

# Depth Learning: When Depth Estimation Meets Deep Learning

Prof. Liang LIN

SenseTime Research & Sun Yat-sen University

# Outline

## Part 1. Introduction

- Motivation
- Stereo Matching
- Single Image Depth Estimation

## Part 2. Depth Learning—Our Progress

- Cascade Residual Learning (CRL)
- Zoom and Learn (ZOLE)
- Single View Stereo Matching (SVS)

## Part 3. Conclusion

# Outline

## Part 1. Introduction

- Motivation
- Stereo Matching
- Single Image Depth Estimation

## Part 2. Depth Learning—Our Progress

- Cascade Residual Learning (CRL)
- Zoom and Learn (ZOLE)
- Single View Stereo Matching (SVS)

## Part 3. Conclusion

# Motivation—Importance of Depth



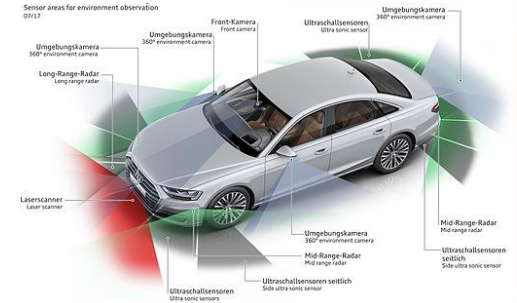
Depth-of-field Rendering (Bokeh)



- Face ID



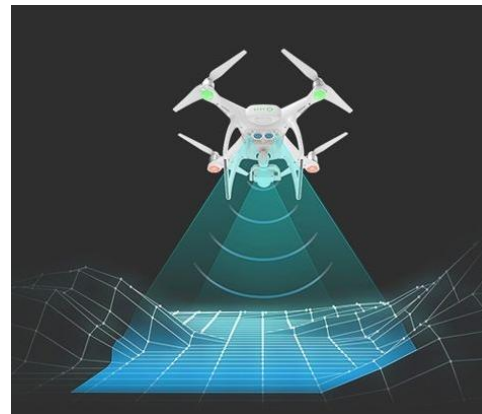
AR / VR



ADAS



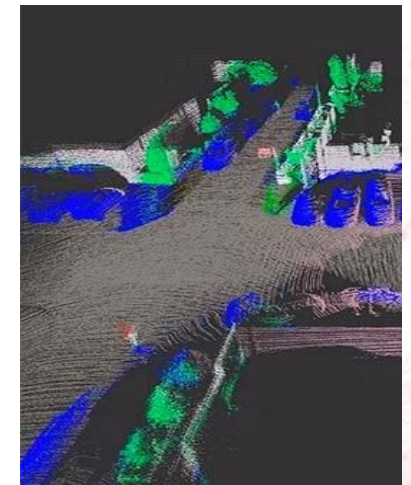
Industrial Automation



Obstacle avoidance



Entertainment



- 3D Reconstruction

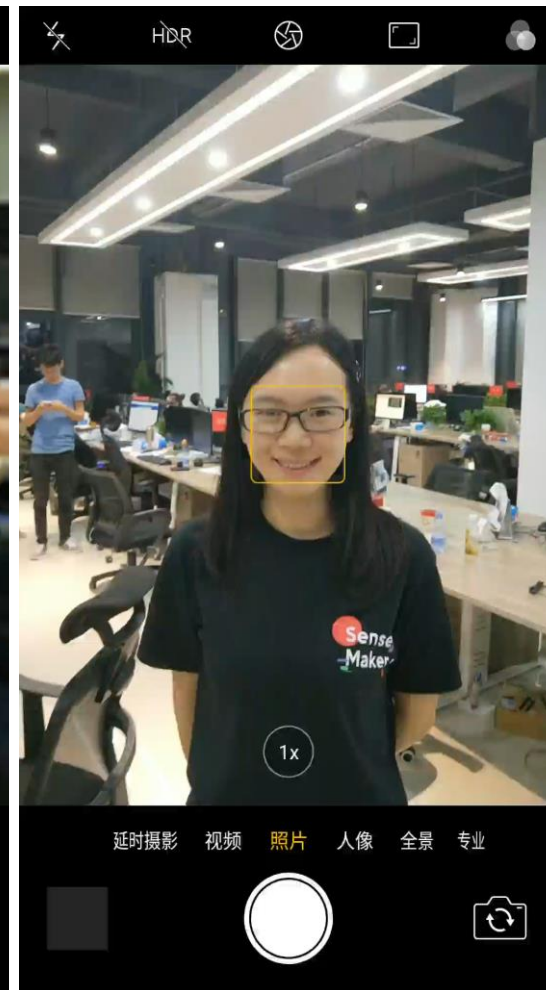
# Application: Bokeh in mobile phone



景物虚化效果  
Object Bokeh



视频重对焦  
Video Refocus

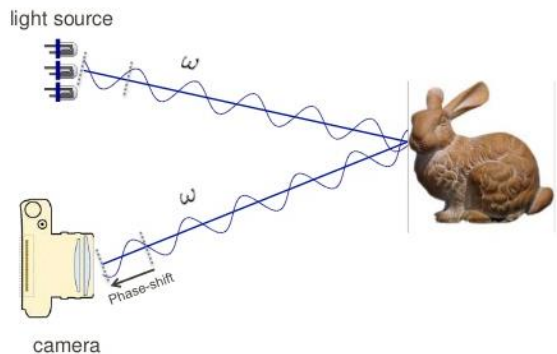


人脸对焦跟随  
Face Tracking

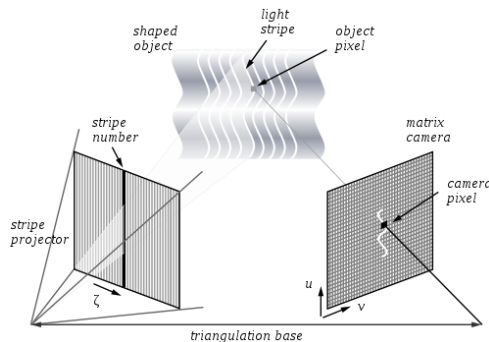
# Applications: AR & Vide Editing



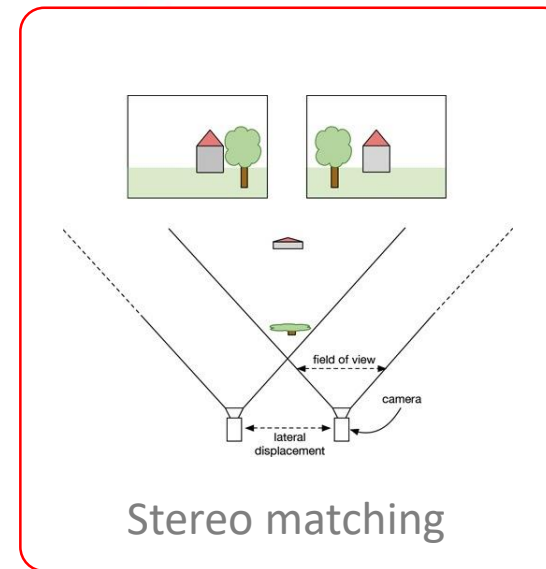
## Multiple ways to estimate depth



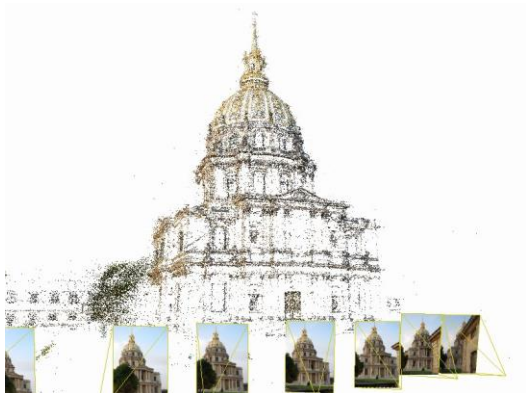
Time of flight



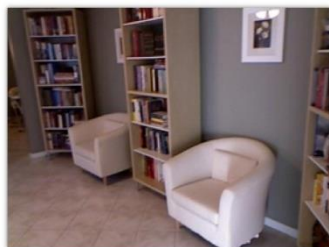
Structure light



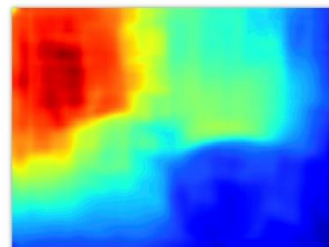
Stereo matching



Structure from motion



Single RGB Image



Depth Map

Single Image Depth Estimation

Most **cost-effective** approaches  
Main focus of this talk

# Outline

## Part 1. Introduction

- Motivation
- Stereo Matching
- Single Image Depth Estimation

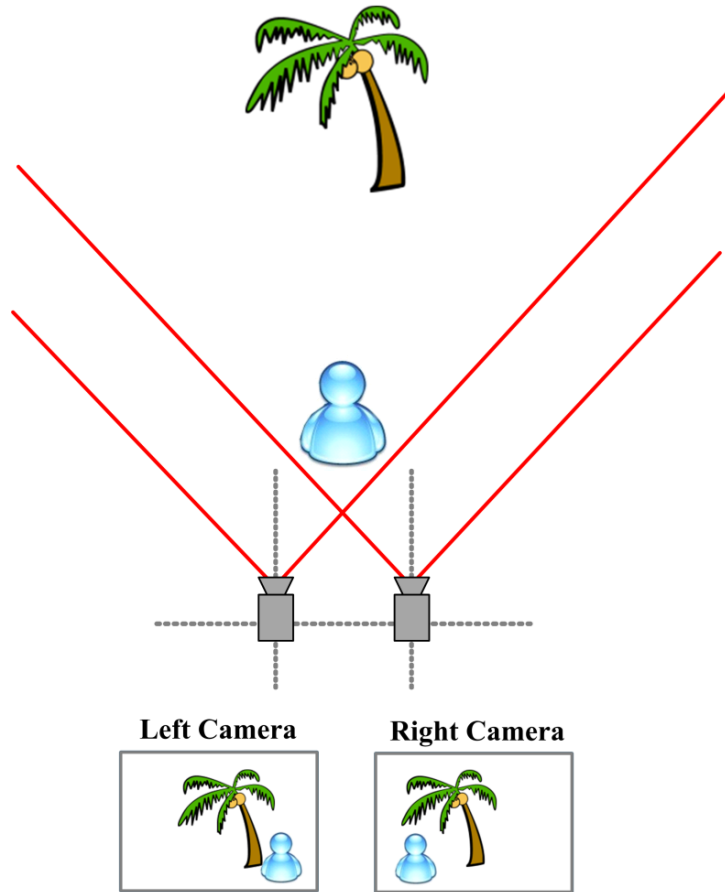
## Part 2. Depth Learning—Our Progress

- Cascade Residual Learning (CRL)
- Zoom and Learn (ZOLE)
- Single View Stereo Matching (SVS)

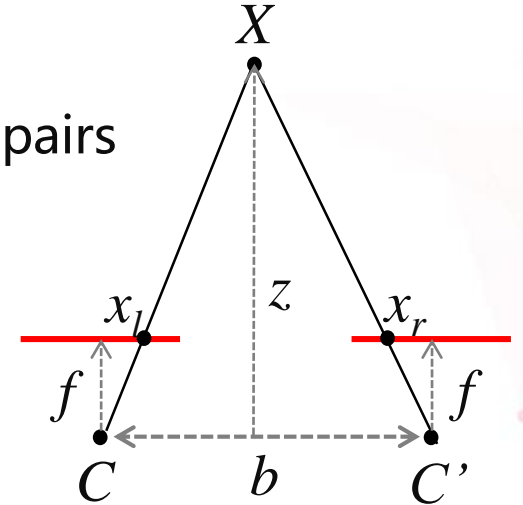
## Part 3. Conclusion



# Stereo Matching



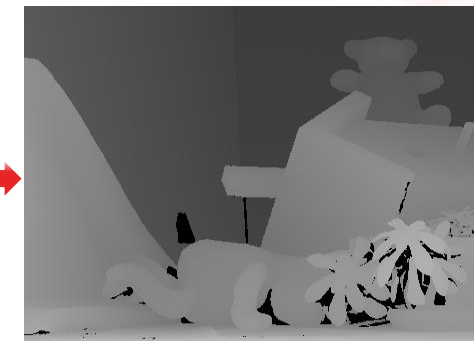
- Estimate **pixel correspondence** in rectified pairs
- Disparity Map  $D(x, y) = x_r - x_l$
- Depth Map  $Z(x, y) = \frac{f \cdot b}{D(x, y)}$



Left image



Right image



Disparity Map

# Stereo Matching—Challenges



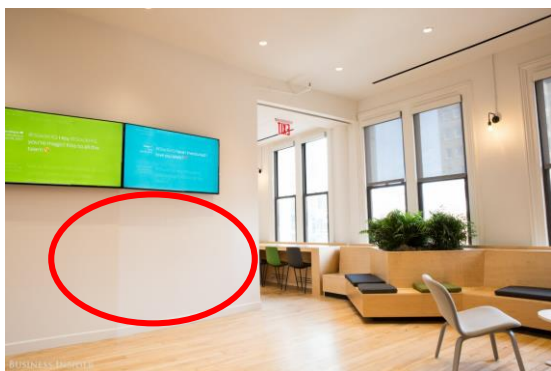
Underdetermined (ill-posed)



Photometric variations



Occlusions



Texture-less Areas

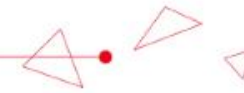


Repetitive patterns



Reflections

# Stereo Matching—Classic Pipeline

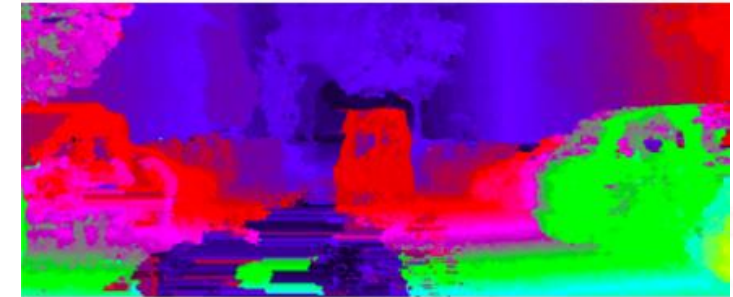


Typical method: Semi-Global Matching (SGM)

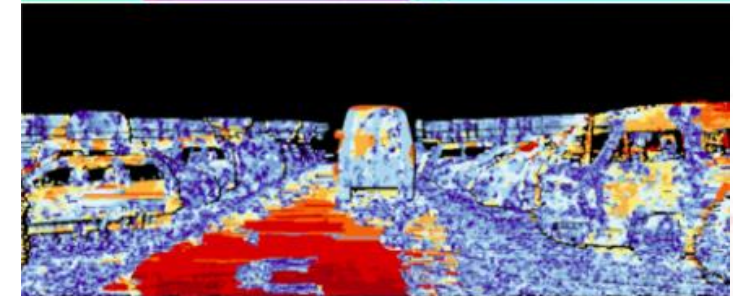
- Matching cost computation
  - Compute 3D cost volume along horizontal scanline
- Cost Aggregation
  - Refine matching cost by aggregating information of neighboring pixels
- Disparity Computation
  - Derive the disparity from the matching cost
- Disparity Refinement
  - Post-processing on the disparity map



Left img



SGBM



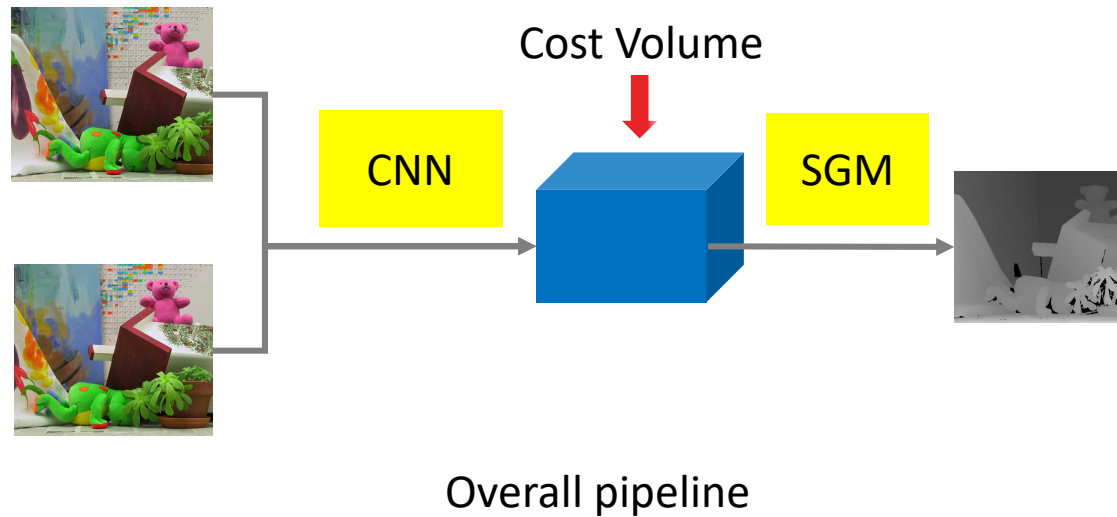
Error

**Erroneous at ill-posed region!**

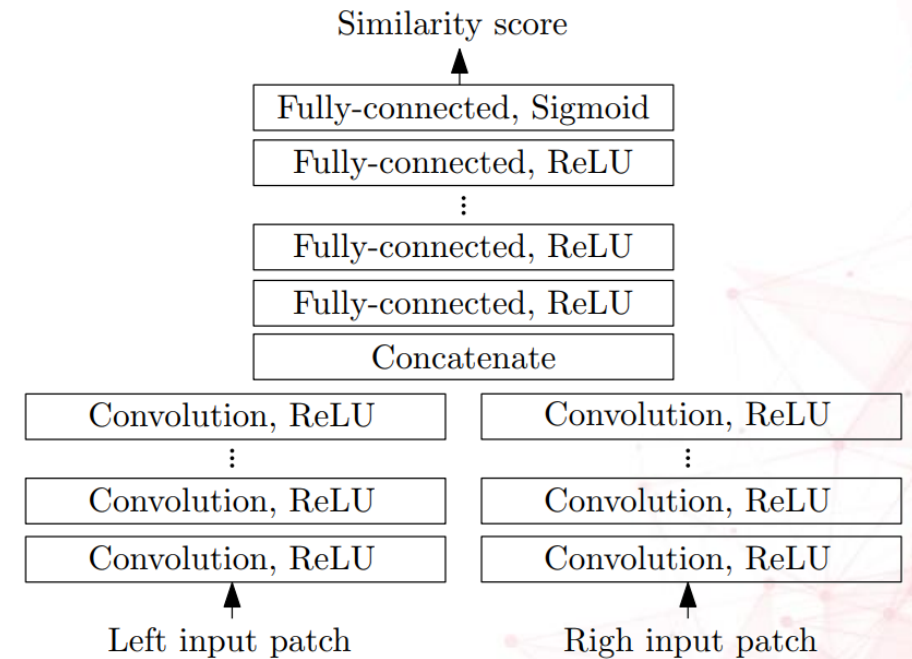
[2] H. Hirschmuller, "Stereo processing by semiglobal matching and mutual information", *IEEE Trans. on pattern analysis and machine intelligence*, 2008.

## Matching cost learning

- Train a model to classify **patches** into two classes (similar and not similar)
- A **small set of image pairs** with ground-truth disparities generate millions of patches
- Depend on the performance of **SGM**



Overall pipeline



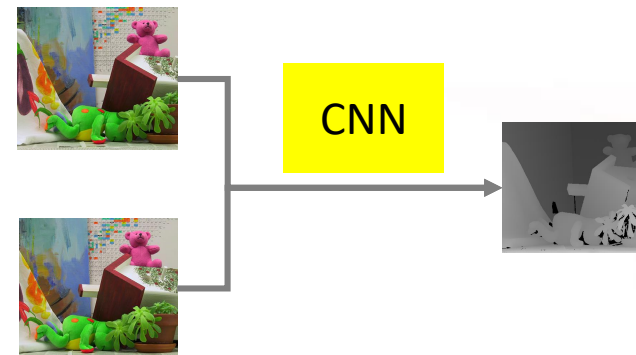
CNN Architecture of [3]

[3] J. Zbontar and Y. LeCun. Stereo matching by training a convolutional neural network to compare image patches. *Journal of Machine Learning Research*, 2016.

# Stereo Matching Meets Deep Learning (II)

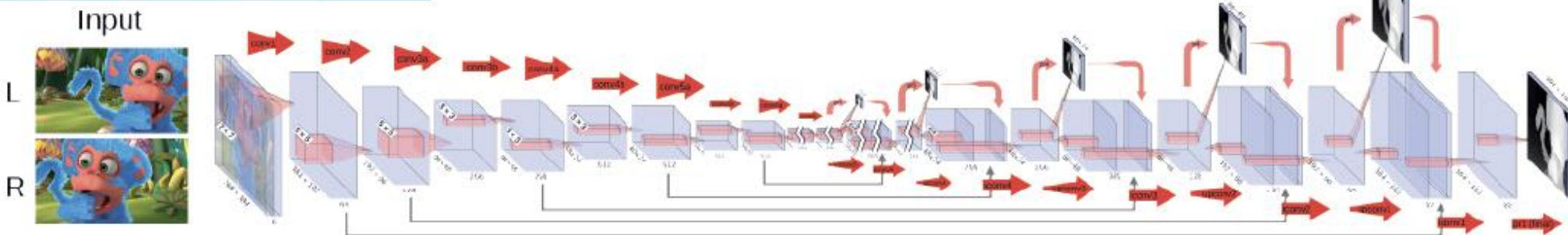
## End-to-end Learning

- Train an end-to-end model to **regress disparity**
- A large-scale dataset is needed to train a good model
- Usually via **hourglass** structure
- Variations:
  - **Correlation** layer to compute cost volume, e.g., DispNetC [4]
  - **Unsupervised** learning with left-right check, e.g., [5]



Overall pipeline

Disparity Network Architecture



[4] N. Mayer, et al., "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation." In *Proc. IEEE CVPR*, 2016.

[5] C. Zhou, H. Zhang, X. Shen, and J. Jia. "Unsupervised Learning of Stereo Matching." In *Proc. IEEE CVPR*, 2017.

# Outline

## Part 1. Introduction

- Motivation
- Stereo Matching
- Single Image Depth Estimation

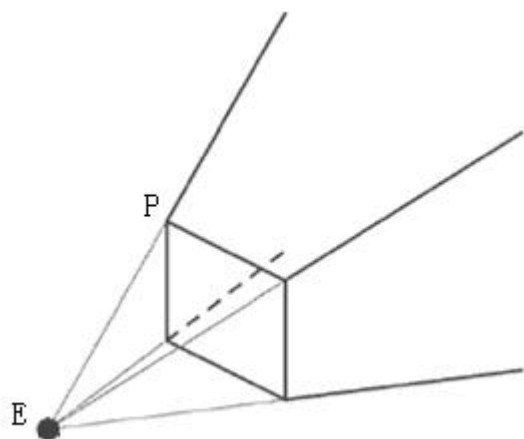
## Part 2. Depth Learning—Our Progress

- Cascade Residual Learning (CRL)
- Zoom and Learn (ZOLE)
- Single View Stereo Matching (SVS)

## Part 3. Conclusion

## Estimate depth from a single image

- Monocular depth estimation:  $D = F(I)$
- Highly-ill posed, infinite configurations
- **Deep learning** to rescue!

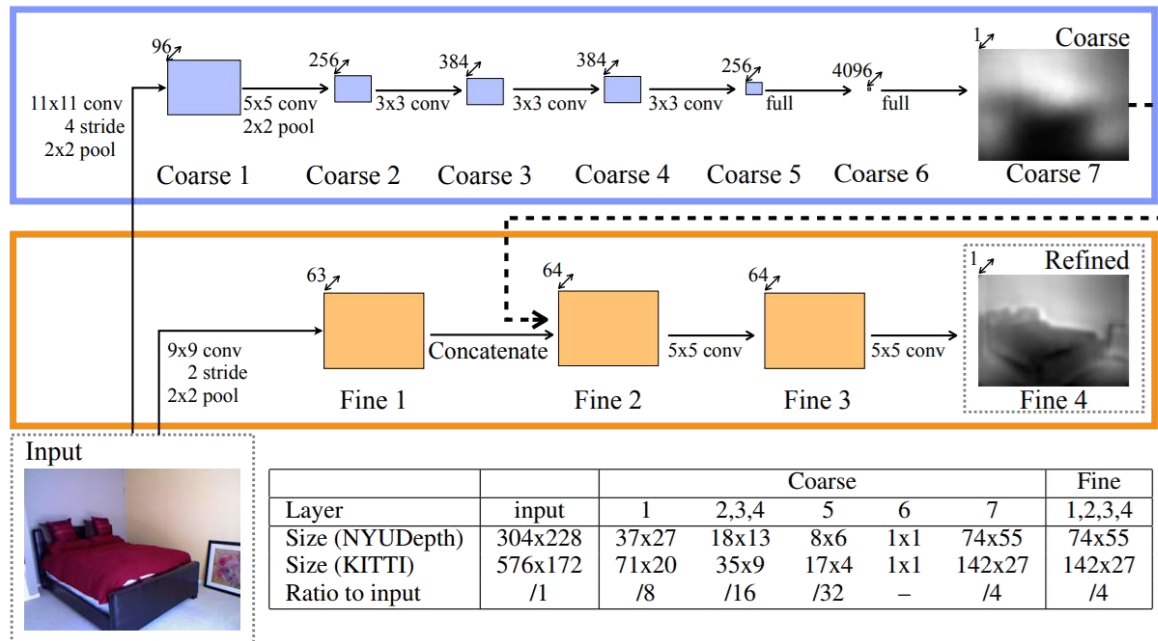


Needs semantic info to understand the scene



## Direct regression with CNN

- Train a model end-to-end to **regress** scene depth directly, e.g., [6]
- Issue: generalization, requires lots of data

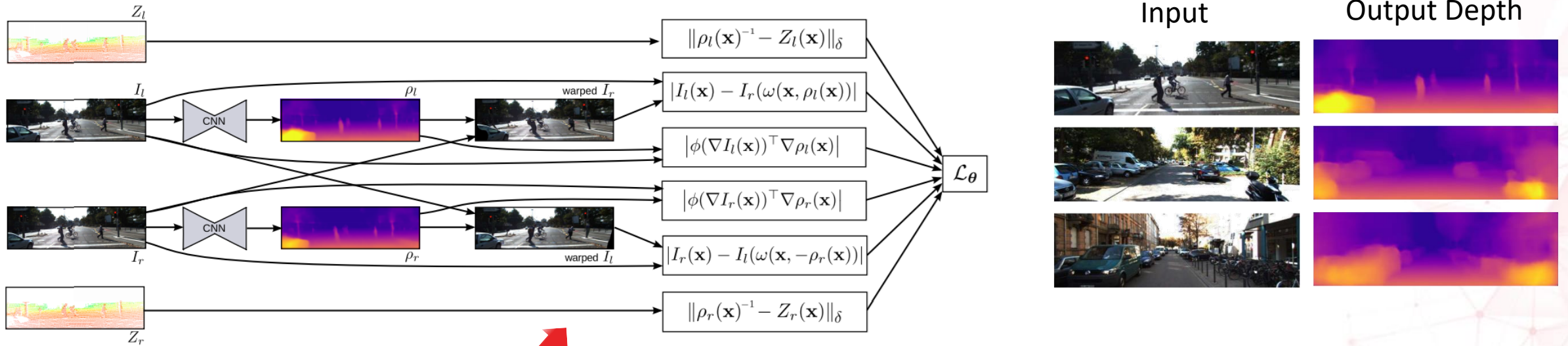


[6] D. Eigen, C. Puhrsch, R. Fergus "Depth map prediction from a single image using a multi-scale deep network." In *Proc. NIPS*, pp. 2366-2374, 2016.



## Semi-supervised/Unsupervised schemes

- Use stereo pairs for training
- Based on left-right consistency and smoothness of depth



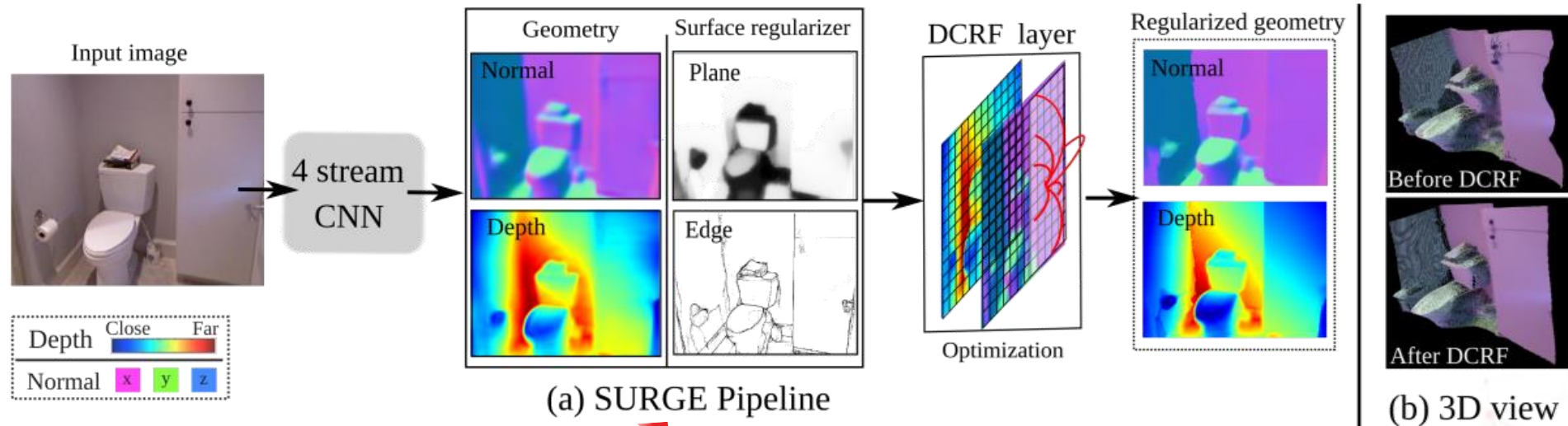
[8] defines a series of loss functions based on left-right consistency and depth smoothness

[7] C. Godard, O. Mac Aodha, G.J. Brostow "Unsupervised monocular depth estimation with left-right consistency," In *Proc. CVPR*, 2017.

[8] Y Kuznetsov, J Stückler and L. Bastian, "Semi-supervised deep learning for monocular depth map prediction," In *Proc. CVPR*, 2017.

## CNN with expressive modeling

- Integrate geometric models (e.g., [9]) or image models (e.g., CRF, [10]) into CNN



[9] classifies pixels into planes and edges, enforce two points on the same plane have same normal

[9] P. Wang, et al., "SURGE: surface regularized geometry estimation from a single image," In *Proc. NIPS*, pp. 172-180, 2016.

[10] F. Liu, et al., "Deep convolutional neural fields for depth estimation from a single image," In *Proc. CVPR*, pp. 5162-5170, 2015.

## Part 2. Depth Learning—Our Progress

# Outline

## Part 1. Introduction

- Motivation
- Stereo Matching
- Single Image Depth Estimation

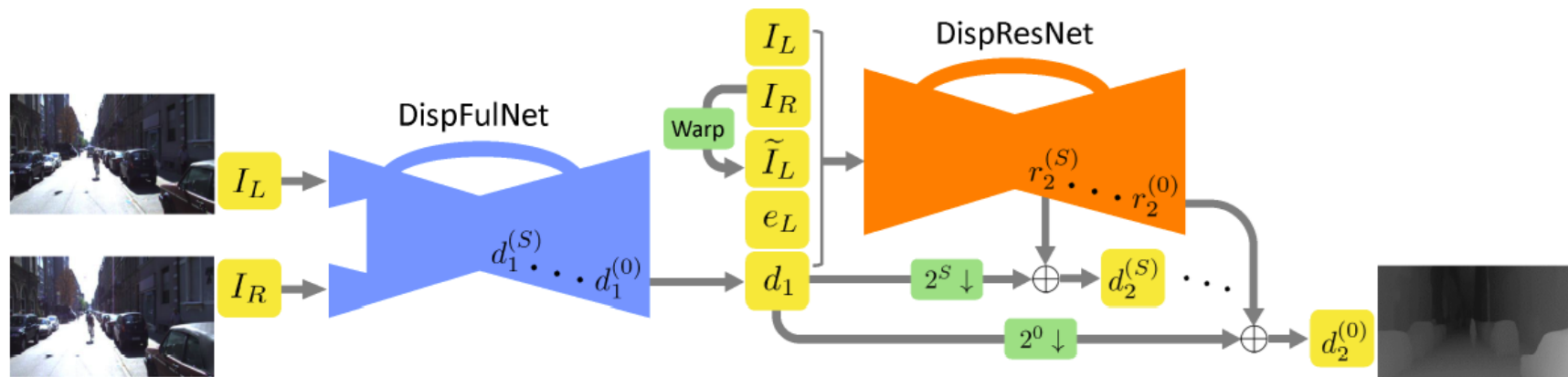
## Part 2. Depth Learning—Our Progress

- Cascade Residual Learning (CRL)
- Zoom and Learn (ZOLE)
- Single View Stereo Matching (SVS)

## Part 3. Conclusion

# Cascade Residual Learning: A Two-stage Convolutional Neural Network for Stereo Matching

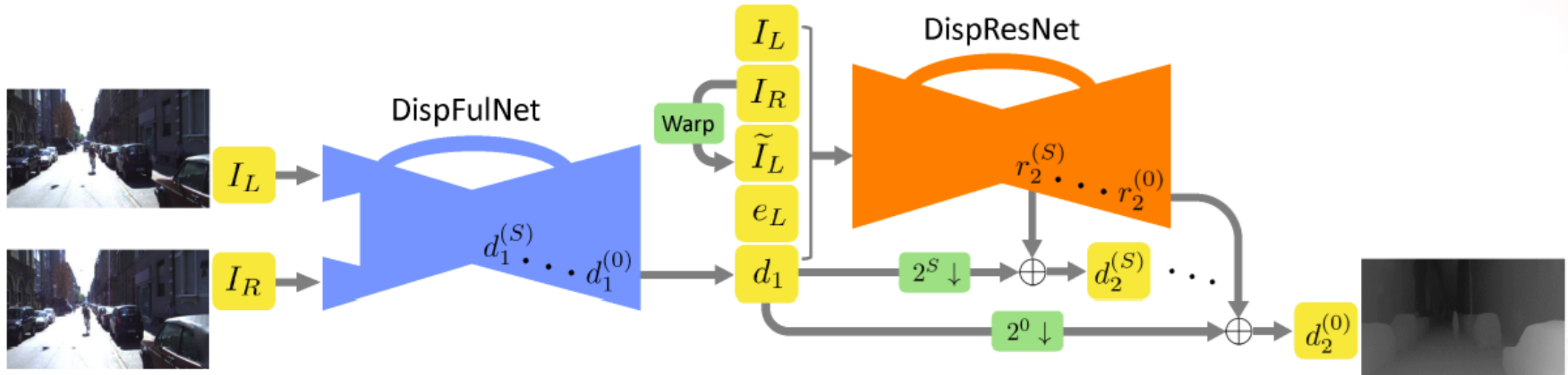
Jiahao Pang, Wenxiu Sun, Jimmy Ren, Chengxi Yang and Yan Qiong,  
SenseTime Research



[11] J. Pang, et al., "Cascade residual learning: A two-stage convolutional neural network for stereo matching," In *Proc. ICCV Workshop*, 2017.

# Cascade Residual Learning (I)

- Cascade residual networks for accurate disparity estimation



Two-stage disparity estimation:

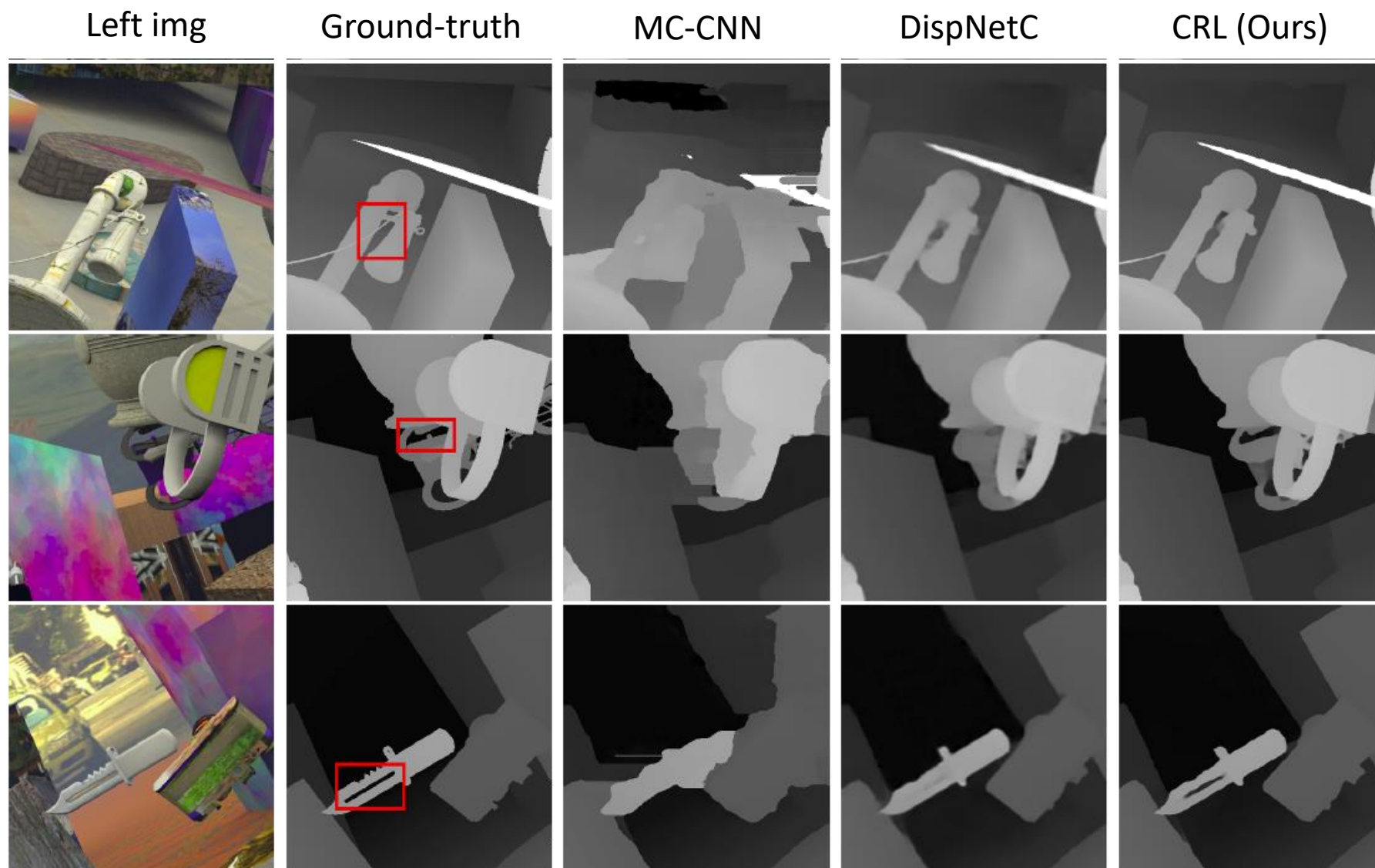
- DispFulNet:** Resembles DispNet, use extra deconvolution modules to obtain full resolution output
- DispResNet:** Inspired by ResNet, the second stage learns the residual signals across multiple scales

# Cascade Residual Learning (II)

- Ranked 1<sup>st</sup> on KITTI Stereo 2015 from Mar. 2017 – Sept. 2017

	Method	Setting	Code	D1-bg	D1-fg	D1-all	Density	Runtime	Environment
1	<b>CRL</b>			2.48 %	3.59 %	2.67 %	100.00 %	0.47 s	Nvidia GTX 1080
J. Pang, W. Sun, J. Ren, C. Yang and Y. Qiong: <a href="#">Cascade Residual Learning: A Two-stage Convolutional Neural Network for Stereo Matching</a> . arXiv preprint arXiv:1708.09204 2017.									
2	<b>GC-NET</b>			2.21 %	6.16 %	2.87 %	100.00 %	0.9 s	Nvidia GTX Titan X
A. Kendall, H. Martirosyan, S. Dasgupta, P. Henry, R. Kennedy, A. Bachrach and A. Bry: <a href="#">End-to-End Learning of Geometry and Context for Deep Stereo Regression</a> . arXiv preprint arxiv:170									
3	<b>DRR</b>			2.58 %	6.04 %	3.16 %	100.00 %	0.4 s	Nvidia GTX Titan X
S. Gidaris and N. Komodakis: <a href="#">Detect, Replace, Refine: Deep Structured Prediction For Pixel Wise Labeling</a> . arXiv preprint arXiv:1612.04770 2016.									
4	<b>L-ResMatch</b>		<a href="#">code</a>	2.72 %	6.95 %	3.42 %	100.00 %	48 s	1 core @ 2.5 Ghz (C/C++)
A. Shaked and L. Wolf: <a href="#">Improved Stereo Matching with Constant Highway Networks and Reflective Loss</a> . arXiv preprint arxiv:1701.00165 2016.									
5	<b>Displets v2</b>		<a href="#">code</a>	3.00 %	5.56 %	3.43 %	100.00 %	265 s	>8 cores @ 3.0 Ghz (Matlab + C/C++)
F. Guey and A. Geiger: <a href="#">Displets: Resolving Stereo Ambiguities using Object Knowledge</a> . Conference on Computer Vision and Pattern Recognition (CVPR) 2015.									
6	<b>D3DNet</b>			2.88 %	6.60 %	3.50 %	100.00 %	0.35 s	Nvidia GTX Titan X

# Cascade Residual Learning (III)





# Outline

## Part 1. Introduction

- Motivation
- Stereo Matching
- Single Image Depth Estimation

## Part 2. Depth Learning—Our Progress

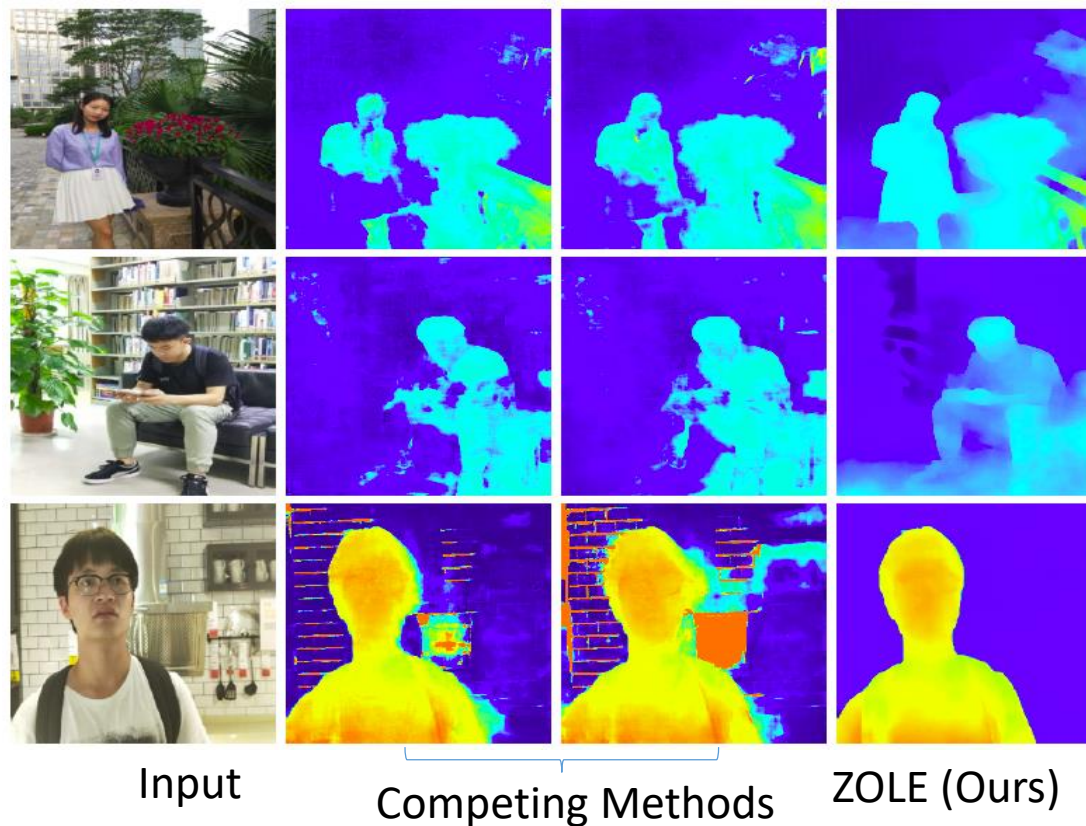
- Cascade Residual Learning (CRL)
- Zoom and Learn (ZOLE)
- Single View Stereo Matching (SVS)

## Part 3. Conclusion

# Zoom and Learn: Generalizing Deep Stereo Matching to Novel Domains

Jiahao Pang<sup>1</sup>, Wenxiu Sun<sup>1</sup>, Chengxi Yang<sup>1</sup>, Jimmy Ren<sup>1</sup>, Ruichao Xiao<sup>1</sup>, Jin Zeng<sup>1</sup>, Liang Lin<sup>1,2</sup>

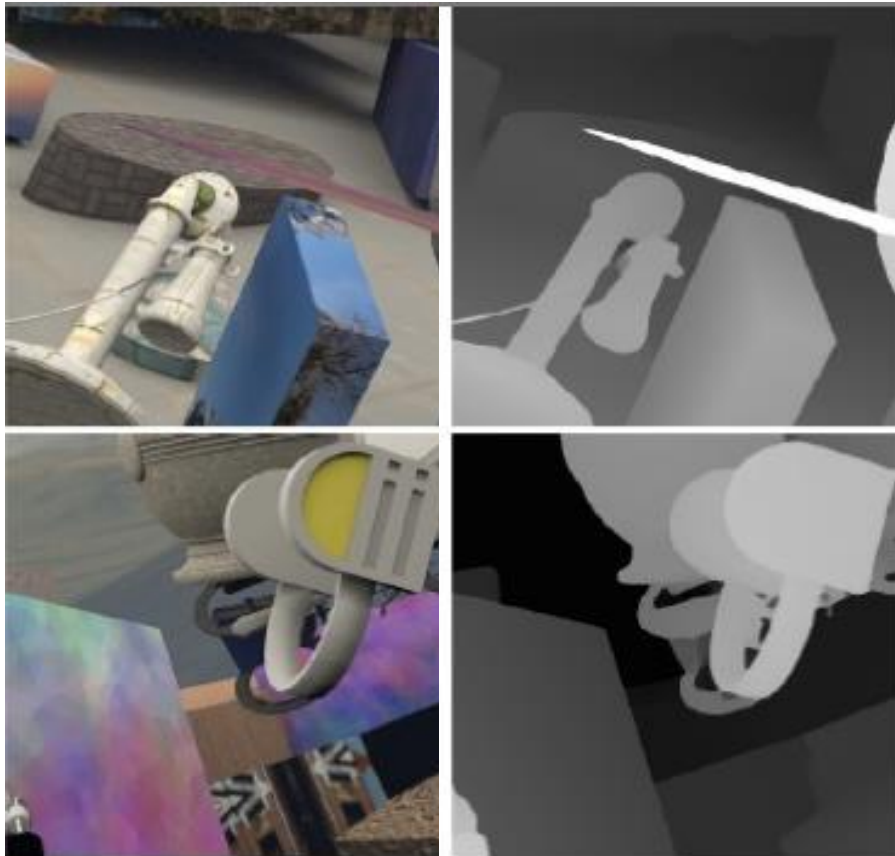
<sup>1</sup> SenseTime Research, <sup>2</sup> Sun Yat-sen University



# Zoom and Learn—Motivation

- **Great** results on pre-trained synthetic dataset

Source domain



However

- Directly apply model to a **new domain** brings **horrible** results

Target domain



- **Finetuning** in the target domain?  
No, hard to collect ground-truth depth in new domain, *e.g.* bokeh
- Our solution: a **self-adaptation** approach

Synthetic data, with  
ground-truth



Stereo pairs in new  
domain, no ground-truth

- **Key observation:** up-sampled stereo pair leads to disparity map with **extra details**
- Stereo matching CNN  $S$  parameterized by  $\Theta$ , stereo pair  $P$ , consider two schemes

Scheme A

$$D = S(P; \Theta)$$

Scheme B

$$D' = \frac{1}{r} \cdot \downarrow_r (S(\uparrow_r(P); \Theta))$$

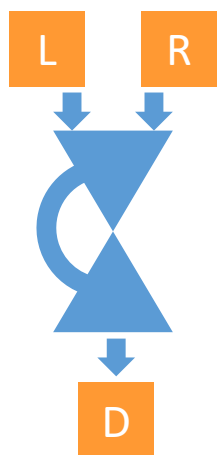


$$D' = \frac{1}{r} \cdot \downarrow_r (S(\uparrow_r(P); \Theta))$$

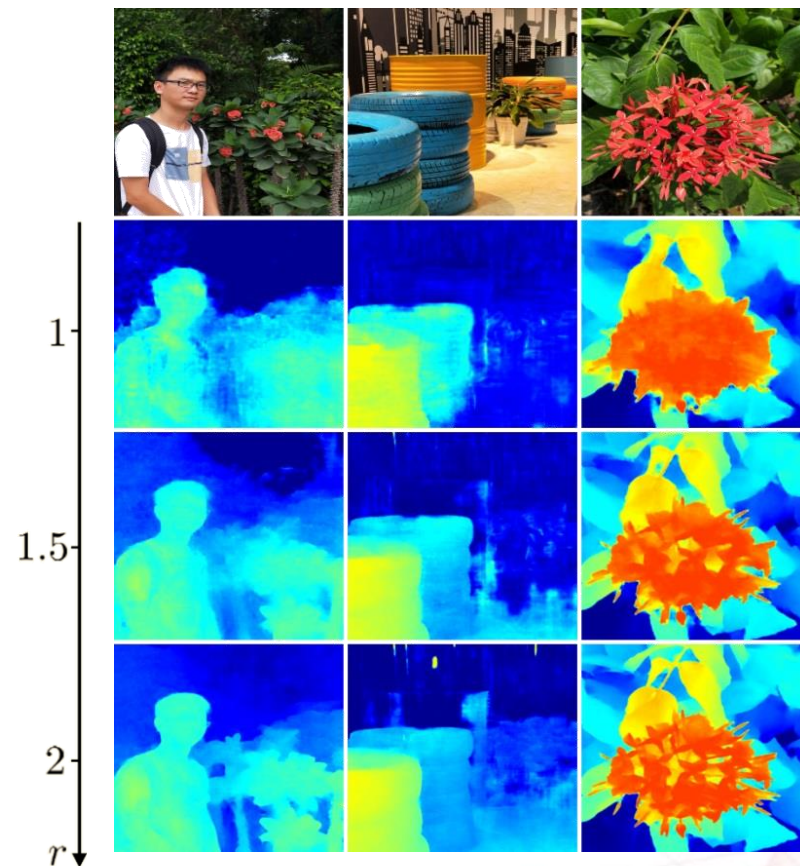
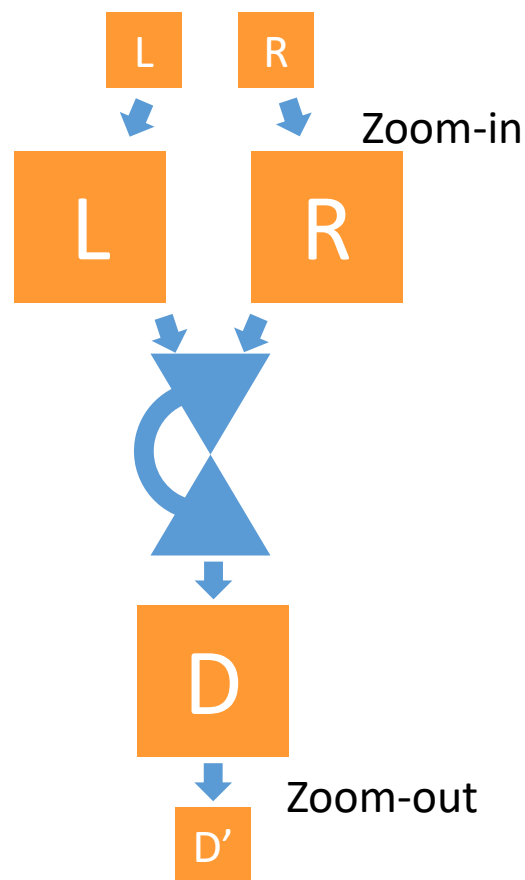
Scheme B brings more details

$$D = S(P; \Theta)$$

Scheme A



Scheme B



- However
  - A bigger  $r$  does not necessarily mean better results
  - Performance first improves then deteriorates

Results of DispNetC on KITTI 2015

Network	Resolution				
	896	1280	1664	2048	2432
DispNetC	14.26%	9.97%	<b>8.81%</b>	9.17%	10.53%
DispNetS	18.95%	11.61%	9.18%	<b>8.64%</b>	9.08%

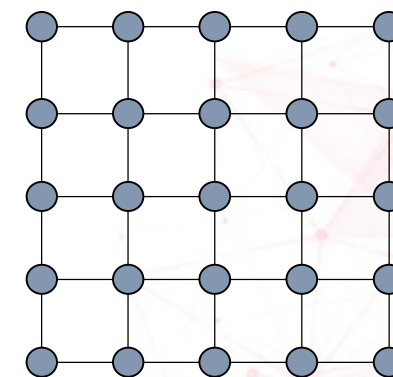
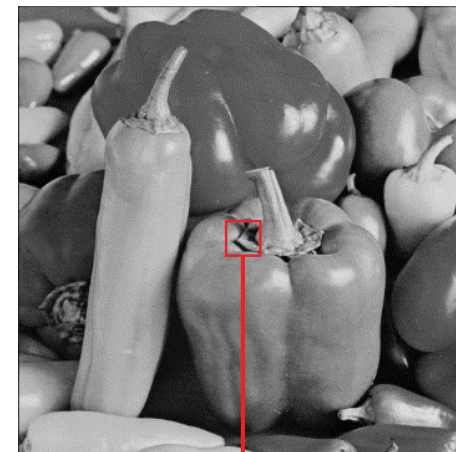
- Analysis
  - Up-sampling means perform stereo matching at subpixel accuracy
  - Larger input translates to smaller receptive field

- Strategy
- Let the CNN **learn its own** higher-res details
  - **Graph Laplacian regularization** to keep desired edges while smoothing out artifacts

- Minimize **graph Laplacian regularizer**: widely used for image restoration

$$\mathbf{s}^T \mathbf{L} \mathbf{s} = \sum_{(i,j) \in \mathcal{E}} w_{ij} (\mathbf{s}(i) - \mathbf{s}(j))^2 \quad \leftarrow \mathbf{L} : \text{Graph Laplacian matrix}$$

- Bigger weights  $w_{ij}$ , higher similarity of pixels  $i$  and  $j$ , force  $\mathbf{s}(i)$  close to  $\mathbf{s}(j)$
- Graph Laplacian regularization—simple & **perform well** on many restoration problems
  - Denoising
  - Super-resolution
  - Deblurring
  - Dequantization of JPEG images
  - Bit-depth enhancement
- We integrate it as **a loss function** in our work



Graph of a  $5 \times 5$  patch,  
(may not be a grid graph)

- Settings

$S(\cdot; \Theta^{(0)})$  : init stereo network pre-trained with synthetic data,  $\Theta^{(0)}$  : model parameter

$N$  stereo pairs  $P_i = (L_i, R_i), 1 \leq i \leq N$  for training

- The first  $N_{\text{dom}}$  pairs are **real** stereo pairs of the target domain, no ground-truth
- The rest  $N_{\text{syn}} = N - N_{\text{dom}}$  pairs are **synthetic** data, they have ground truth disparities  $D_i$

- We solve for a new set of model parameter  $\Theta^{(k+1)}$  at iter.  $k$

- First create a set of **pseudo "ground-truths"** for the  $N_{\text{dom}}$  real stereo pairs by zooming (up-sampling)

$$D_i = \frac{1}{r} \cdot \downarrow_r (S(\uparrow_r(P_i); \Theta^{(k)})), 1 \leq i \leq N_{\text{dom}}$$

- $D_i \in \mathbb{R}^{Mm}$ , we divide a disparity map  $D_i$  into  $M$  square patches
- Matrix extracting the  $j$ -th patch from  $D_i$  is denoted as  $\mathbf{R}_j$ ,  $\mathbf{R}_j \in \mathbb{R}^{m \times Mm}$



- Formulation

Data term: drives  $\mathbf{s}_{ij}$  to be similar to  $\mathbf{d}_{ij}$

Smoothness term: graph Laplacian regularization

$$\Theta^{(k+1)} = \arg \min_{\Theta} \sum_{i=1}^{N_{\text{dom}}} \sum_{j=1}^M \|\mathbf{s}_{ij} - \mathbf{d}_{ij}\|_1 + \lambda \cdot \mathbf{s}_{ij}^T \mathbf{L}_{ij}^{(k)} \mathbf{s}_{ij} + \tau \cdot \sum_{i=N_{\text{dom}}+1}^N \|S(P_i; \Theta) - D_i\|_1,$$

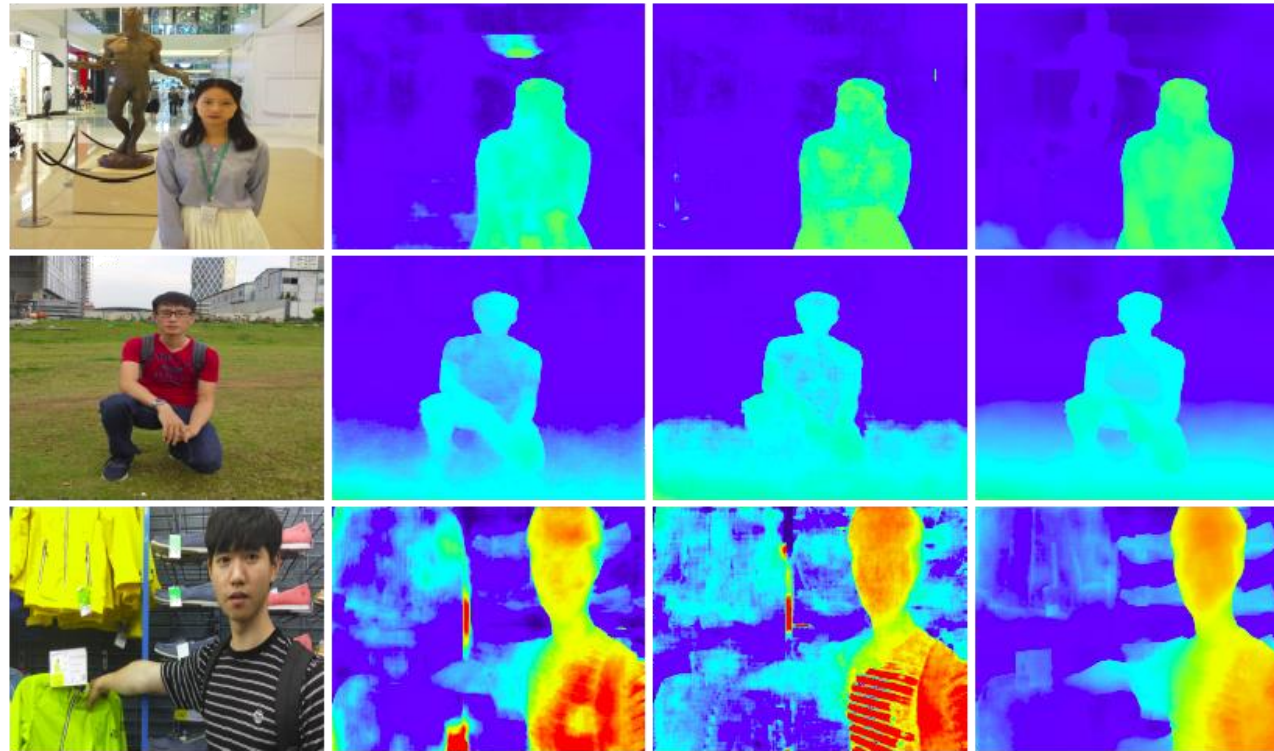
$$\text{s.t. } \mathbf{s}_{ij} = \mathbf{R}_j \cdot \text{vec}(S(P_i; \Theta)), \quad \mathbf{d}_{ij} = \mathbf{R}_j \cdot \text{vec}(D_i)$$

Feasibility: a stereo model works well for the target domain should perform reasonably on the synthetic data

- graph Laplacian  $\mathbf{L}_{ij}^{(k)}$  are pre-computed, based on the left image and the predictions
- In practice, we resolve the optimization with stochastic gradient descent

# Zoom and Learn—Results (I)

- Refer to our paper for the detailed setting

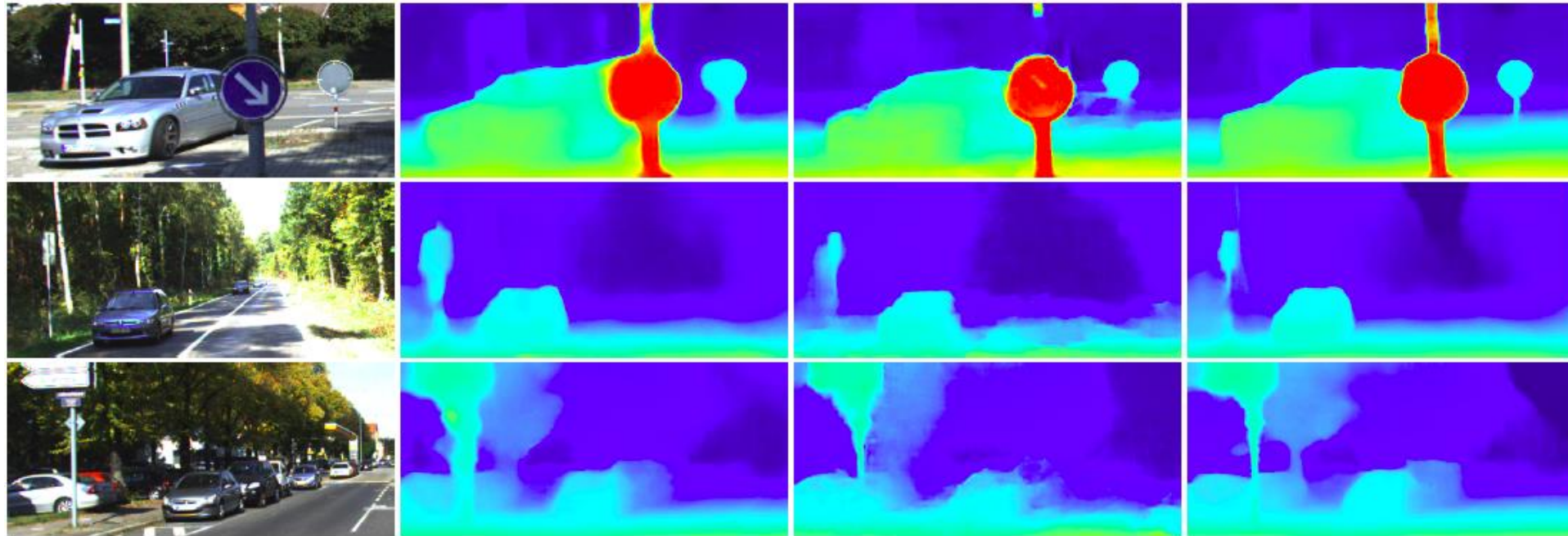


Left img    Tonioni *et al.* [14]    DispNetC    ZOLE (Ours)

Dataset	Metric		Model							
			Tonioni		DispNetC		DispNetC-80		ZOLE	
Smartphone	PSNR	SSIM	22.92	0.845	21.99	0.790	22.39	0.817	<b>23.12</b>	<b>0.855</b>
FlyingThings3D-80	EPE	3ER	1.08	6.79%	1.03	5.63%	<b>0.93</b>	<b>5.11%</b>	1.11	6.54%

# Zoom and Learn—Results (II)

- Refer to our paper for detailed setting



Input

Tonioni *et al.*

DispNetC

ZOLE (Ours)

Metric	Model		
	Tonioni	DispNetC	ZOLE
EPE	1.27	1.64	<b>1.25</b>
3ER	7.06%	11.41%	<b>6.76%</b>

# Outline

## Part 1. Introduction

- Motivation
- Stereo Matching
- Single Image Depth Estimation

## Part 2. Depth Learning—Our Progress

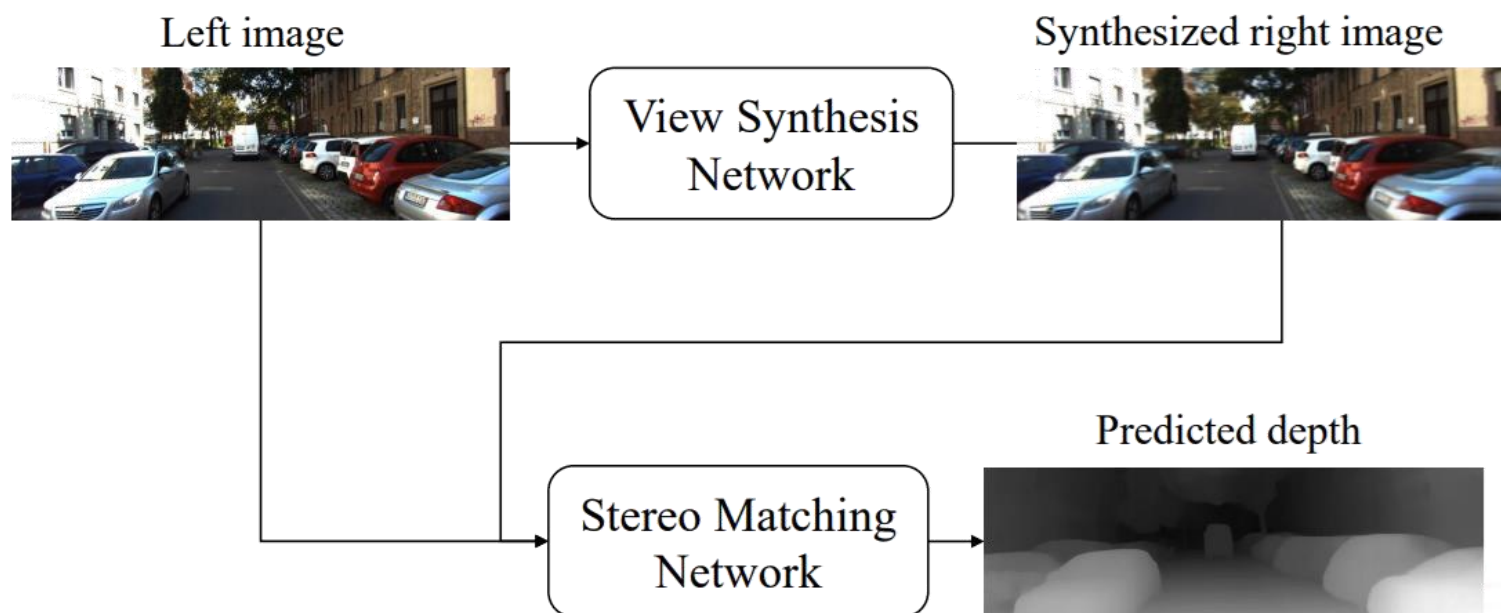
- Cascade Residual Learning (CRL)
- Zoom and Learn (ZOLE)
- Single View Stereo Matching (SVS)

## Part 3. Conclusion

# Single View Stereo Matching

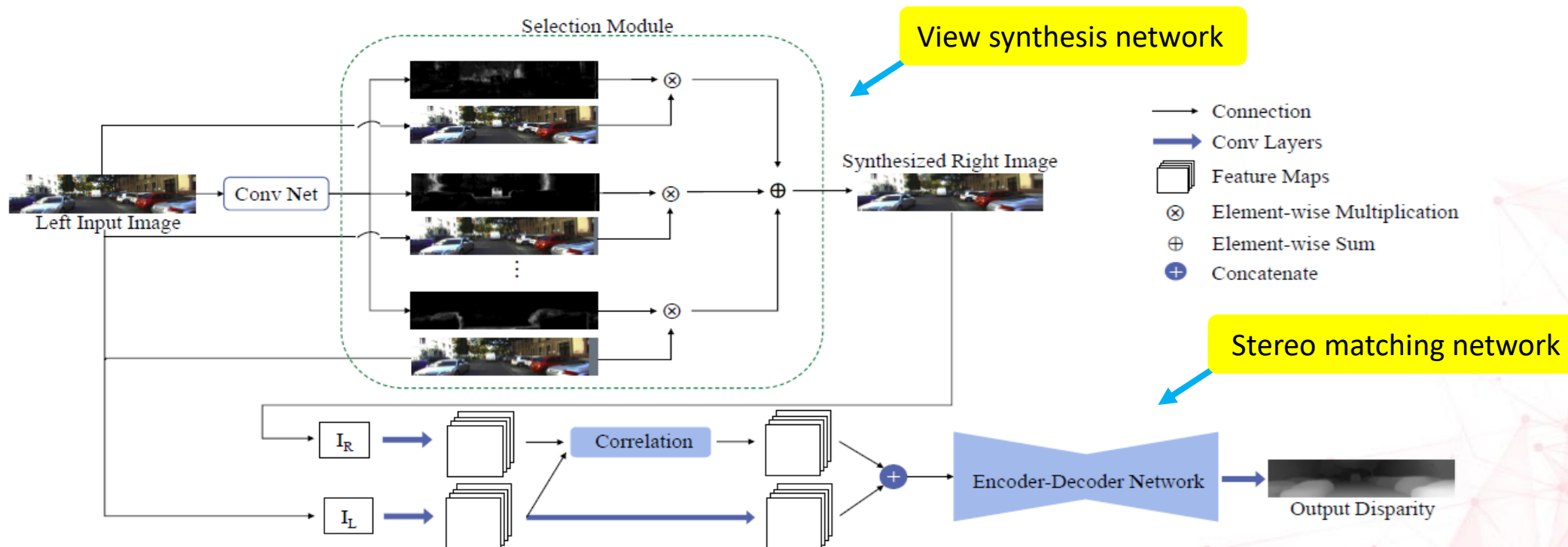
Yue Luo<sup>1</sup>, Jimmy Ren<sup>1</sup>, Mude Lin<sup>1,2</sup>, Jiahao Pang<sup>1</sup>, Wenxiu Sun<sup>1</sup>, Hongsheng Li<sup>3</sup>, Liang Lin<sup>1,2</sup>

<sup>1</sup>SenseTime Research, <sup>2</sup>Sun Yat-sen University, <sup>3</sup>CUHK



# Single View Stereo Matching—Method (I)

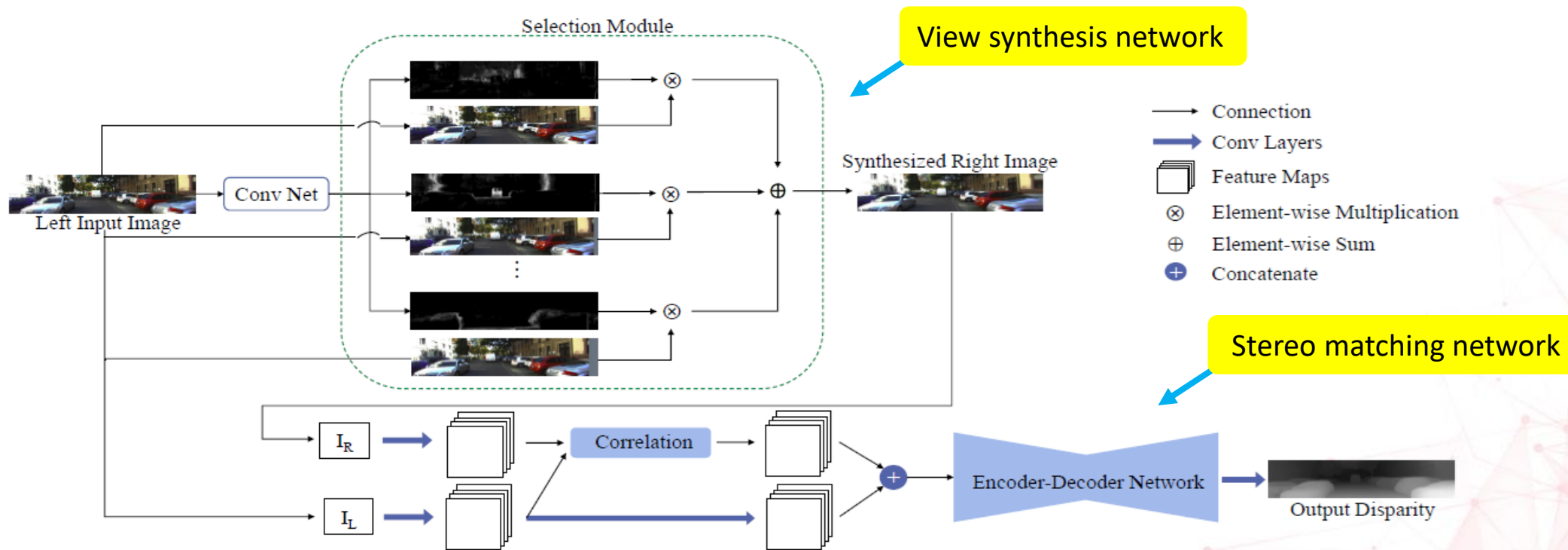
- Decomposed into two stages: **view synthesis** & **stereo matching**
- **View synthesis network**: based on Deep3D [16]
  - Formulate the left-to-right transformation with differentiable selection module
- **Stereo matching network**: based on DispFulNet of our CRL



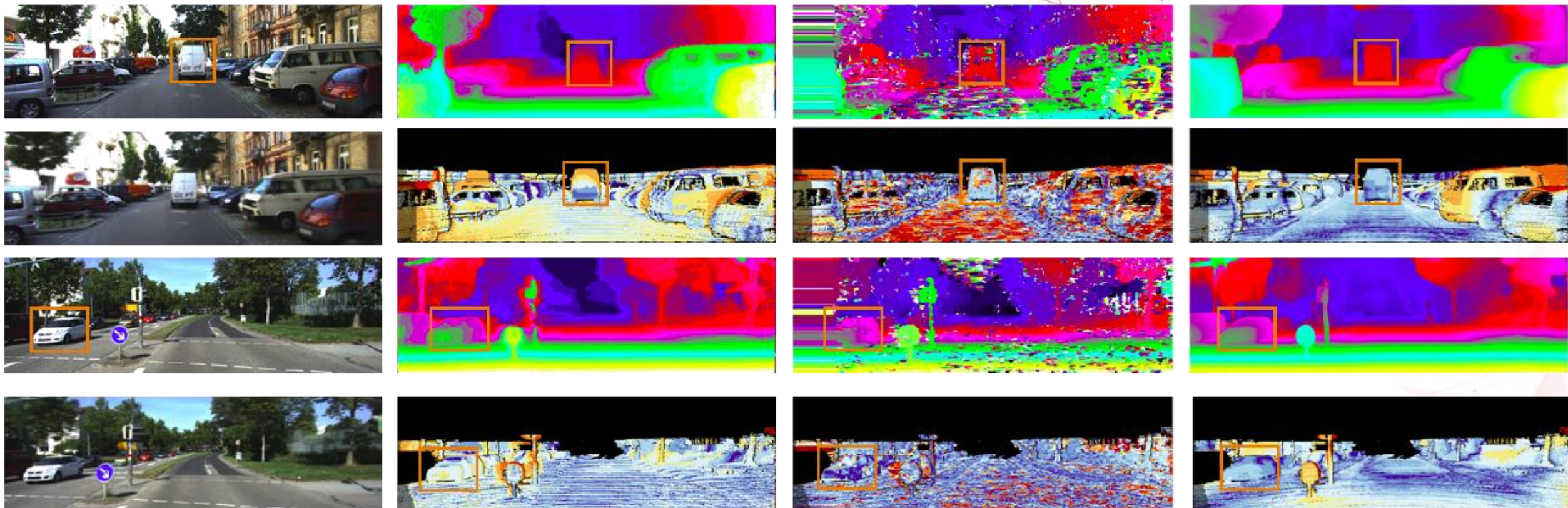
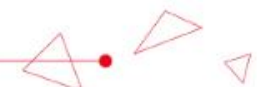
# Single View Stereo Matching—Method (II)



- Benefits
  - Explicitly encode the **geometric transformation** within individual networks to better solve the problem
  - Demand on labelled depth data is greatly **alleviated**



# Single View Stereo Matching—Results (I)



Input

Gordard *et al.* [7]

OCV-BM (Stereo)

SVS (Ours)

Method	D1-bg	D1-fg	D1-all
Godard <i>et al.</i>	27.00	28.24	27.21
OCV-BM	<b>24.29</b>	30.13	25.27
Ours	25.18	<b>20.77</b>	<b>24.44</b>

- Ours produce **sharp** edges
- **Outperforms stereo block matching method!**



# Single View Stereo Matching—Results (I)



Qualitative results on KITTI Eigen test set



Qualitative results on Make3D dataset (top two rows) and Cityscapes dataset (bottom two rows) using the model trained on KITTI

# Single View Stereo Matching—Results (II)



# Conclusion

## Part 1. Introduction

- Motivation
- Stereo Matching
- Single Image Depth Estimation

## Part 2. Depth Learning—Our Progress

- Cascade Residual Learning (CRL)
- Zoom and Learn (ZOLE)
- Single View Stereo Matching (SVS)

## Part 3. Conclusion

## Part 3. Conclusion

# Conclusion

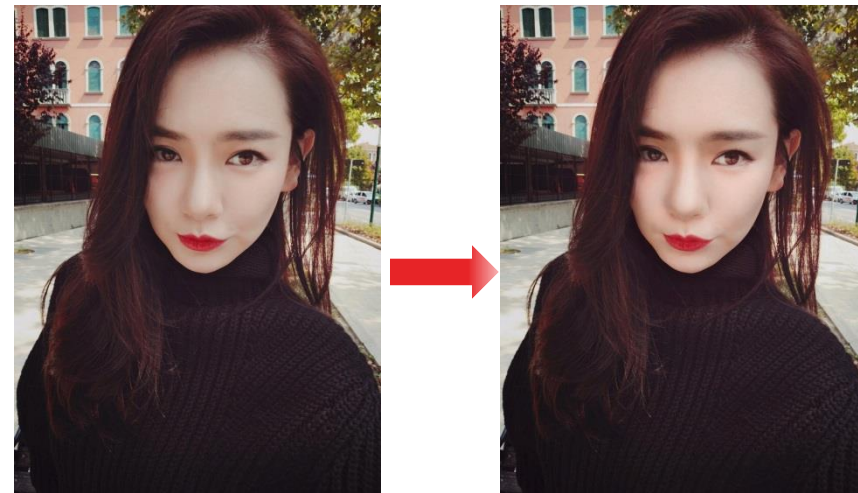
- Depth estimation is important in many areas
- Deep learning brings a new era for depth estimation
- Pure data-driven approach can be fragile (e.g., single image depth estimation)
- Combine data-driven approaches and model-based approaches is a new trend
- Welcome to SenseTime to MAKE IT HAPPEN!

## About Our Team: Image Restoration and Enhancement

Depth estimation, computational photography, image processing on smartphones



Our bokeh solution on recently  
launched VIVO V9



3D relighting on smartphones

**Get involve!** Internship and full-time employee are welcome

Location: Hong Kong, Shenzhen, Beijing, Shanghai, etc

Thank You!

Q & A