geometry meets semantics for semi-supervised monocular depth estimation. In *Computer Vision—ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III 14*, pages 298–313. Springer, 2019.
- [129] Pierluigi Zama Ramirez, Fabio Tosi, Luigi Di Stefano, Radu Timofte, Alex Costanzino, Matteo Poggi, Samuele Salti, Stefano Mattoccia, et al. NTIRE 2025 challenge on hr depth from images of specular and transparent surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2025.
- [130] Pierluigi Zama Ramirez, Fabio Tosi, Matteo Poggi, Samuele Salti, Stefano Mattoccia, and Luigi Di Stefano. Open challenges in deep stereo: The booster dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21168–21178, June 2022.
- [131] Jure Zbontar, Yann LeCun, et al. Stereo matching by training a convolutional neural network to compare image patches. *J. Mach. Learn. Res.*, 17(1):2287–2318, 2016.
- [132] Oliver Zendel, Angela Dai, Xavier Puig Fernandez, Andreas Geiger, Vladen Koltun, Peter Kontschieder, Adam Kortylewski, Tsung-Yi Lin, Torsten Sattler, Daniel Scharstein, Hendrik Schilling, Jonas Uhrig, and Jonas Wulff. The robust vision challenge (<http://www.robustvision.net/>), 2018, 2020, 2022.
- [133] Jiayi Zeng, Chengtang Yao, Lidong Yu, Yuwei Wu, and Yunde Jia. Parameterized cost volume for stereo matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 18347–18357, October 2023.
- [134] Feihu Zhang, Victor Prisacariu, Ruigang Yang, and Philip HS Torr. GA-Net: Guided aggregation net for end-to-end stereo matching. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

- [135] Chaoqiang Zhao, Matteo Poggi, Fabio Tosi, Lei Zhou, Qiyu Sun, Yang Tang, and Stefano Mattoccia. Gasmono: Geometry-aided self-supervised monocular depth estimation for indoor scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16209–16220, 2023.
- [136] Chao Zhou, Hong Zhang, Xiaoyong Shen, and Jiaya Jia. Unsupervised learning of stereo matching. In *The IEEE International Conference on Computer Vision (ICCV)*. IEEE, October 2017.