

779 **7. Appendix**

780 In this supplementary material, we provide detailed data  
781 processing methods and statistical details in Section A. Sec-  
782 tion B elaborates on the parameter settings for SIOU. Ad-  
783 ditionally, Section C presents supplementary experimental  
784 results.

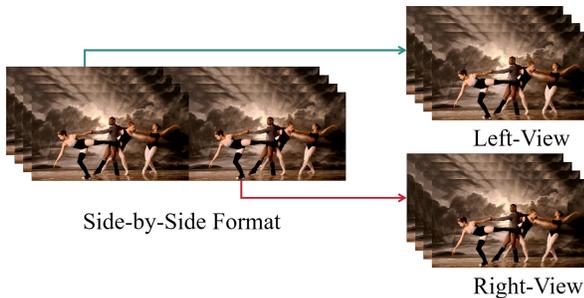
785 **A. Data Curation and Preprocessing.**

Figure 7. Input (Left-View) and ground truth (Right-View) images produced by splitting a video frame.

786 We collect a substantial amount of 3D content in the left-  
787 right format from movies and videos, as illustrated in Fig. 7.  
788 This format necessitates specific viewing equipment, such  
789 as 3D glasses, to ensure that the left and right eyes per-  
790 ceive the corresponding Left-View and Right-View images,  
791 respectively. By dividing these images from the middle,  
792 we create two distinct perspectives: the Left-View and the  
793 Right-View images. Conversion to other stereoscopic for-  
794 mats can be achieved by applying appropriate processing  
795 techniques to this image pair.

```
# define label
texts_complexity = [
    "a simple scene with fewer than three objects",
    "a complex scene with many objects"]
texts_distance = [
    "an indoor scene",
    "an outdoor scene"]
```

Figure 8. CLIP scene categories.

796 For data statistics, we employ CLIP [34] as a scene clas-  
797 sifier. Specifically, we feed text prompts and images into  
798 text encoder and image encoder of CLIP, respectively. We  
799 then calculate the cosine similarity between the resulting  
800 embeddings, assigning the category with the highest simi-  
801 larity as the classification result. Fig. 8 showcases the spe-  
802 cific text prompts used. We perform pairwise statistics for  
803 the four categories (indoor, outdoor, simple, and complex).  
804 Additionally, we analyze the scene distribution within the  
805 dataset. As illustrated in Fig. 9, Mono2Stereo encompasses

common indoor environments like living rooms and bed-  
rooms, as well as more unique settings such as underwa-  
ter scenes, cliffs, and rivers. For overall scene category sta-  
tistics in Fig. 9, we utilize prompts in the format of “a/an  
[category] scene.”

806  
807  
808  
809  
810

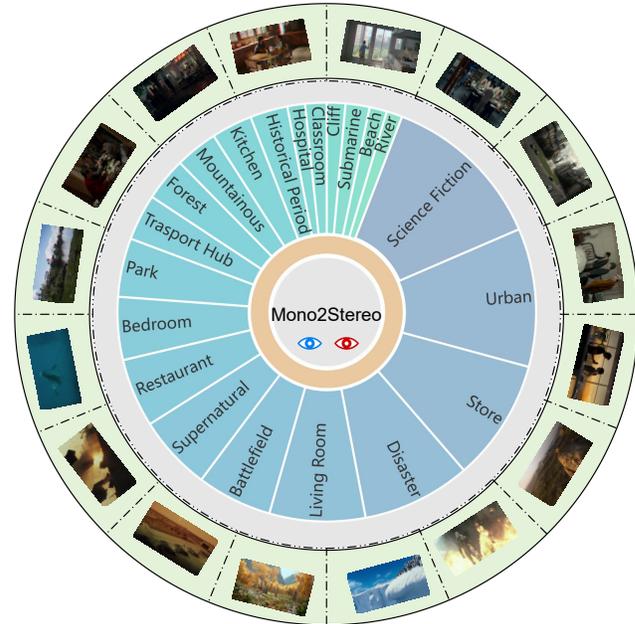


Figure 9. Distribution Characteristics of the Mono2Stereo Dataset.

**B. Parameter Settings for SIOU.**

811

812 To illustrate the individual roles of the two terms within  
813 SIOU, we conduct separate human subjective evaluations,  
814 with the results presented in Tab. 8. As shown, both terms  
815 contribute to achieving good consistency, suggesting that  
816 each reflects stereo quality to a certain extent by primar-  
817 ily focusing on the true disparity regions between the Left-  
818 View and Right-View images. Regarding the balancing pa-  
819 rameter  $\alpha$  in Eq. (1), we randomly divide 1100 image pairs  
820 into two sets: 500 pairs for optimal parameter and threshold  
821 searching, and 600 pairs for generalization validation. We  
822 experiment with various parameter settings for  $\alpha$ , includ-  
823 ing 0.25, 0.5, 0.7, 0.75, and 0.8. Our findings indicate that  
824 these settings yield better consistency compared to a single  
825 item. Notably,  $\alpha = 0.75$  demonstrates the highest level of  
826 consistency. Therefore, we set  $\alpha$  to 0.75 for final SIOU. Val-  
827 idation on a set of 600 pairs, as illustrated in Tab. 8, shows  
828 no significant signs of overfitting.

829 For IoU2, employing a lower threshold ensures greater  
830 sensitivity to discrepancies, encompassing both disparity  
831 and pixel shifts. When validating across 500 sample pairs,  
832 we observe that a threshold of 5 yields the highest consis-  
833 tency. Attempting to decrease this threshold further actually

Table 8. Correlation with human judgements of stereo quality. IoU1 and IoU2 are components of the proposed SIoU metric. Both demonstrate correlation with human perception. Combining these components into SIoU yields even higher correlation scores. The results are based on a validation set of 600 pairs.

Metric	SIoU	IoU1	IoU2
Spearman Rank	<b>0.84</b>	0.81	0.80
Kendall Rank	<b>0.73</b>	0.70	0.68

```

1  def compute SIoU(left, right, pred):
2
3      # Convert RGB images into gray scale
4      left_gray = RGB2Gray(left)
5      right_gray = RGB2Gray(right)
6      pred_gray = RGB2Gray(pred)
7
8      # detect edges of the right and generated images
9      right_edges = Canny(right_gray)
10     pred_edges = Canny(pred_gray)
11
12     # compute the differences
13     diff_r1 = abs(right_gray - left_gray)
14     diff_g1 = abs(pred_gray - left_gray)
15     logical_r1 = Zeroslike(diff_r1.shape)
16     logical_g1 = Zeroslike(diff_g1.shape)
17     logical_r1[diff_r1 > 5] = 1
18     logical_g1[diff_g1 > 5] = 1
19
20     IoU1 = IoU(pred_edges, right_edges)
21     IoU2 = IoU(logical_g1, logical_r1)
22     SIoU = 0.75 * IoU1 + 0.25 * IoU2
23
24     return SIoU

```

Figure 10. Pseudocode for the SIoU calculation process.

834 reduces consistency. This occurs because a lower threshold  
835 ( $< 5$ ) incorporates more pixels into consideration, includ-  
836 ing those in areas that do not significantly impact the stereo  
837 effect, which is undesirable.

## 838 C. Supplementary Experimental Results.

### 839 C.1. Detailed Analysis of Two Conditions

840 In this paper, we define the complete Left-View image as the  
841 geometric condition, while the warped version of the Left-  
842 View image serves as the viewpoint condition. These condi-  
843 tions correspond to the inputs of single-stage and two-stage  
844 models, respectively. This section provides further clarifi-  
845 cation. As depicted in Fig. 11, the Left-View is a complete  
846 natural image, offering comprehensive geometric structure  
847 and texture details. Conversely, the Warped image, derived  
848 from the Left-View image through disparity warping, ex-

hibits a perspective closer to the Right-View image. There-  
849 fore, the Left-View image provides richer geometric infor-  
850 mation, while the Warped image explicitly offers an obser-  
851 vational viewpoint, spatially aligning it closer to the target.  
852 This distinction forms the basis for our naming convention  
853 and motivates our design of the dual-condition model, lever-  
854 aging the complementary strengths of both conditions. 855

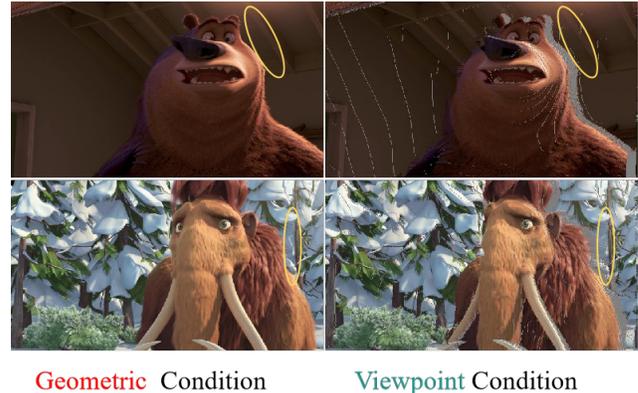


Figure 11. Visualization of Dual-Condition. The yellow circles highlight the differences in key spatial relationships, while the gray areas represent geometric differences.

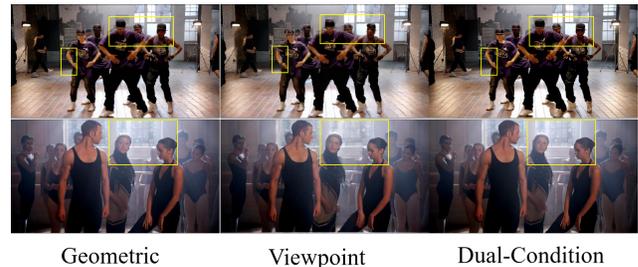


Figure 12. The influence of identical conditions on the output results. Areas with significant differences are highlighted by yellow boxes.

856 Furthermore, we present results under three different  
857 conditions, as illustrated in Fig. 12. Both the “Geometric”  
858 and “Viewpoint” conditions exhibit artifacts to varying de-  
859 grees, with the “Viewpoint” condition displaying more pro-  
860 nounced artifacts due to its partially occluded input. In con-  
861 trast, the “Dual-Condition” yields superior image quality.

### 862 C.2. Evaluating Performance in Various Scenes

863 To gain a deeper understanding of the performance across  
864 different scenarios, we evaluate models separately on five  
865 distinct scenes from the Mono2Stereo test dataset. As  
866 shown in Tab. 11, we observe that the model struggles in  
867 pairwise comparisons involving indoor, complex, and ani-  
868 mation scenes. We hypothesize that this is due to limita-  
869 tions in three key areas where the model requires further

870 improvement: disparity range estimation accuracy, geomet-  
871 ric understanding, and color distribution handling. Conse-  
872 quently, we suggest that future research should focus on ad-  
873 dressing these aspects. Finally, Mono2Stereo also provides  
874 20 video clips for evaluating models. Despite our method  
875 being single-frame based, it still achieves promising results.

### 876 C.3. Why Velocity Edges?

877 Regarding the edge consistency constraint, the most intu-  
878 itive approach appears to be constraining the edges within  
879 the latent space. Visualization of the feature maps, as illus-  
880 trated in Fig. 13, confirms that both the latent and velocity  
881 exhibit positional correlation with the image. However, dur-  
882 ing training, we observe that predicting the latent or noise  
883 results in significantly slower convergence and even opti-  
884 mization failure, while velocity prediction does not suffer  
885 from these issues. Consequently, we opt to constrain the  
886 edges of the velocity field.

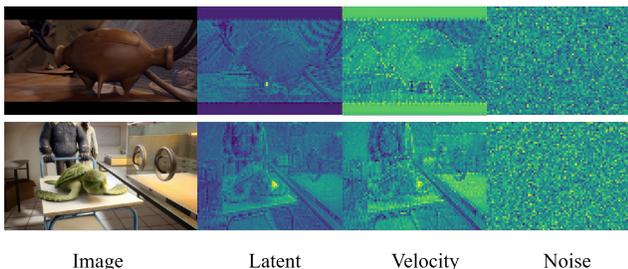


Figure 13. Feature maps of latent and velocity.

### 887 C.4. Ablation Study on Inria 3DMovie.

888 To further validate the effectiveness of the Edge Consis-  
889 tency loss, we conduct out-of-domain performance evalu-  
890 ations using the Inria 3DMovie dataset, which comprises  
891 2,727 stereoscopic image pairs. As shown in Tab. 9, incor-  
892 porating the Edge Consistency constraint consistently im-  
893 proves performance across all three tested conditions. This  
894 suggests that the benefits of this constraint are not limited  
895 to specific datasets, demonstrating its potential for general-  
896 ization.

### 897 C.5. Ablation Study on Edge Consistency Loss

898 When applying the Edge Consistency loss, we conduct ex-  
899 periments to validate the impact of different  $\alpha$  values in  
900 Eq. (3) within a small range. Using the dual-condition dif-  
901 fusion model, we experiment with  $\alpha$  values of 0.75, 1, and  
902 1.25, while  $\alpha = 0$  represents the absence of the edge con-  
903 straint. As Tab. 10 illustrates, applying the edge consis-  
904 tency constraint at varying strengths consistently leads to  
905 improvements in SIOU, indicating that the constraint term  
906 is not overly sensitive to the specific  $\alpha$  value. We offer an  
907 additional analysis: when  $\alpha$  is 0, all pixels in the image are

Table 9. Impact of LEC Loss across three conditions on Inria 3DMovie dataset.

Geo.	View.	LEC Loss	Inria 3DMovie			
			SIOU $\uparrow$	RMSE $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$
✓			0.2836	7.47	30.66	0.693
✓		✓	0.2949	7.46	30.68	0.693
	✓		0.3147	7.61	30.50	0.678
	✓	✓	0.3145	7.52	30.61	0.684
✓	✓		0.3147	7.44	30.70	0.691
✓	✓	✓	0.3186	7.31	30.85	0.697

Table 10. Impact of EC loss across three conditions. EC loss consistently improves performance, with notable gains in SIOU, the metric for perceived stereo quality.

LEC Loss	Mono2Stereo			
	SIOU $\uparrow$	RMSE $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$
0	0.2588	6.90	31.35	0.721
0.75	0.2608	6.83	31.45	<b>0.725</b>
1	<b>0.2619</b>	<b>6.82</b>	<b>31.45</b>	0.721
1.25	0.2615	6.88	31.38	0.719

optimized equally. The edge constraint, in essence, imposes  
a stricter penalty on regions that genuinely influence

908  
909

Table 11. Evaluating models across various scenes.

Method	Indoor		Outdoor		Complex		Simple		Animation		Video	
	SIoU $\uparrow$	RMSE $\downarrow$										
StereoDiffusion [46]	0.2387	7.48	0.2441	7.68	0.2182	7.78	0.2571	6.17	0.2296	8.01	0.1992	8.38
Geometric Condition	0.2505	5.31	0.2543	5.74	0.2561	5.94	0.2791	4.28	0.2525	5.73	0.2610	5.61
Viewpoint Condition	0.2761	5.71	0.2824	6.02	0.2713	6.62	0.2986	5.76	0.2764	6.49	0.2735	5.95
Dual Condition	0.2819	5.21	0.2969	5.65	0.2894	5.78	0.3095	4.29	0.2999	5.76	0.2817	5.50