

DMS: Diffusion-Based Multi-Baseline Stereo Generation for Improving Self-Supervised Depth Estimation

Zihua Liu<sup>1</sup>, Yizhou Li<sup>2</sup>, Songyan Zhang<sup>3</sup> and Masatoshi Okutomi<sup>1</sup>



SONY



1.Institute of Science Tokyo 2. Sony Semiconductor Solutions Group 3.Nanyang Technological University

## Method Overview

Limitations of Existing Methods

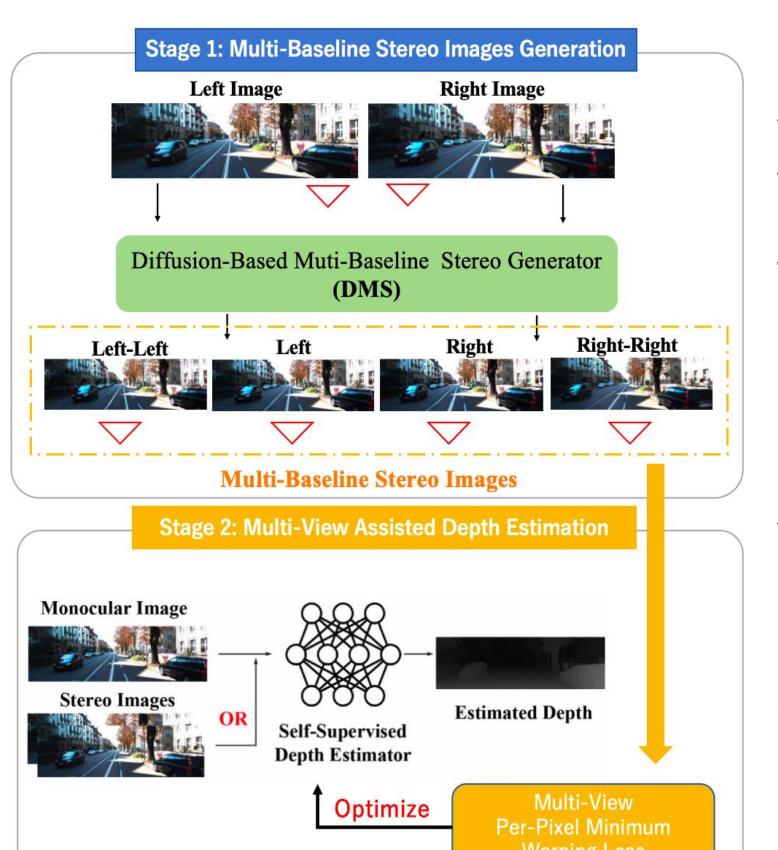
In self-supervised depth estimation, the photometric warping loss is

bounded by the two-views information, yielding only limited gains

**Occlusions** 

Out of Views

Two-Stage Training Pipeline.



in occluded and out-of-view regions.

Self-Supervise

Stage 1: Start from given views, using diffusion model to generated multi-baseline images.

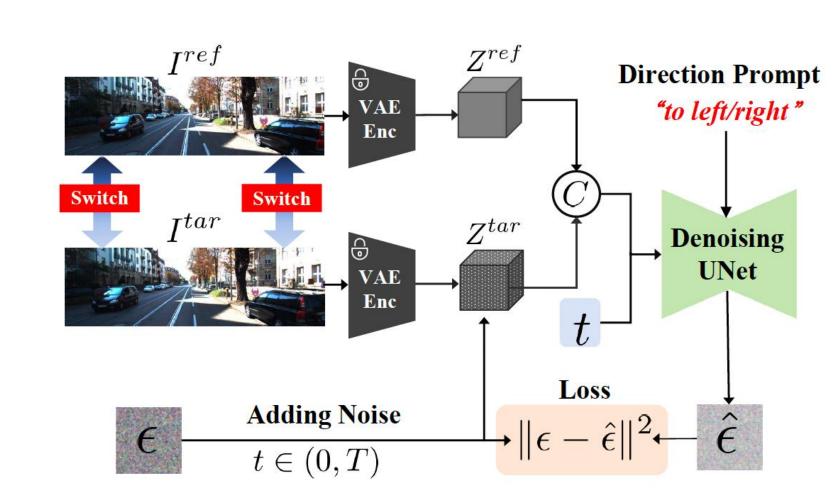
Stage 2: Use multibaseline images to provide extra geometry clues to assist the selfsupervised depth estimation.

Occlusion Map

Out of View Map

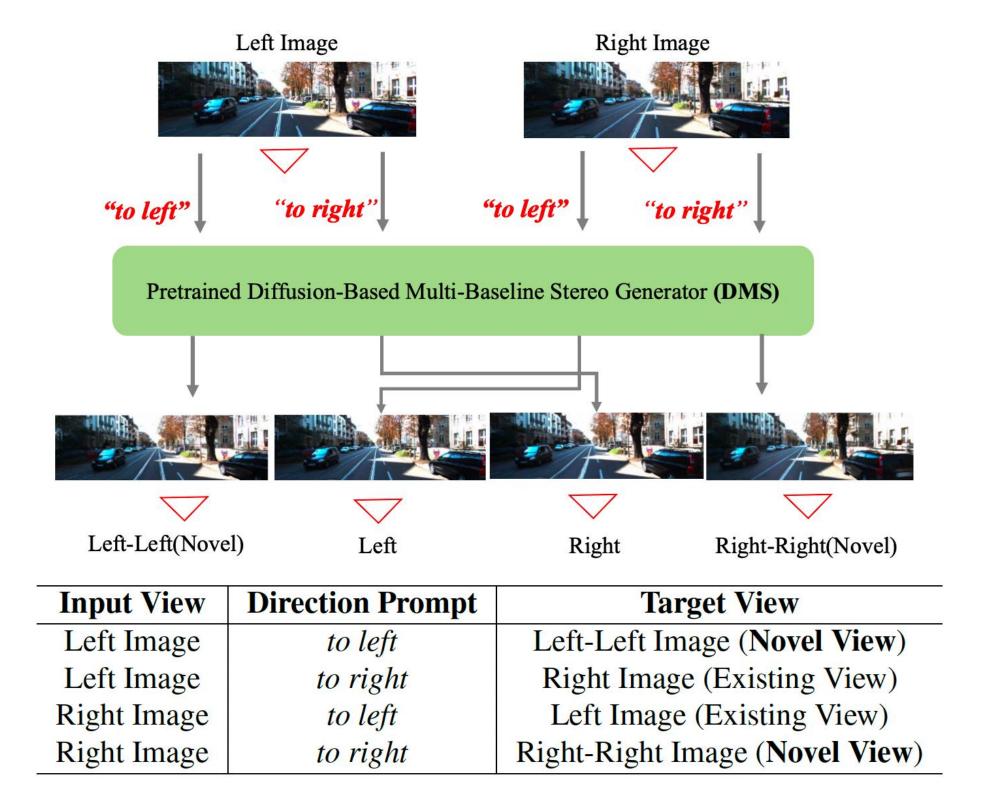
#### DMS: Diffusion-Based Multi-Baseline Stereo Generation

Training Phase of the DMS



- (a) Training Phase of the DMS
- Finetune Stable Diffusion to produce the missing stereo view from a reference image, either left→right or right→left.
- Guide the model with directional prompts—'to left' or 'to right'—to indicate view offsets

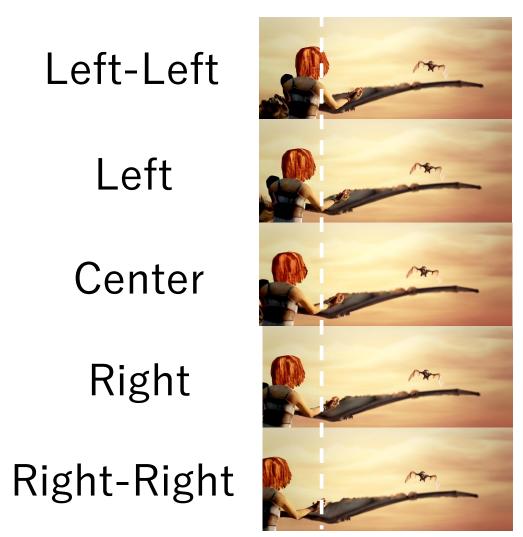
Inference Phase of the DMS



 At inference stage, extend the baseline by prompting: left + 'to left' → left-left, right + 'to right' → right-right.

## **Experimental Results**

Multi-Baseline Stereo Images Generation.





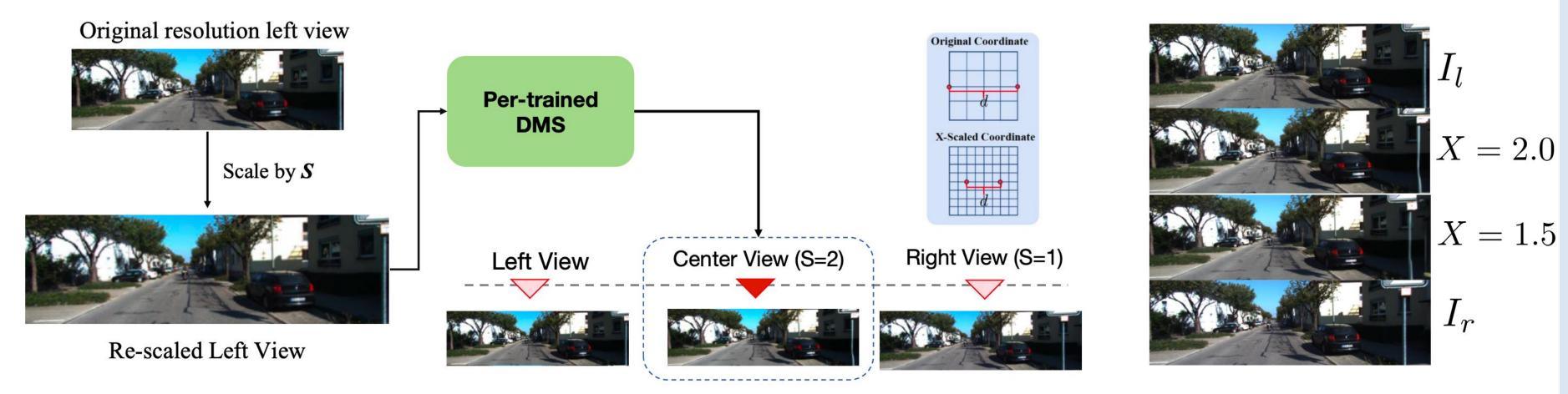
\* The white line marking the same horizontal position.

Improving Self-Supervised Stereo Matching

<u>(-</u>	SceneFlow[47]						KITTI 2015 [48]					
Method	EPE↓			>3px%↓			EPE↓			>3px%↓		
	All	Occ	Oof	All	Occ	Oof	All	Occ	Oof	All	Occ	Oof
Baseline	4.09	22.66	10.83	13.5	83.6	42.8	1.48	4.38	9.26	7.7	39.6	64.2
+ ll + rr	2.45	11.45	6.16	9.0	51.0	25.0	1.34	3.83	7.82	6.5	34.4	42.8
+ <i>c</i>	3.75	21.70	7.73	12.5	81.3	38.4	1.36	4.14	7.64	6.7	37.0	49.6
+ ll + rr + c	2.32	11.15	5.57	8.4	49.2	23.4	1.24	3.56	7.31	5.8	32.3	39.7
Left	Image	,	GTD	isparit	ty	Е	Baselii	1e	Ва	aselii	ne+DI	<i>NS</i>
Occlus Occlus	ion M	Tap Oi	ıt-of-F	rame .	Map	En	ror M	lap		Erro	or Maj	
KITTI 2015											652 (2)	

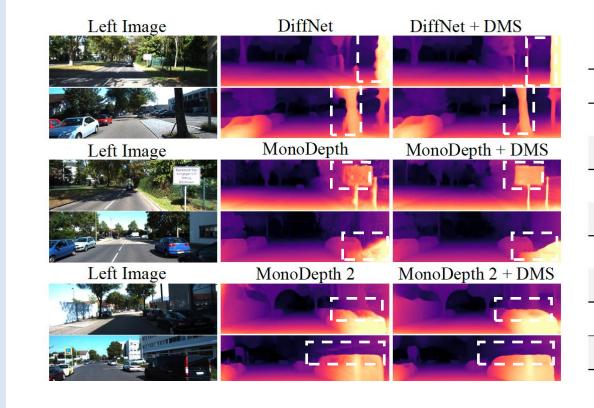
	D1-bg	D1-fg	D1-all					
Non-Learning Based Method								
SGM <sup>†</sup> [28]	8.95	20.55	10.88					
$SGM^{\dagger} + DMS$	<b>7.96(11.1%</b> ↓)	$16.68 (18.8\%\downarrow)$	9.41(13.5% \_)					
.]	Learning-Based	Methods						
Zhou <i>et.al</i> [95]	-	=	9.91					
SegStereo [82]	-	-	8.79					
OASM[36]	6.89	19.42	8.98					
Flow2Stereo[41]	5.01	14.62	6.61					
PASMnet[71]	5.41	16.36	7.23					
PASMnet + DMS	5.24(3.1%↓)	$13.96(14.7\%\downarrow)$	6.69(7.5%↓)					
StereoNet*[34]	7.31	17.77	9.05					
StereoNet* + DMS	4.68(36.0%↓)	$12.06(32.1\%\downarrow)$	<b>5.91(34.7%</b> ↓					
CFNet* [59]	7.22	18.54	9.11					
CFNet* + DMS	4.64(35.7%↓)	12.33(33.5%\( )	<b>5.92</b> (35.0%↓)					
RaftStereo* [39]	3.38	13.62	5.08					
RaftStereo* + DMS	2.95(12.7%↓)	<b>6.88</b> (50.9%↓)	3.60(29.1%\1)					
IGEVStereo* [80]	3.76	11.14	4.98					
IGEVStereo* + DMS	2.80(25.5%↓)	6.37(45.27%↓)	3.40(31.72%					
MCStereo* [18]	3.01	13.38	4.73					
MCStereo* + DMS	2.83(6.0%↓)	<b>7.03(47.46%</b> ↓)	3.67(22.41%)					

# Intermediate View Approximation with Scaling Operations



We observe that scaling the reference view makes the view shift inversely proportional to the original baseline ( $\mathbf{d}' \approx \mathbf{d}/\mathbf{scale}$ ), enabling an approximate estimate of intermediate views

#### Improving Self-Supervised Monocular Depth Estimation



M-41 J	Ala Dall	C. D.II	DMCE	DMCEL
Method	Abs Rel↓	Sq Rel↓	RMSE↓	RMSE_log↓
MonoDepth [21]	0.120	1.041	5.272	0.217
MonoDepth + DMS	0.109	0.860	5.004	0.202
MonoDepth2 [22]	0.109	0.873	4.960	0.209
MonoDepth2 + DMS	0.105	0.811	4.850	0.200
SDFANet* [97]	0.104	0.997	4.583	0.186
SDFANet* + DMS	0.097	0.643	4.218	0.181
DiffNet [96]	0.104	0.809	4.766	0.201
DiffNet + DMS [96]	0.098	0.726	4.606	0.191