

Supplementary Materials

S²M²: Scalable Stereo Matching Model for Reliable Depth Estimation

1. Implementation

1.1. Training Strategy

Our training methodology consists of two phases: first, establishing baseline models for our ablation studies, and second, a sequential fine-tuning process to develop our final benchmark models, which is directly built upon the first phase.

Phase 1: Baseline Models for Ablation Studies

The models used in our ablation studies were trained following a complete, single-stage procedure. Each model was trained on image patches of 768×448 px with a batch size of 4 for 2,000,000 gradient updates. For optimization, we used the AdamW optimizer with a one-cycle learning rate scheduler, decaying the learning rate from an initial value of 1×10^{-5} to 1×10^{-6} .

The training was performed on a foundational set of synthetic datasets—including SceneFlow [11], TartanAir [19], and CREStereo [9]—amounting to approximately 0.5M image pairs. This process yielded our fully trained ablation models, which serve as the starting point for the next phase.

Phase 2: Sequential Fine-tuning for Benchmark Submission

Instead of training from scratch, our final benchmark models were developed by sequentially fine-tuning the fully trained ablation models from Phase 1. This process involves two subsequent steps:

Step 1: Fine-tuning on an Expanded Dataset: First, we took a converged ablation model and continued its training at the same fixed resolution (768×448). The key difference in this step was the use of a significantly expanded data pool, totaling up to 2M image pairs. This collection includes widely-used datasets in combination with additional sources such as FoundationStereo [23], IRS [18], Falling Things [17], Virtual KITTI 2 [6], DrivingStereo [28], and SMD-Net [16], among others. This step allows the model to consolidate its learned features and generalize across a much more diverse range of data distributions.

Step 2: Multi-Resolution Fine-tuning with Gradient Checkpointing: In the final step, the model from Step 1 was further fine-tuned across three resolutions simultaneously (768×448 , 1440×1024 , and 1920×1440). This was performed on an H100×8 GPU server with a total batch size of 56.

A significant technical hurdle at this stage is the high memory requirement for high-resolution training. Processing images at 1920×1440 px exceeds the 80GB memory capacity of a single H100 GPU. To resolve this, we employed PyTorch [13]’s gradient checkpointing feature. This technique trades computation for memory by recomputing intermediate activations during the backward pass instead of storing them all. Gradient checkpointing was essential for enabling stable, high-resolution training, allowing our final model to effectively learn features across a wide spectrum of scales.

1.2. Data Augmentation

To address the practical challenges of real-world stereo matching, we employed a targeted data augmentation strategy beyond standard color and cropping adjustments [9]. Specifically, we introduce a random horizontal shift when cropping the left and right images. This technique is designed to simulate imperfect rectification and, importantly, mimic the design of benchmark datasets like Middlebury v3 [14].

In these datasets, the stereo pairs are intentionally cropped with a sample-specific horizontal offset to maximize the overlapping field of view and manage occlusions. By replicating this characteristic in our training data, our model becomes

robust to the subtle misalignments frequently present in real-world or benchmark stereo pairs. This robustness is critical for preventing performance degradation from a domain shift during testing; without it, a model trained on ideal data learns to associate strong perspective cues with large disparity values. When a test sample violates this learned correlation due to an artificial offset, the model’s predictions fail as its internal priors conflict with the ground truth.

Furthermore, a key consequence of this horizontal shift is the introduction of negative disparity values for some pixels. While unconventional, we leverage this as a deliberate training strategy. Training with these negative disparities encourages the model to learn depth relationships in a more object-centric manner. Instead of relying solely on the common assumption of positive disparity, the network is forced to more comprehensively understand the relative spatial arrangement of objects and surfaces. This, in turn, enables more flexible and accurate depth estimation, particularly around object boundaries and in scenes with complex foreground-background interactions.

2. Loss Function Details

The final loss function consists of four components: disparity loss, occlusion loss, confidence loss, and PMC loss, each weighted by predefined hyperparameters:

$$\mathcal{L}_{\text{total}} = \lambda_D \mathcal{L}_D + \lambda_O \mathcal{L}_O + \lambda_C \mathcal{L}_C + \lambda_{\text{PMC}} \mathcal{L}_{\text{PMC}}. \quad (1)$$

The hyperparameters are set as follows: $\lambda_D = 1$, $\lambda_O = 0.1$, $\lambda_C = 0.1$, and $\lambda_{\text{PMC}} = 1$. Below, we detail the disparity, occlusion, and confidence loss formulations.

2.1. Disparity Loss

Disparity loss supervises both the initial global disparity estimate and the refined disparities at each iteration. We use a smooth L1 loss for the initial disparity, applied only to non-occluded pixels, while the refined disparities are optimized using L1 loss over all pixels. This ensures that the network learns reliable disparity predictions across all refinement stages.

$$\mathcal{L}_D = \mathcal{L}_D^{\text{init}} + \sum_{t=1}^T \mathcal{L}_D^{(t)}, \quad (2)$$

where

$$\mathcal{L}_D^{\text{init}} = \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} \text{Smooth-}\ell_1(D_{ij}^{\text{init}}, D_{ij}^{\text{gt}}), \quad (3)$$

$$\mathcal{L}_D^{(t)} = \frac{1}{H \times W} \sum_{(i,j)} |D_{ij}^{(t)} - D_{ij}^{\text{gt}}|. \quad (4)$$

Here, $\Omega = \{(i,j) \mid O(i,j) = 1\}$ is the set of non-occluded pixels, and $H \times W$ represents all pixels in the image. The initial disparity D^{init} is supervised only on non-occluded pixels, while the refined disparities $D^{(t)}$ (for $t = 1, \dots, T$) are trained on all pixels.

2.2. Occlusion and Confidence Loss

The occlusion and confidence maps are learned using an L1 loss rather than binary cross-entropy to ensure numerical stability. The losses for occlusion and confidence maps follow a similar formulation:

$$\mathcal{L}_O = \mathcal{L}_O^{\text{init}} + \sum_{t=1}^T \mathcal{L}_O^{(t)}, \quad (5)$$

$$\mathcal{L}_C = \mathcal{L}_C^{\text{init}} + \sum_{t=1}^T \mathcal{L}_C^{(t)}. \quad (6)$$

Specifically, the occlusion losses are computed over all pixels:

$$\mathcal{L}_O = \frac{1}{H \times W} \sum_{(i,j)} |O_{ij} - O_{ij}^{\text{gt}}|. \quad (7)$$

In contrast, the initial confidence loss is computed only over non-occluded pixels, while the confidence losses in the refinement steps are optimized over all pixels:

$$\mathcal{L}_C^{\text{init}} = \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} |\Gamma_{ij}^{\text{init}} - \Gamma_{ij}^{\text{gt}}|. \quad (8)$$

$$\mathcal{L}_C^{(t)} = \frac{1}{|H \times W|} \sum_{(i,j)} |\Gamma_{ij}^{(t)} - \Gamma_{ij}^{\text{gt}}|. \quad (9)$$

Ground Truth Definition - Confidence Map (Γ): Defined as 1 if the absolute disparity error is below 4 pixels, which corresponds to one pixel in the 1/4 resolution feature space; otherwise, 0.

- **Occlusion Map (O):** Defined by left-right consistency check, where a pixel is non-occluded if the left-right disparity warping error is below 2 pixels.

This loss formulation ensures that occlusion and confidence estimates remain reliable throughout training while preventing numerical instability.

3. Synthetic Data Generation

3.1. Train Dataset

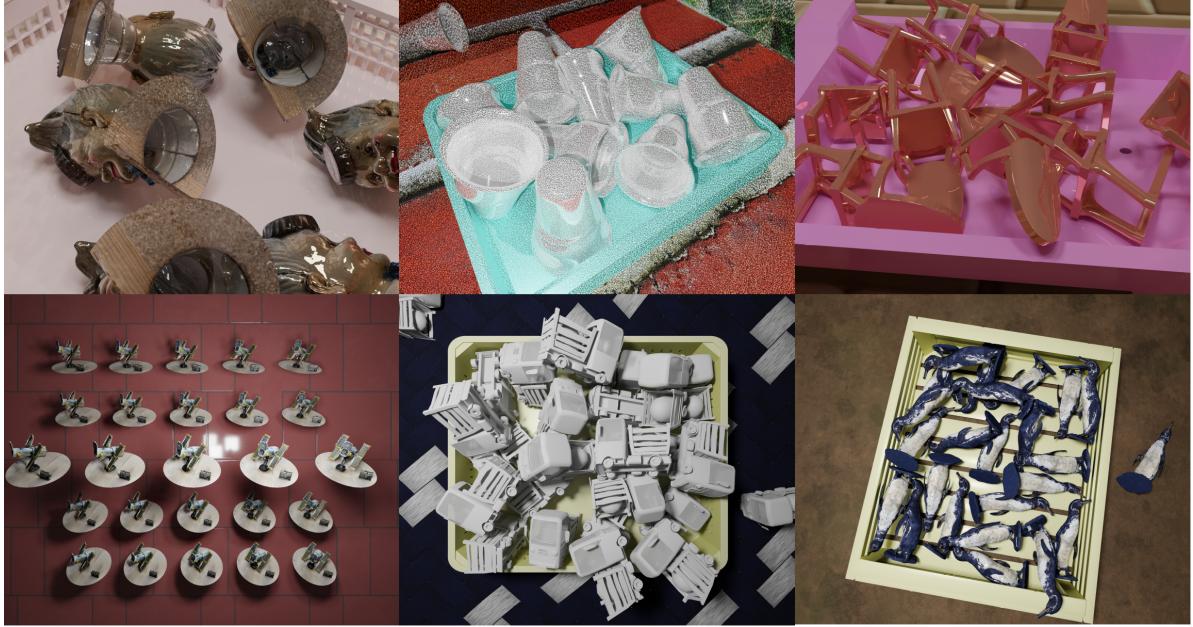


Figure 1. Configurations of scenes in our synthetic training dataset: Objects from Objaverse are randomly dropped into a bin or arranged in structured layouts, with variations in surface materials such as reflective surfaces. Some scenes also include projected random patterns to simulate active stereo scenarios.

Our basic training dataset comprises a diverse set of available synthetic datasets [6, 9, 11, 16, 17, 18, 19, 23, 28]. To further enhance the training dataset, we generated a specialized synthetic dataset to complement the existing data by leveraging 3D assets from Objaverse [4] and BlenderKit [1] within the Blender environment. The dataset was constructed through physics-based simulations where uniform object types were randomly dropped into a bin to generate realistic stacked configurations. Scenes were rendered using a stereo camera system to obtain passive RGB images, depth maps, and occlusion maps. To increase diversity, we varied object sizes, lighting conditions, bin shapes, and camera configurations. The camera was randomly positioned within the space between the top and the sides of the bin. Camera parameters were assigned random values within the following ranges for each scene generation. Field of View (FOV) was selected from 35° to 45°, and baseline was selected from 0.25m to 0.35m. The bin shape was randomly selected from a set of dozens of predefined shapes, and

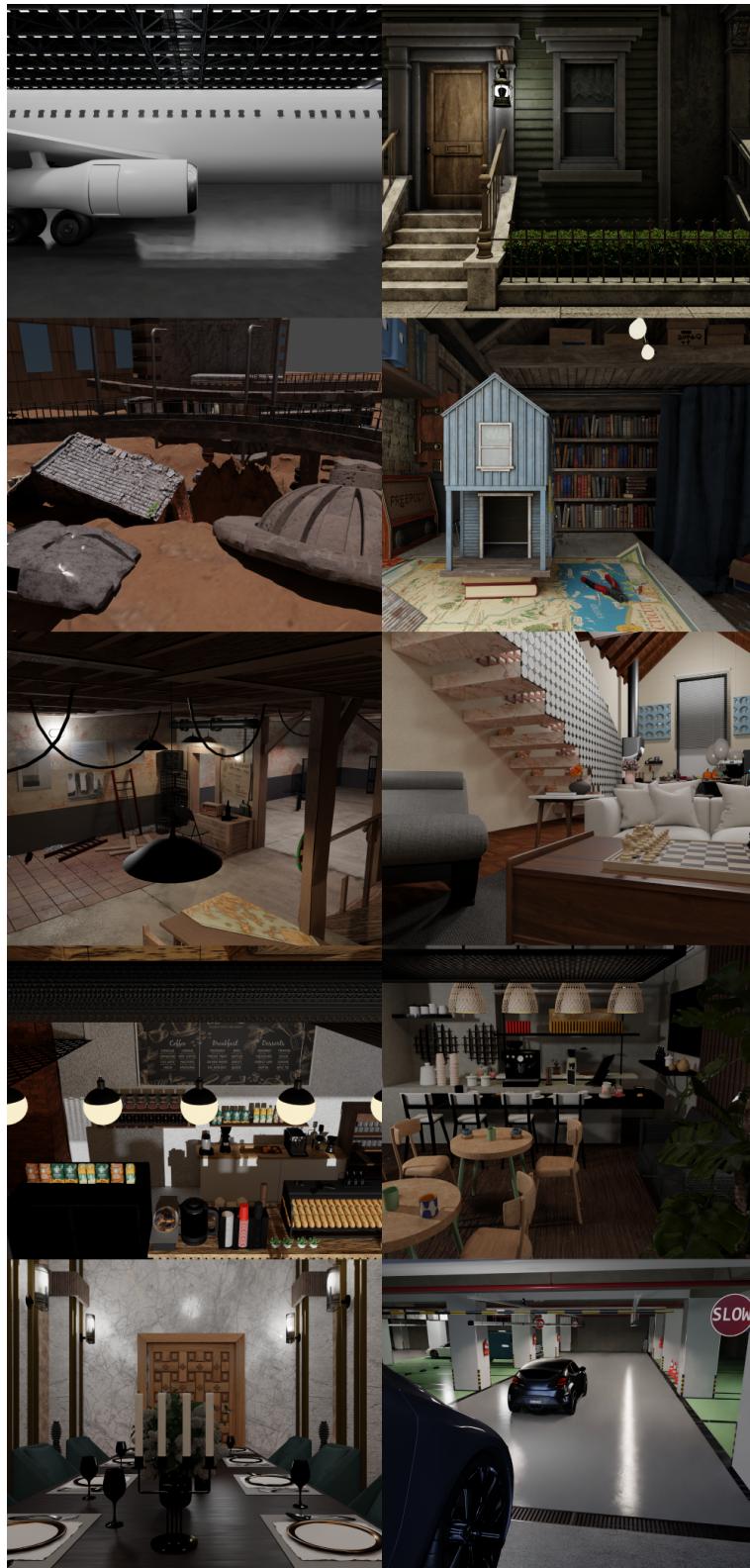


Figure 2. Configurations of scenes in our synthetic benchmark dataset: Our benchmark synthetic dataset was generated across 10 environments, encompassing indoor and outdoor settings. These include house interiors, cafes, parking lots, aircraft hangers, residential areas and abandoned urban areas, comprehensively covering various real-world spatial contexts.

its side lengths were randomly determined within the range of 0.6m to 1.2m. The materials for the floor and walls were randomly chosen from a collection of 2,176 materials provided from BlenderKit. We randomly set 2 to 5 lights assigning random values to their locations and energy. The object size was normalized such that the bounding bin volume of the object was between 10% and 20% of the bin volume. All objects within a single scene had the same size. A total of 8 to 30 objects were randomly dropped from above the bin. The scenes were rendered 3 seconds after the objects were dropped.

Additionally, we varied material properties to include challenging surface types, such as reflective and transparent materials, making the dataset more diverse and representative of real-world conditions. Alpha, roughness, and metallic values were randomly selected within the range of 0.1 to 0.7, and the index of refraction (IOR) values were randomly assigned between 1.5 and 2.5. Both the objects and the bin were designed to feature reflective and transparent surfaces.

One of the key challenges in stereo matching arises when objects exhibit repetitive structures, which can lead to ambiguities in disparity estimation. To address this, we created additional datasets containing horizontally aligned identical objects, forcing the network to learn more robust correspondences in such scenarios. Identical objects were arranged in a grid format of either 5×5 or 6×6 with uniform spacing. The objects within each horizontal row were positioned at the same height, ensuring consistent gaps between them.

Our final synthetic training dataset consists of approximately 40,000 scenes, complementing existing synthetic datasets while specifically addressing limitations such as repetitive textures, low-texture regions, and disparity ambiguities.

3.2. Our Synthetic Benchmark Dataset

Our benchmark dataset was constructed utilizing scenes and 3D models sourced from BlenderKit. As shown in Figure 2, the dataset comprises 10 distinct environmental categories spanning both indoor and outdoor settings under diverse illumination conditions. A total of 130 scenes were generated across these environments to address challenges in stereo depth estimation. The dataset specifically incorporates optically complex surfaces including reflective and transparent objects, horizontal repetitive structures, and scenes containing sharp depth discontinuities where distant background elements are juxtaposed with proximate foreground objects. This dataset is designed to evaluate the challenging aspects of stereo depth estimation through controlled introduction of real-world complexities.

3.2.1. Analysis

Each scene in our benchmark dataset has a resolution of 2448×2048 pixels, with stereo baselines varied between 0.1 m and 0.5 m. The disparity distribution spans from 2.19 pixels to over 1000 pixels as shown in Figure 3. Since many existing models struggle with extremely large disparity values, we performed benchmark evaluations in the following section with maximum disparity limits set at 320 pixels and 512 pixels. These limits covered 82.304% and 93.442% of the entire benchmark dataset, respectively, ensuring that the analysis remains representative of the dataset's overall distribution.

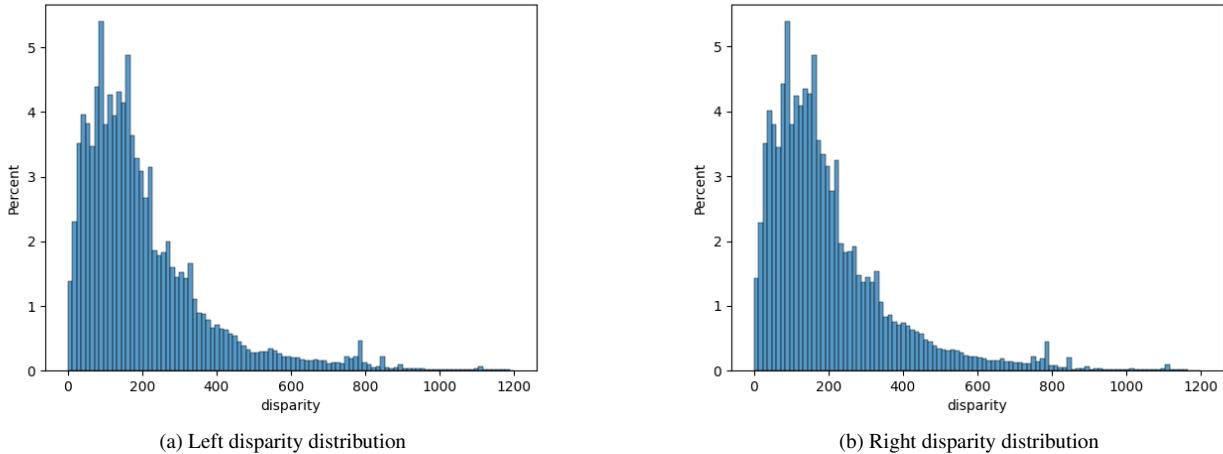


Figure 3. Benchmark dataset disparity distribution (x-axis: disparity (pixels), y-axis: proportion of the total dataset (%)): The figure illustrates the disparity distribution when the maximum disparity is set to 1200 pixels, covering 99.591% of the entire benchmark dataset.

While synthetic data generation assumes ideal stereo camera calibration, achieving perfect calibration in real-world sce-

narios is challenging, often leading to geometric misalignment. This misalignment introduces a vertical shift in the images. To rigorously examine its impact, we generated an imperfect synthetic dataset that simulates realistic calibration errors under identical conditions as the perfectly calibrated case. Instead of applying a simple vertical shift to the images, the right camera was deliberately offset perpendicular to the baseline axis by 0.4–1.8 mm during rendering, inducing physically consistent y-disparities. Across all 130 scenes, the average y-disparity measured 0.9964 pixels, with 95% of values distributed between 0.8561 and 1.1302 pixels per scene (Table 1).

| | y-disparity |
|---------|-----------------------------|
| Average | 0.8561 pixel - 1.1302 pixel |
| Minimum | 0.0132 pixel - 0.6887 pixel |
| Maximum | 0.8561 pixel - 9.9977 pixel |

Table 1. Y-disparity analysis on imperfect calibration dataset with value range for 95%

4. Critical Analysis of Fine-Tuning Effects on the KITTI Benchmark

While fine-tuning on the KITTI dataset [7, 12] is a common practice to boost benchmark scores, its impact on a model’s true generalization capability warrants a deeper investigation. In this section, we present additional experiments arguing that the observed improvements in error metrics after fine-tuning may not represent a genuine enhancement in stereo accuracy, but rather an overfitting to the inherent noise and systematic biases present in the LiDAR-derived ground truth. Specifically, the noise often originates from the LiDAR sensor’s own physical limitations, while the systematic biases can be traced to more structured issues, such as spatio-temporal misalignments caused by moving vehicles or subtle calibration errors between the stereo rig and the LiDAR sensor.

Quantitative Analysis: Contradiction in Metrics

Our initial analysis begins with a quantitative comparison of our model before and after fine-tuning on the KITTI 2015 [12] training set. As shown in Table 2, fine-tuning leads to a significant reduction in standard error metrics such as End-Point-Error (EPE) and Bad-2.0 on the training data.

However, this apparent improvement is contradicted by a notable decrease in photometric consistency, which we measure using the Structural Similarity Index (SSIM) [21]. To calculate SSIM, we warp the right image to the left view’s perspective using the estimated disparity map and measure their similarity. This comparison was performed exclusively on non-occluded regions, which are identified using a predicted occlusion map from the model. While we acknowledge that SSIM is not a perfect representation of photometric consistency due to the presence of non-Lambertian surfaces (e.g., reflections), a strong positive correlation between improvements in error metrics and photometric consistency is typically expected. The observed divergence—where disparity metrics improve while SSIM deteriorates—is a strong indication of an underlying issue with the fine-tuning process. It strongly suggests that the fine-tuned model is producing disparity maps that are less geometrically consistent with the scene, potentially by overfitting to artifacts in the ground truth labels rather than reconstructing the true scene structure. This effect has also been observed in other works [29], where fine-tuning on specific real-world datasets degraded performance on other domains.

| Fine Tuning | Split | EPE | Bad-2.0 | SSIM |
|-------------|-------|-------|---------|-----------------------------|
| No | Train | 0.885 | 1.100 | 0.883 (-) |
| Yes | Train | 0.461 | 0.280 | 0.840 (\downarrow 0.043) |
| No | Test | - | - | 0.879 (-) |
| Yes | Test | - | - | 0.844 (\downarrow 0.035) |

Table 2. Impact of fine-tuning on the KITTI 2015 training and test sets. While fine-tuning significantly improves disparity accuracy metrics (EPE, Bad-2.0) on the training set, the photometric consistency (SSIM) decreases on both training and test sets. This suggests potential overfitting to dataset-specific noise and a reduction in generalization.

Qualitative Inspection: Evidence of Ground Truth Bias

To understand the source of this metric-consistency discrepancy, we performed detailed 3D visualizations and manual inspections of the disparity maps. We conducted a comprehensive comparison across three different architectures:

- **Our Model**, evaluated in both its zero-shot and fine-tuned states.
- **Selective-IGEV** [20], also evaluated in its zero-shot and fine-tuned states.
- **FoundationStereo** [23], evaluated in its zero-shot configuration.

Our investigation revealed a consistent trend among all three zero-shot models, as visualized in Figure 4. Our model and FoundationStereo [23] produced disparity maps that, when visualized as 3D point clouds, were structurally sound, visually clean, and highly consistent with the image content. The zero-shot output from Selective-IGEV [20], albeit with slightly lower visual quality, was broadly consistent with this finding.

Crucially, we observed spatially varying biases between the outputs of all these zero-shot models and the provided KITTI ground truth. We manually verified biases in areas with clear geometric features, such as the sharp, well-defined edges of vehicle license plates, where our inspection confirmed that the ground truth disparity exhibits a bias of up to 2 pixels compared to our manually calculated values.

In contrast, the fine-tuned versions of both our model and Selective-IGEV exhibited the same clear signs of adapting to this ground truth noise. These models produce disparity maps with blurry edges and structures that, while numerically closer to the noisy ground truth, are less consistent with the actual image content. This demonstrates that the models have effectively learned to replicate the dataset’s inherent bias.

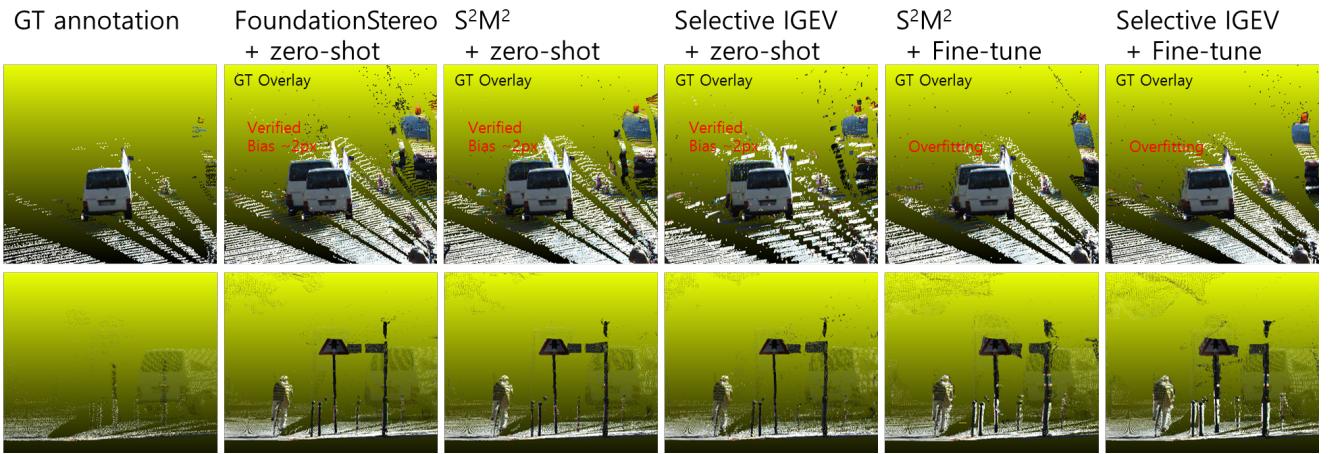


Figure 4. Qualitative comparison illustrating the negative effects of fine-tuning on a KITTI 2015 sample. **Top row:** Focuses on a vehicle area to highlight the bias between the zero-shot and fine-tuned models. The ground truth annotation from the training data is overlaid to visualize how the fine-tuned model adapts to the noisy label. **Bottom row:** Shows a pillar structure where the fine-tuned model estimates the pillar to be significantly thicker than its actual structure, a distortion not present in the zero-shot output.

The combined quantitative and qualitative evidence strongly suggests that the numerical improvements from fine-tuning on KITTI are deceptive. The models are effectively learning the noise and systematic biases of the LiDAR-derived ground truth. The concurrent drop in photometric consistency (SSIM) and the qualitative evidence of fitting to noisy labels indicate that this fine-tuning process harms the model’s generalization capabilities, even as it achieves a lower error score on the benchmark’s training set.

5. Comprehensive Analysis on Our Synthetic Benchmark

To gain deeper insight into the capabilities and limitations of stereo matching models, we conducted a comprehensive analysis on our synthetic benchmark. We investigate model robustness by systematically evaluating the impact of three key factors: **stereo calibration quality**, **input image resolution**, and **maximum disparity range**. In the following sections, we vary one factor at a time to isolate its effect on model accuracy and robustness. Our analysis also examines the trade-off between accuracy and processing speed, offering a thorough assessment of model efficiency under different conditions.

We benchmarked a total of 16 models, including three variants of S^2M^2 , on a single NVIDIA H100 GPU with 80GB of VRAM. To ensure a fair comparison, we adhered to the hyperparameter configurations and refinement iterations proposed in each model’s original publication. Furthermore, our synthetic data contains unreliable values in regions of large depth, which correspond to small disparity values. To mitigate the impact of this, we excluded disparities below 10 from our evaluation metrics.

5.1. Analysis of Calibration Error Impact

As expected, the introduction of stereo calibration errors results in a performance decrease across all benchmarked models. However, the magnitude of this degradation is not uniform across the board. As illustrated in Figure 5, we observed a particularly severe drop in accuracy for models designed to be lightweight. For instance, models such as FastACV [26] and CoEx [2] exhibited a more pronounced performance decline compared to more robust, heavier architectures. This finding highlights a critical trade-off: while lightweight models provide computational efficiency, they may lack resilience to the inevitable calibration imperfections found in real-world scenarios.

5.2. Analysis of Input Resolution Impact

Ideally, higher input resolution should lead to more accurate disparity predictions. However, our analysis reveals that this is not universally the case. To investigate this, we evaluated all models by combining two resolutions (full, half) with two maximum disparity settings ($D_{max} = 320$ and $D_{max} = 960$). The results, plotted as accuracy versus processing speed, are presented in Figure 6 and Figure 7.

As expected, all models run significantly faster at half resolution. A distinct trend was observed in lightweight models, which not only gained speed but also often achieved superior accuracy at half resolution, especially with a smaller disparity range ($D_{max} = 320$). It remains inconclusive whether this behavior stems from inherent architectural limitations or a domain gap, as we utilized publicly available pre-trained weights which were likely trained on lower-resolution datasets.

In stark contrast, our S²M² models demonstrate robust scalability, consistently benefiting from higher-resolution inputs and larger disparity ranges to achieve optimal performance. Furthermore, they exhibit an excellent accuracy-speed trade-off. Even at half resolution, our models, particularly S²M²-S, provide a highly competitive balance between fast inference and high accuracy, positioning them as an effective solution for various application scenarios.

5.3. Analysis of Maximum Disparity Range Impact

This analysis focuses on how the maximum disparity (D_{max}) hyperparameter affects models that require it for cost volume construction. Other models that do not use this parameter were excluded from this test, with the exception of our S²M² model, which is included as a performance baseline. Ideally, for the models under evaluation, increasing the D_{max} value should improve accuracy, as a larger range covers more ground truth disparities. We investigated this behavior by plotting performance against four D_{max} values—320, 512, 768, and 960—at both full and half resolutions, as shown in Figure 8 and Figure 9.

However, contrary to this expectation, particularly lightweight models exhibit performance degradation once the D_{max} value surpasses an optimal threshold. Even a recent model like MoCha [3] showed significant vulnerability to this parameter change. This suggests that an excessively large and sparsely supervised disparity search space can introduce matching ambiguities that some architectures struggle to resolve. It is difficult to definitively conclude whether this is an inherent architectural limitation or a domain gap issue arising from the use of pre-trained weights on datasets with different disparity distributions.

As previously noted, our S²M² model is unaffected by the D_{max} parameter. Its consistent performance is therefore shown as a solid line in the figures, serving as a stable reference for comparison.

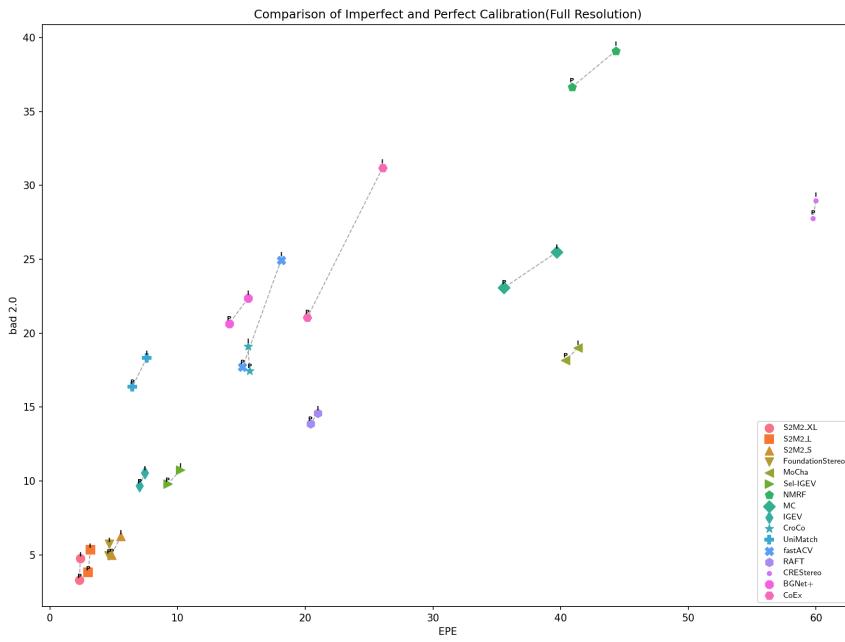


Figure 5. Comparison of model performance under perfect (P) and imperfect (I) calibration settings.

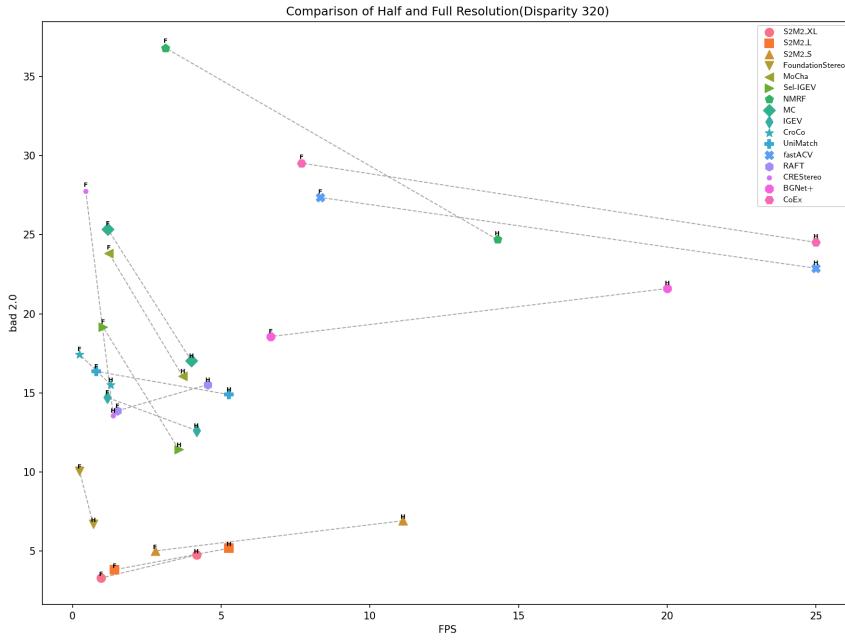


Figure 6. Accuracy vs. inference speed comparison between full (F) and half (H) resolution, with the maximum disparity set to 320. Inference speed was measured with float16 precision using PyTorch's Automatic Mixed Precision (AMP).

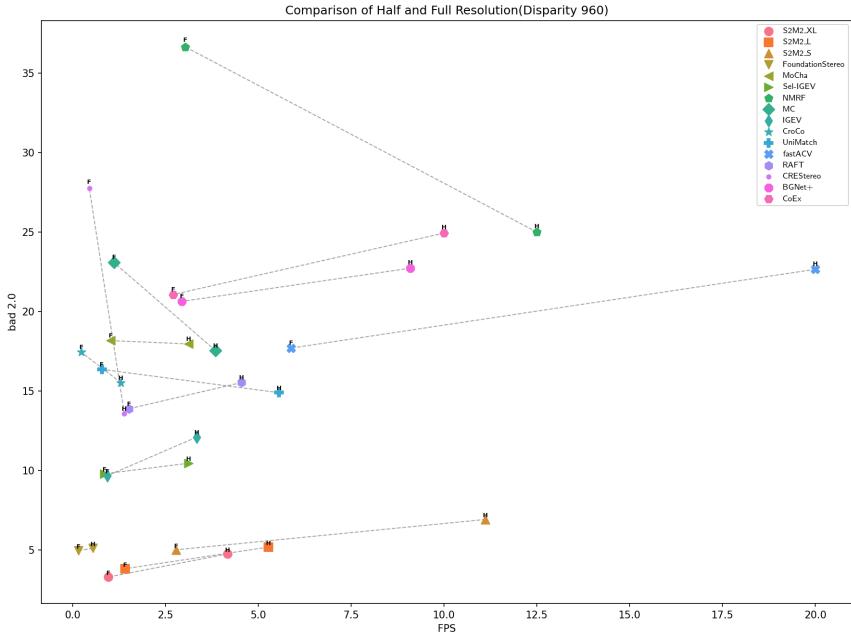


Figure 7. Accuracy vs. inference speed comparison between full (F) and half (H) resolution, with the maximum disparity set to 960. Inference speed was measured with float16 precision using PyTorch’s Automatic Mixed Precision (AMP).

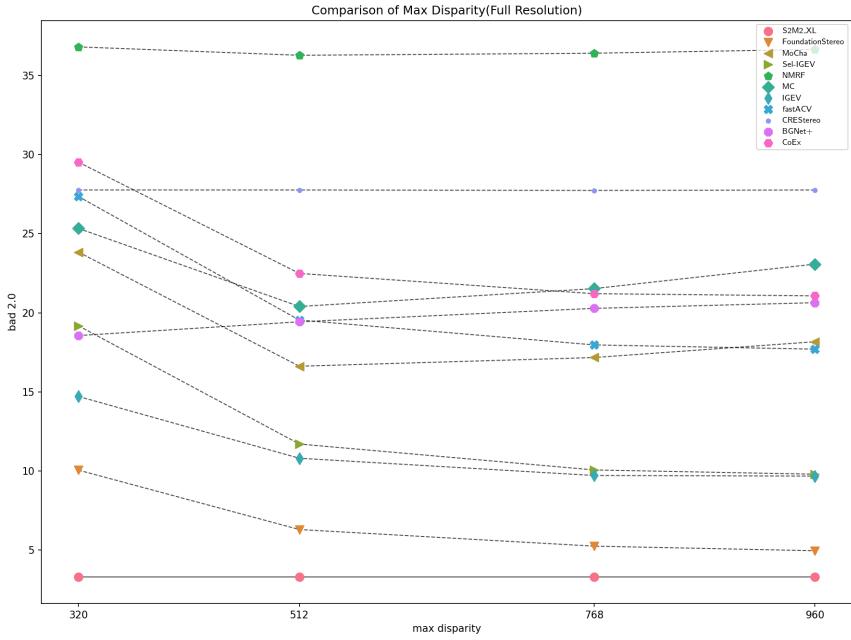


Figure 8. Performance comparison across four maximum disparity (D_{max}) settings at full resolution. The solid line indicates our S²M² model, which is unaffected by this parameter and included for reference.

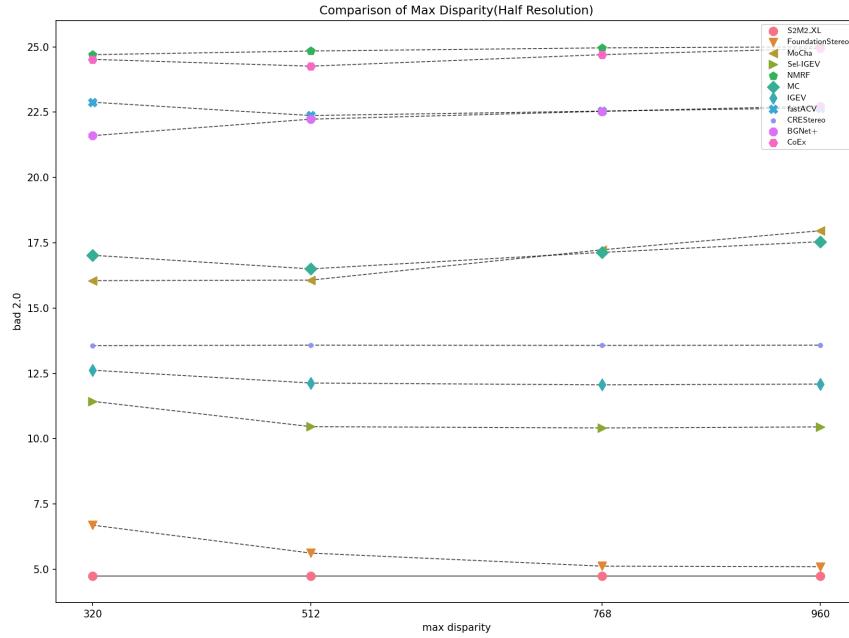


Figure 9. Performance comparison across four maximum disparity (D_{max}) settings at half resolution. The solid line indicates our S²M² model, which is unaffected by this parameter and included for reference.

References

- [1] Blenderkit. <https://www.blenderkit.com/> [Accessed: (2024)].
- [2] Antyanta Bangunharcana, Jae Won Cho, Seokju Lee, In So Kweon, Kyung-Soo Kim, and Soohyun Kim. Correlate-and-excite: Real-time stereo matching via guided cost volume excitation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3542–3548. IEEE, 2021.
- [3] Ziyang Chen, Wei Long, He Yao, Yongjun Zhang, Bingshu Wang, Yongbin Qin, and Jia Wu. Mocha-stereo: Motif channel attention network for stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27768–27777, 2024.
- [4] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13142–13153, 2023.
- [5] Miaoje Feng, Junda Cheng, Hao Jia, Longliang Liu, Gangwei Xu, and Xin Yang. Mc-stereo: Multi-peak lookup and cascade search range for stereo matching. In *International Conference on 3D Vision (3DV)*, pages 344–353. IEEE, 2024.
- [6] Adrien Gaidon, Qiao Wang, Yohann Cabon, and Eleonora Vig. Virtual worlds as proxy for multi-object tracking analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4340–4349, 2016.
- [7] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2012.
- [8] Tongfan Guan, Chen Wang, and Yun-Hui Liu. Neural markov random field for stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2024.
- [9] 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.
- [10] Lahav Lipson, Zachary Teed, and Jia Deng. Raft-stereo: Multilevel recurrent field transforms for stereo matching. In *International Conference on 3D Vision (3DV)*, pages 218–227. IEEE, 2021.
- [11] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4040–4048, 2016.
- [12] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3061–3070, 2015.
- [13] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [14] 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 *Pattern Recognition: 36th German Conference, GCPR 2014, Münster, Germany, September 2-5, 2014, Proceedings 36*, pages 31–42. Springer, 2014.
- [15] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3260–3269, 2017.
- [16] Fabio Tosi, Yiyi Liao, Carolin Schmitt, and Andreas Geiger. Smd-nets: Stereo mixture density networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8942–8952, 2021.
- [17] Jonathan Tremblay, Thang To, and Stan Birchfield. Falling things: A synthetic dataset for 3d object detection and pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 2038–2041, 2018.
- [18] Qiang Wang, Shizhen Zheng, Qingsong Yan, Fei Deng, Kaiyong Zhao, and Xiaowen Chu. Irs: A large naturalistic indoor robotics stereo dataset to train deep models for disparity and surface normal estimation. In *IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2021.
- [19] Wenshan Wang, Delong Zhu, Xiangwei Wang, Yaoyu Hu, Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Sebastian Scherer. Tartanair: A dataset to push the limits of visual slam. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4909–4916. IEEE, 2020.
- [20] Xianqi Wang, Gangwei Xu, Hao Jia, and Xin Yang. Selective-stereo: Adaptive frequency information selection for stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19701–19710, 2024.
- [21] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- [22] Philippe Weinzaepfel, Thomas Lucas, Vincent Leroy, Yohann Cabon, Vaibhav Arora, Romain Brégier, Gabriela Csurka, Leonid Antsfield, Boris Chidlovskii, and Jérôme Revaud. Croco v2: Improved cross-view completion pre-training for stereo matching and optical flow. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17969–17980, 2023.

- [23] Bowen Wen, Matthew Trepte, Joseph Aribido, Jan Kautz, Orazio Gallo, and Stan Birchfield. Foundationstereo: Zero-shot stereo matching. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5249–5260, 2025.
- [24] Bin Xu, Yuhua Xu, Xiaoli Yang, Wei Jia, and Yulan Guo. Bilateral grid learning for stereo matching networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12497–12506, 2021.
- [25] 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.
- [26] Gangwei Xu, Yun Wang, Junda Cheng, Jinhui Tang, and Xin Yang. Accurate and efficient stereo matching via attention concatenation volume. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(4):2461–2474, 2023.
- [27] 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.
- [28] Guorun Yang, Xiao Song, Chaoqin Huang, Zhidong Deng, Jianping Shi, and Bolei 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*, pages 899–908, 2019.
- [29] Jiawei Zhang, Jiahe Li, Lei Huang, Xiaohan Yu, Lin Gu, Jin Zheng, and Xiao Bai. Robust synthetic-to-real transfer for stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20247–20257, 2024.