

ActiveStereoNet

End-to-End Self-Supervised Learning for Active Stereo Systems

Presenter: Yinda Zhang

Yinda Zhang^{1,2}, Sameh Khamis¹, Christoph Rhemann¹, Julien Valentin¹, Adarsh Kowdle¹, Vladimir Tankovich¹, Michael Schoenberg¹, Shahram Izadi¹, Thomas Funkhouser^{1,2}, Sean Fanello¹

Google Inc.¹ Princeton University²

Project Webpage: <http://asn.cs.princeton.edu/>

Goal: Sensing 3D Geometry



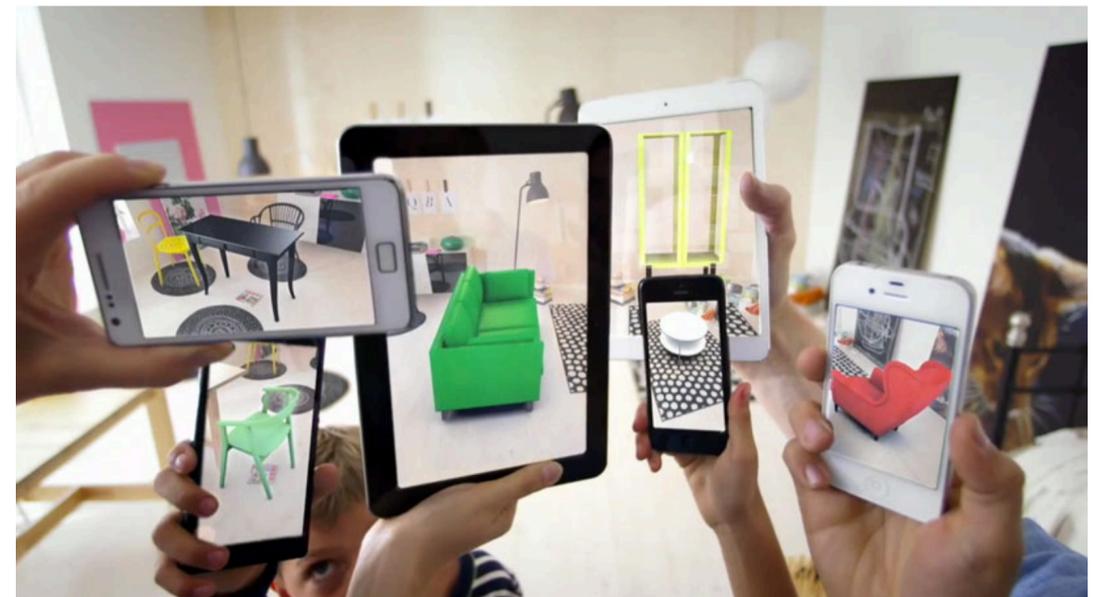
Biometric authentication



Autonomous driving



Indoor Robotics

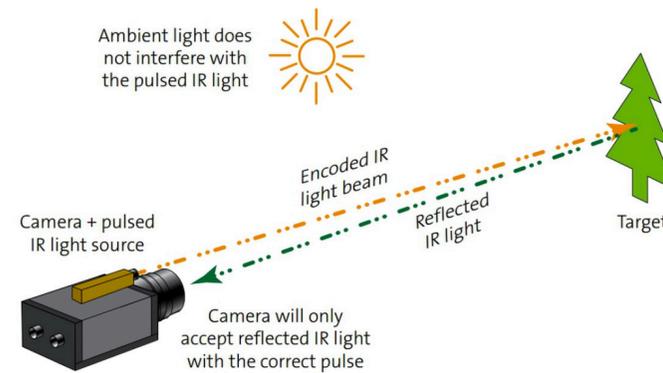


Augmented Reality

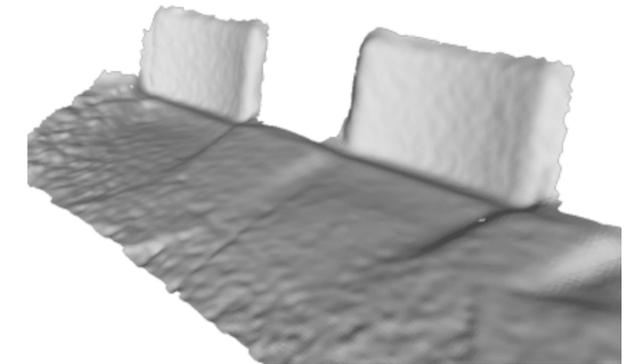
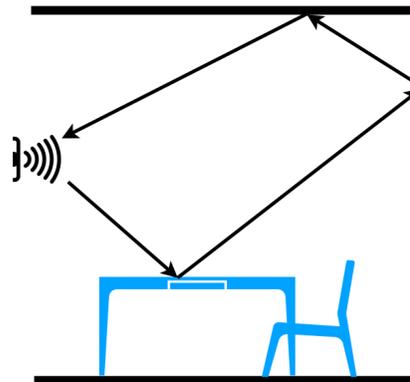
Active Depth Sensing

- Time of Flight

- ~~X~~ Fast motion
- ~~X~~ Multi-Path



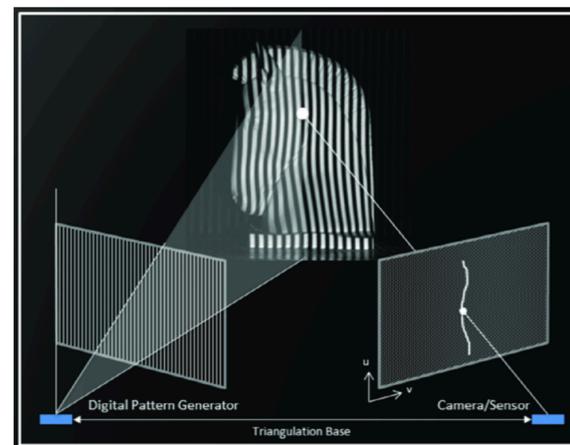
Time of Flight



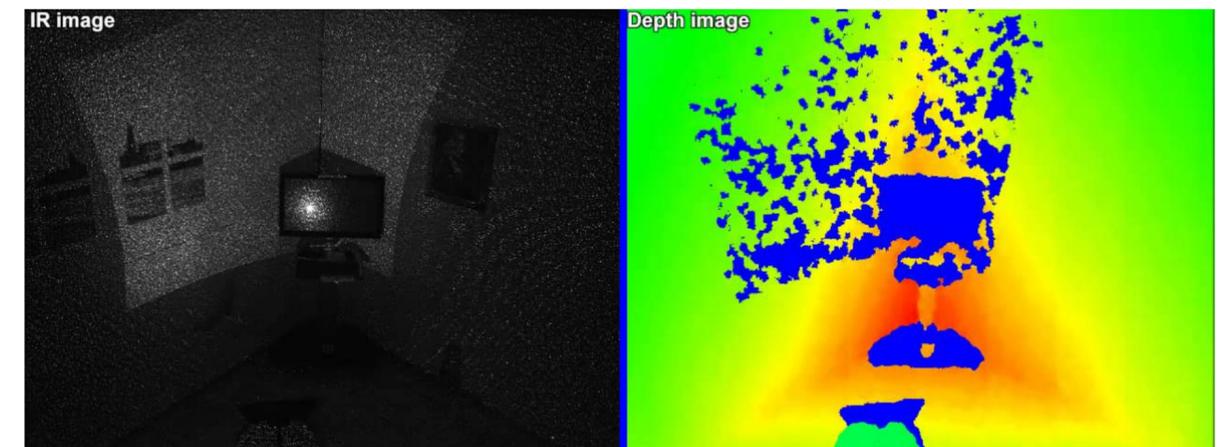
Multi-Path Interference

- Structured Light

- ~~X~~ Calibration
- ~~X~~ Multi-Device



Structured Light



Multi-Device Interference

Passive Depth Sensing

- Stereo Matching

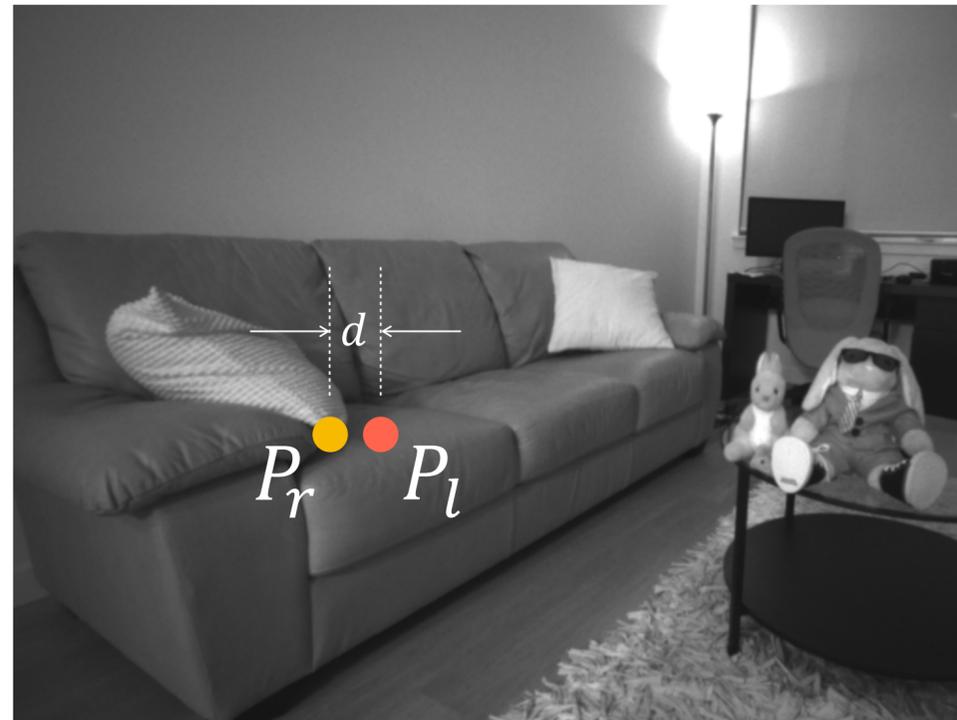
$$Z = \frac{bf}{d}$$

b : distance between camera centers
 f : camera focal length
 d : disparity
 Z : depth

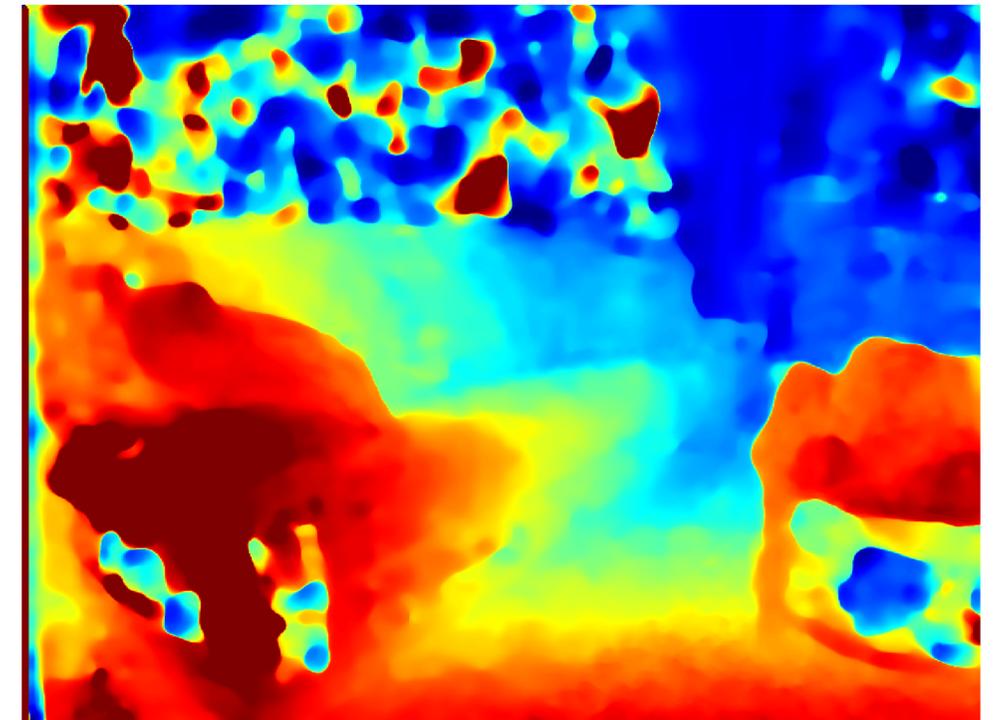
X Texture-less Region



Left View



Right View



Disparity

Active Stereo System

- Stereo Matching

$$Z = \frac{bf}{d}$$

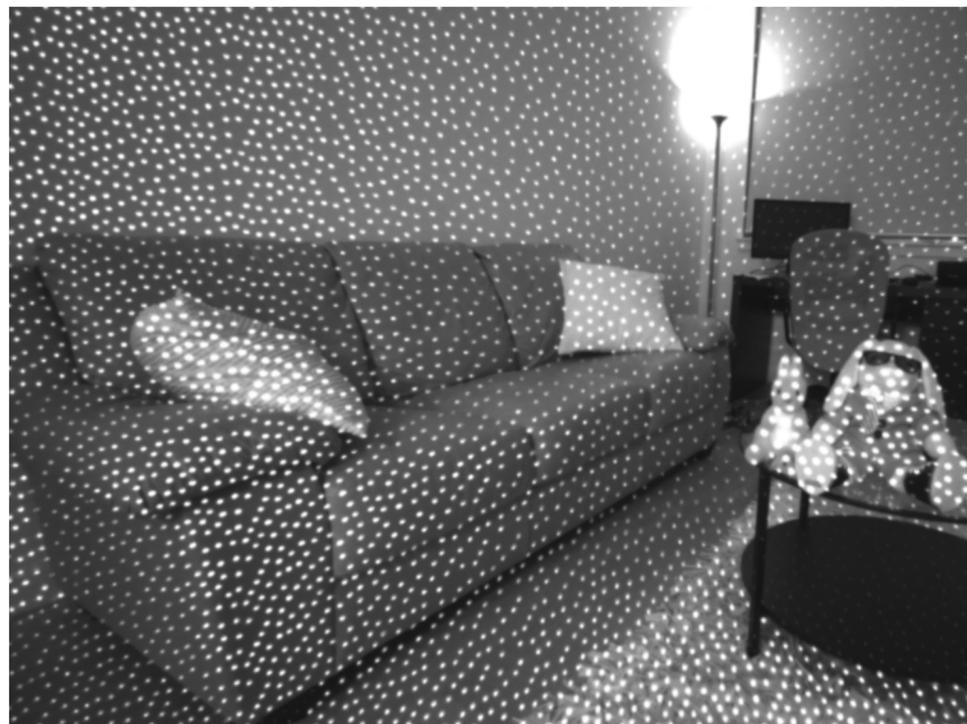
b: distance between camera centers

f: camera focal length

d: disparity

Z: depth

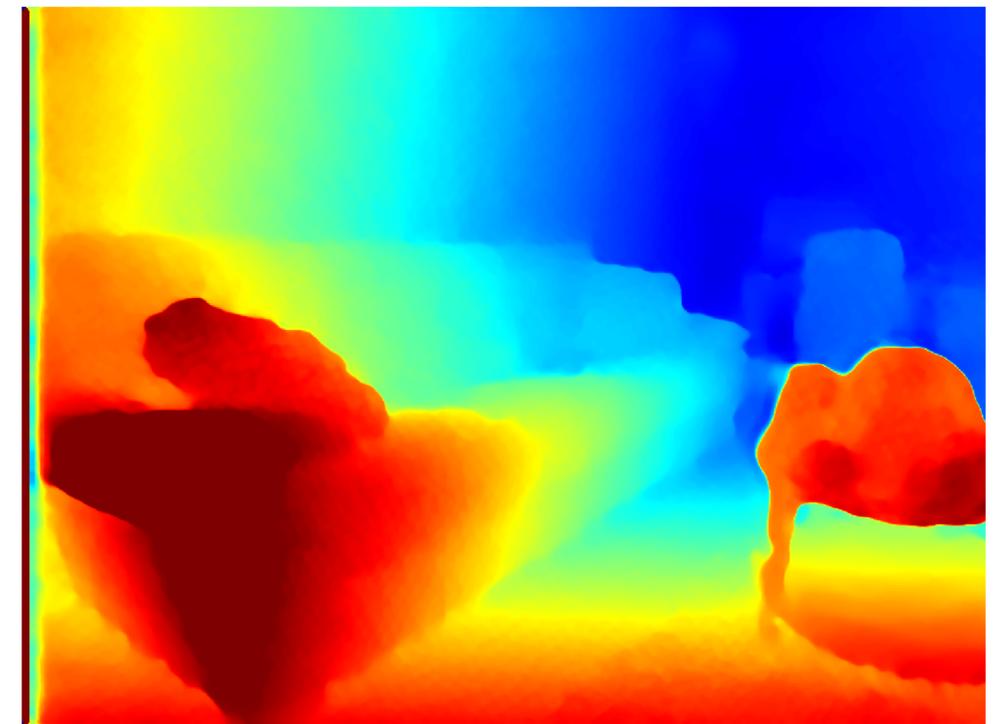
~~X Texture less Region~~



Left View

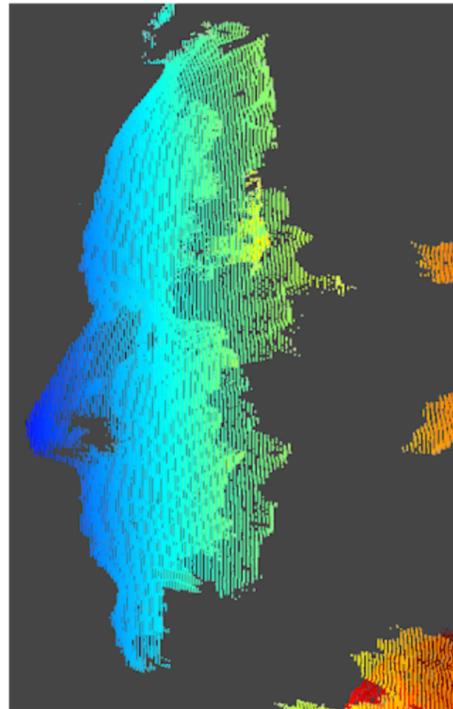
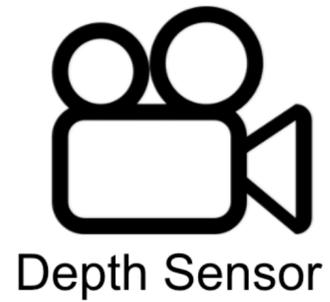


Right View

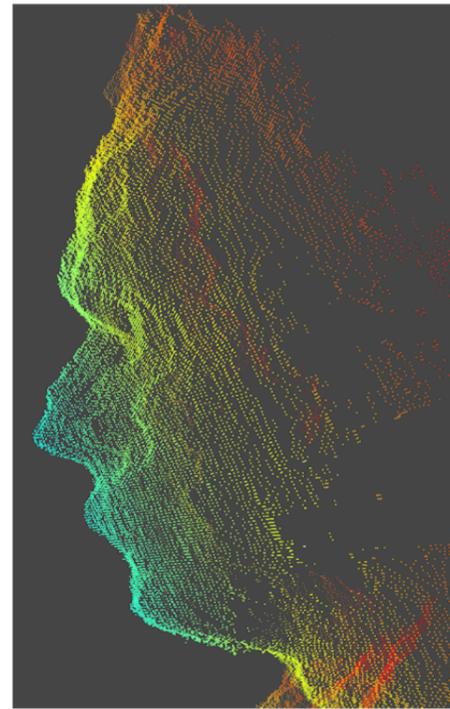


Disparity

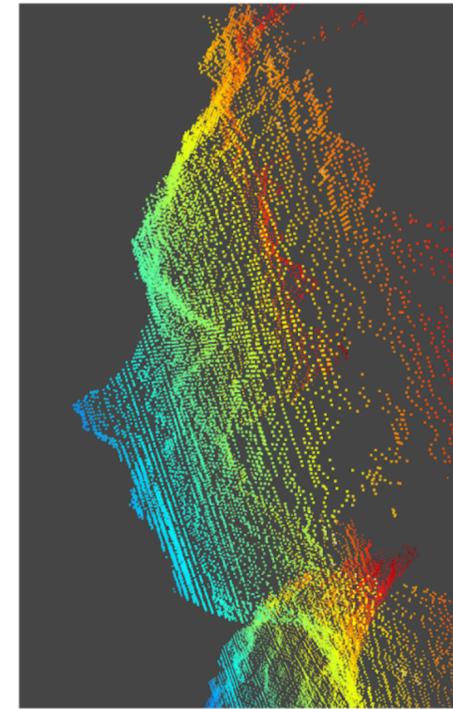
Active Stereo System



300mm



500mm



700mm



2500mm

 **Use deep learning!**

 **No ground truth...**



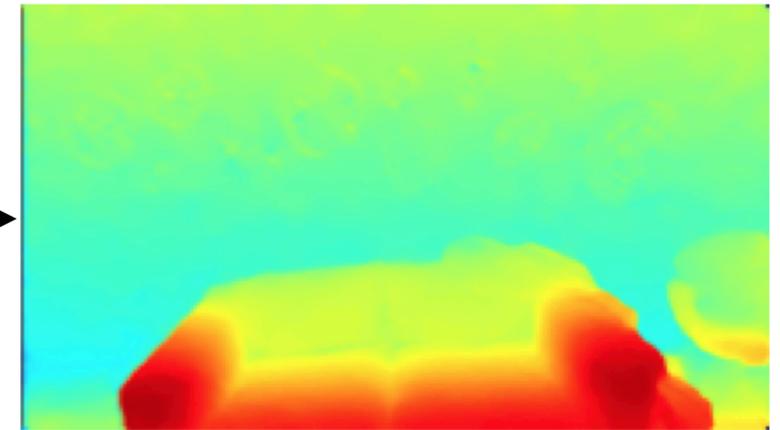
ActiveStereoNet



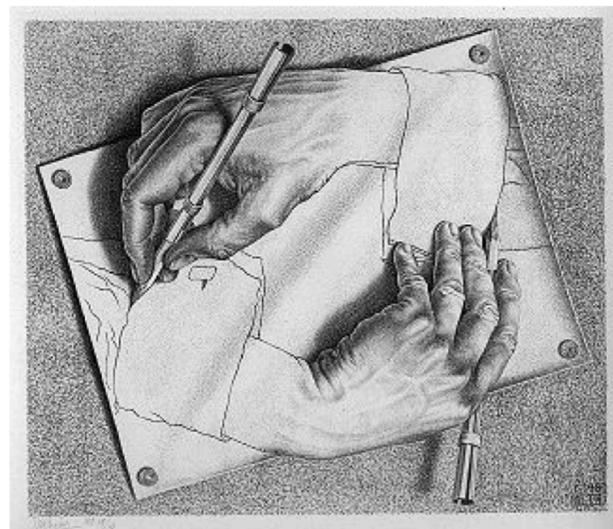
Input: Left/Right View



End-to-End System



Output: Disparity



Self-supervised Learning

=



Annotation



Supervision



Just keep running...

Self-Supervised Learning

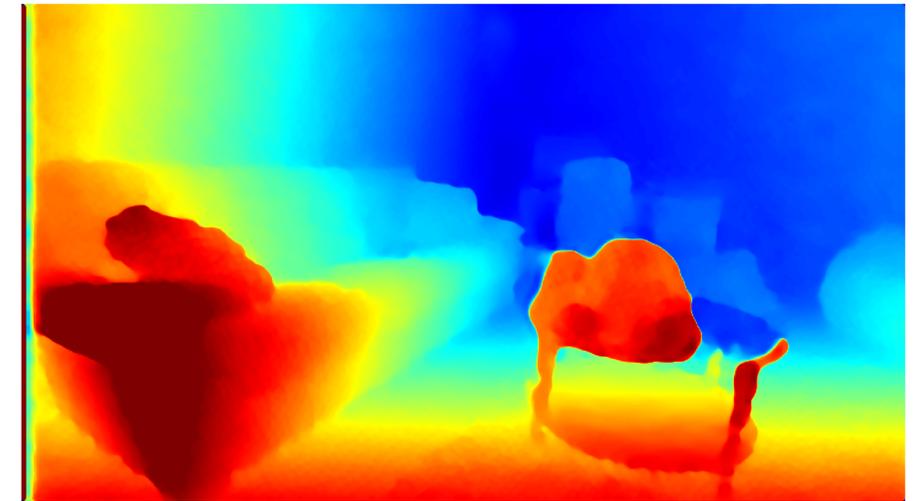


Left View



Right View

Neural Network



Estimated Disparity

Self-Supervised Learning

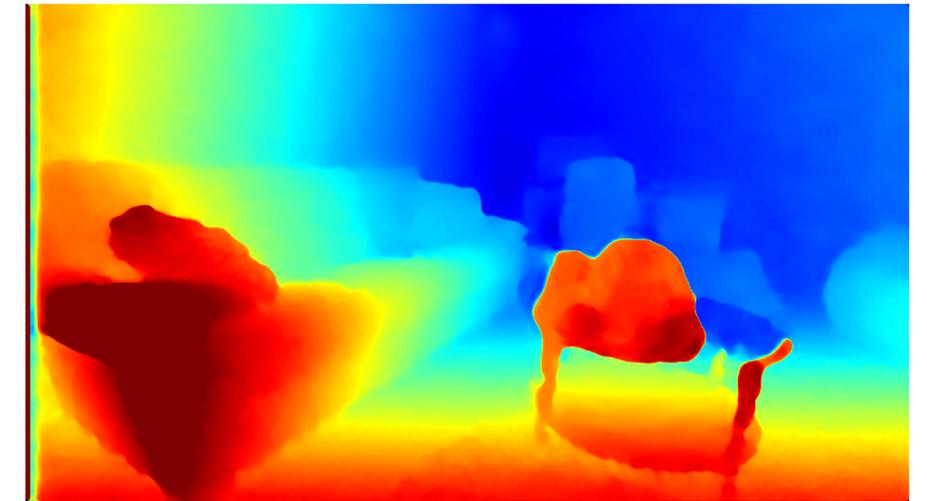


Left View



Right View

Neural Network



Estimated Disparity



Left View



Right View

Self-Supervised Learning

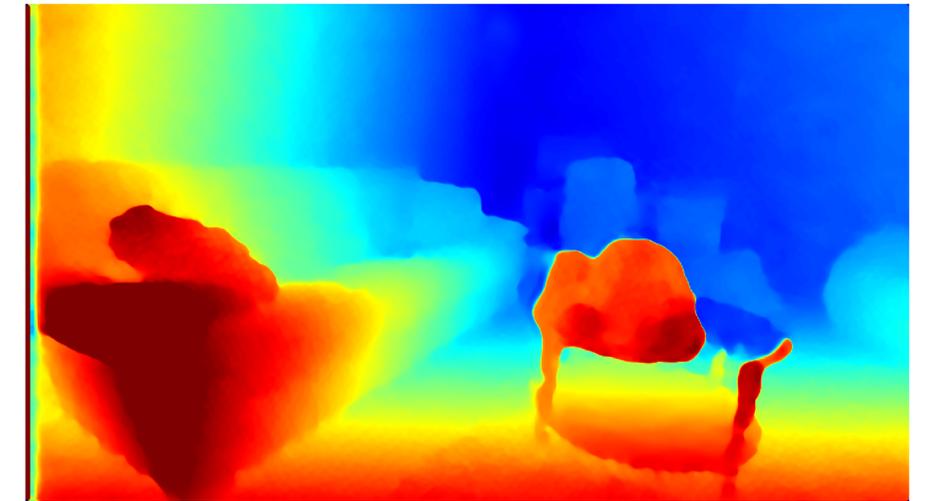


Left View

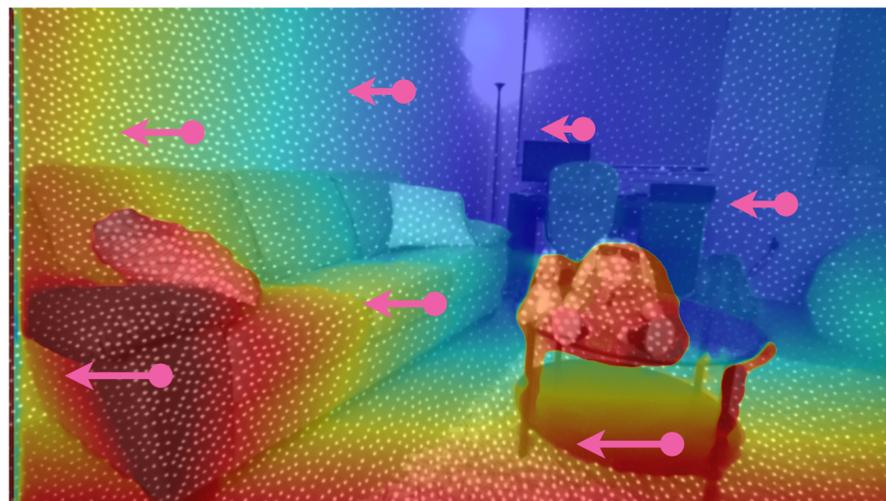


Right View

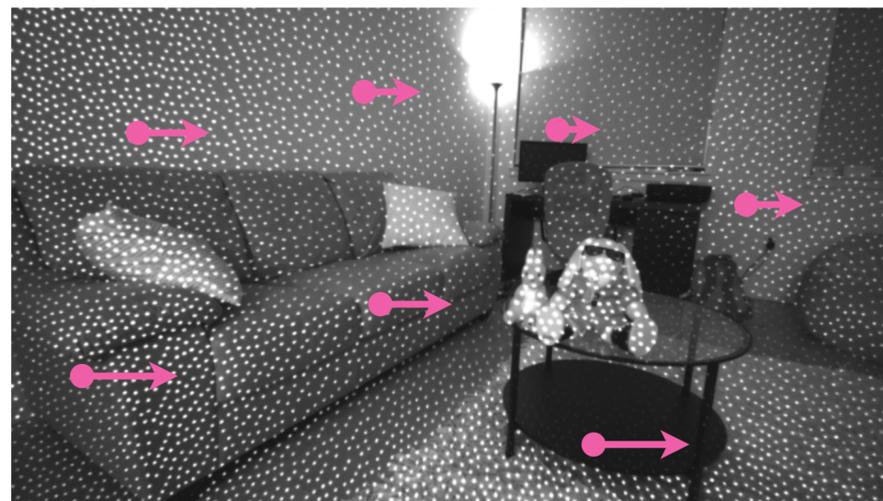
Neural Network



Estimated Disparity



Left View



Right View

Self-Supervised Learning

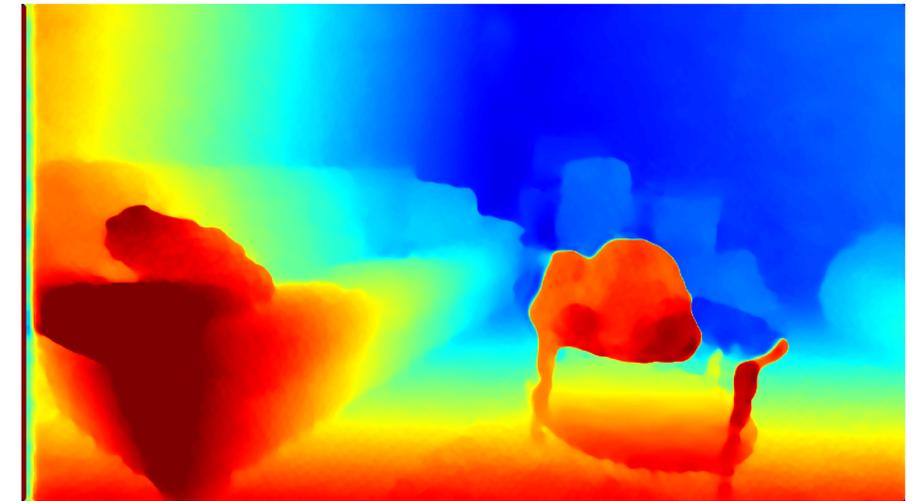


Left View

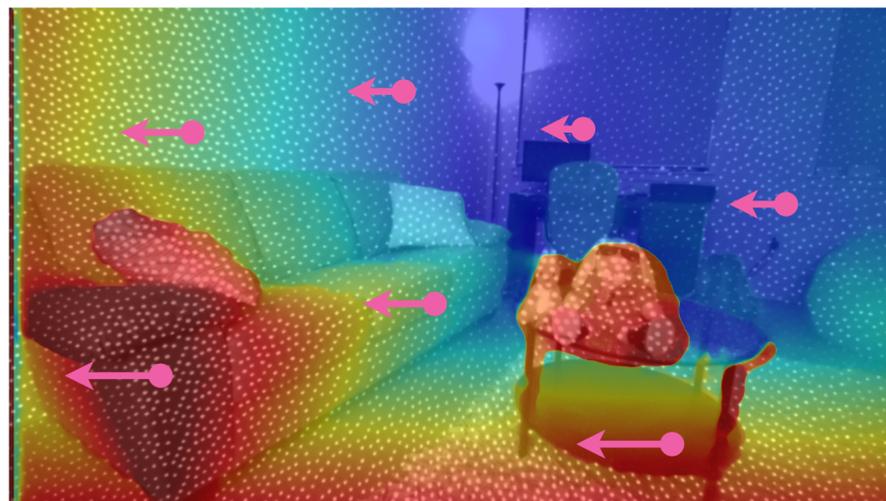


Right View

Neural Network



Estimated Disparity



Left View



Reconstructed Left View

Self-Supervised Learning

$$\text{Photometric Loss} = | \text{Left View} - \text{Reconstructed Left View} |$$

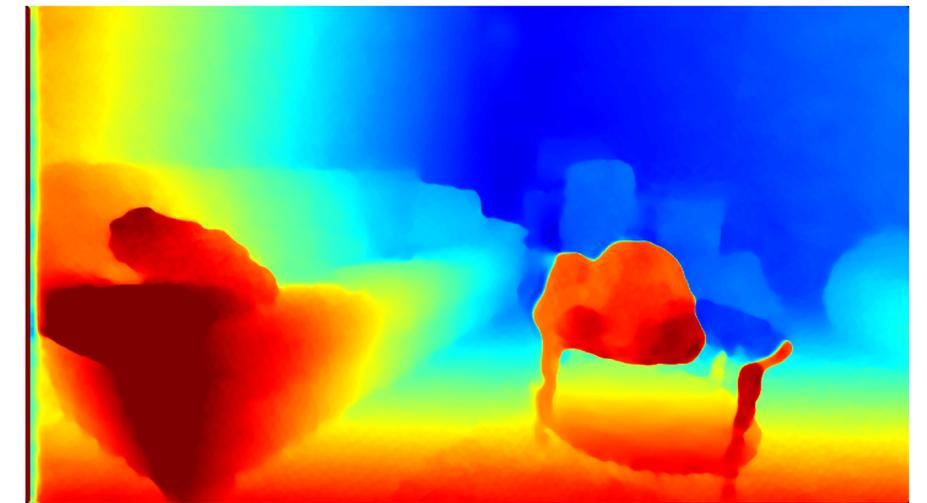


Left View

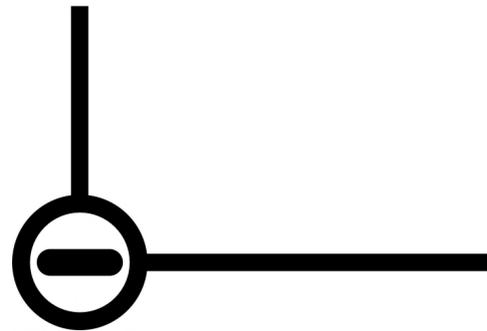


Right View

Neural Network

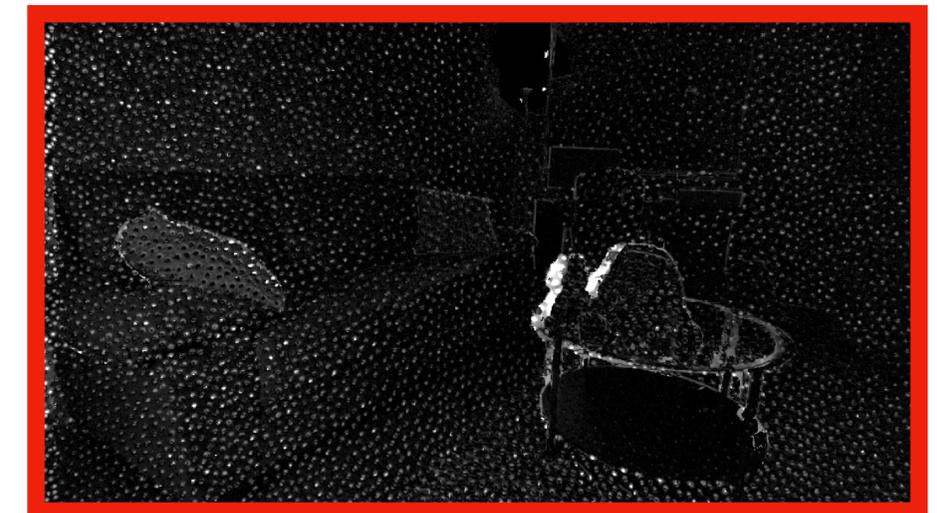


Estimated Disparity



Reconstructed Left View

=



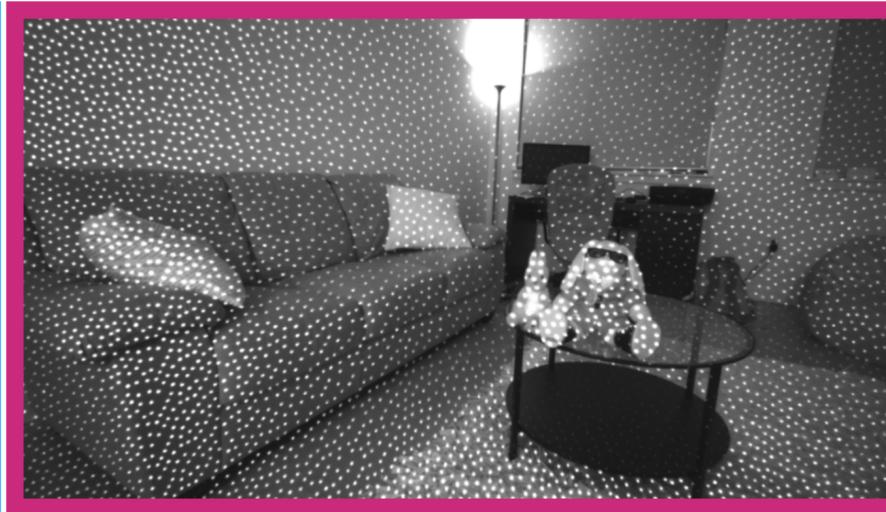
Photometric Loss

Self-Supervised Learning

$$\text{Photometric Loss} = | \text{Left View} - \text{Warping}(\text{Right View}, \text{Left Disparity}) |$$

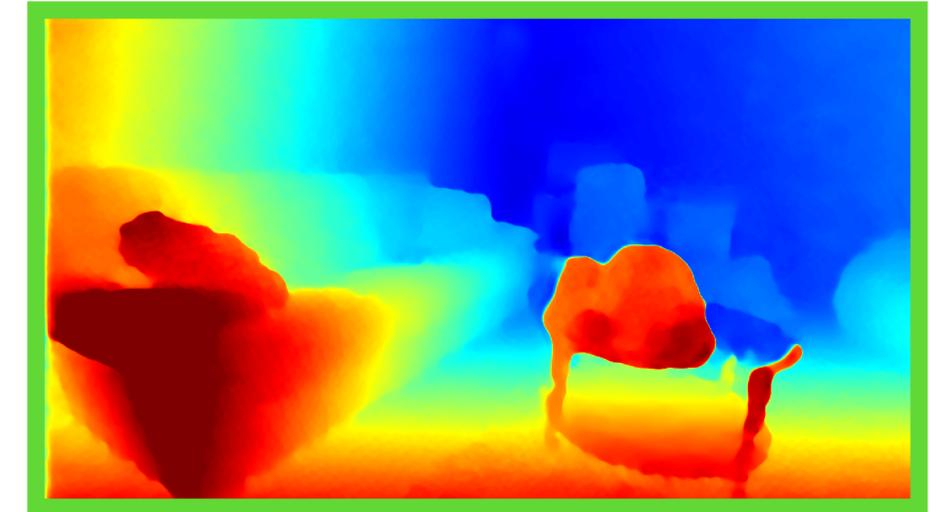


Left View

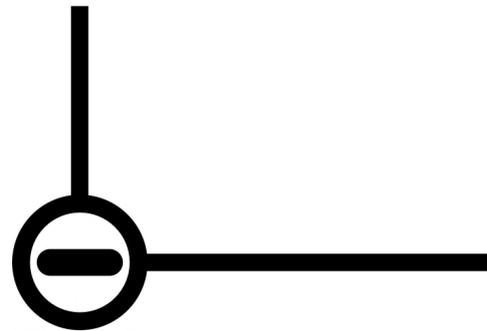


Right View

Neural Network

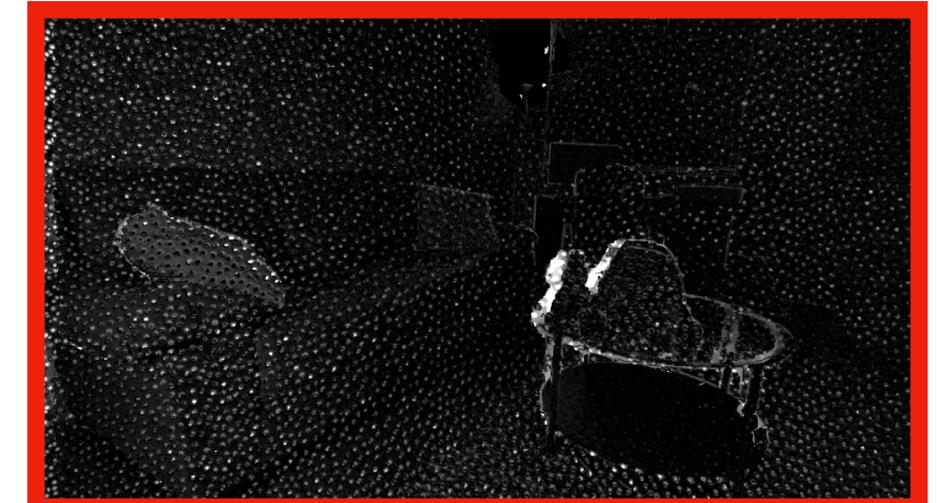


Estimated Disparity



Reconstructed Left View

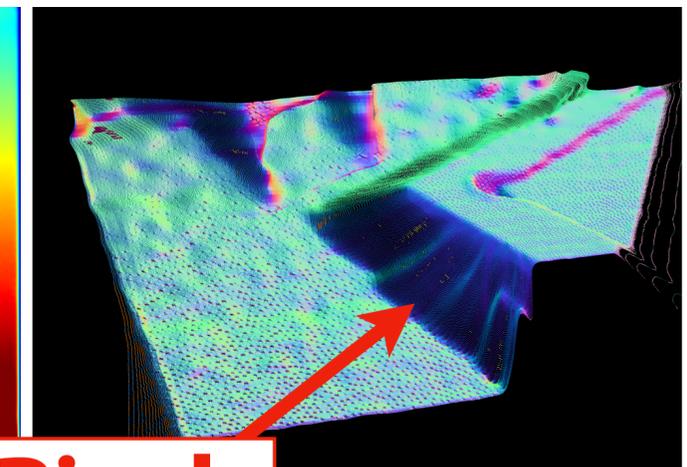
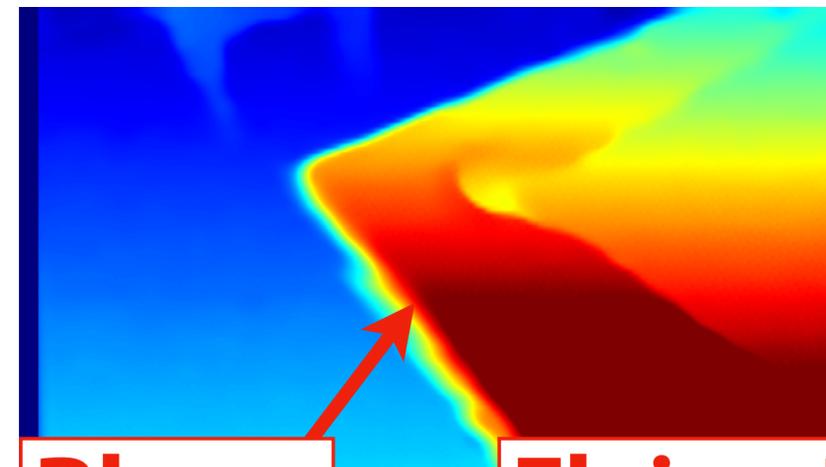
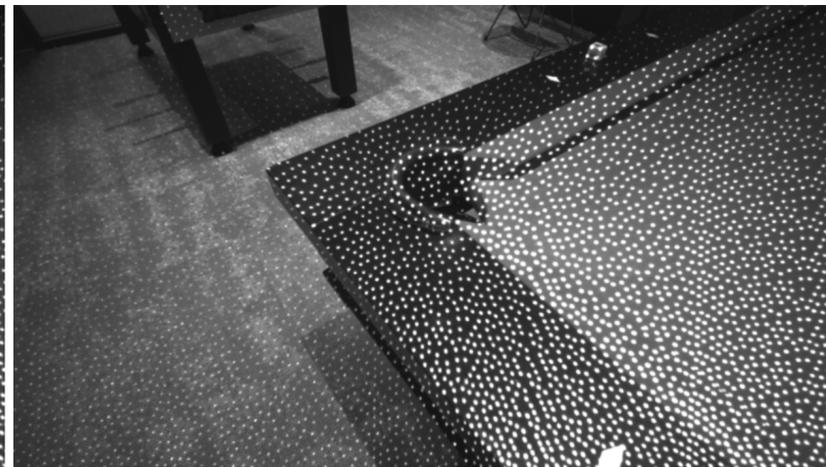
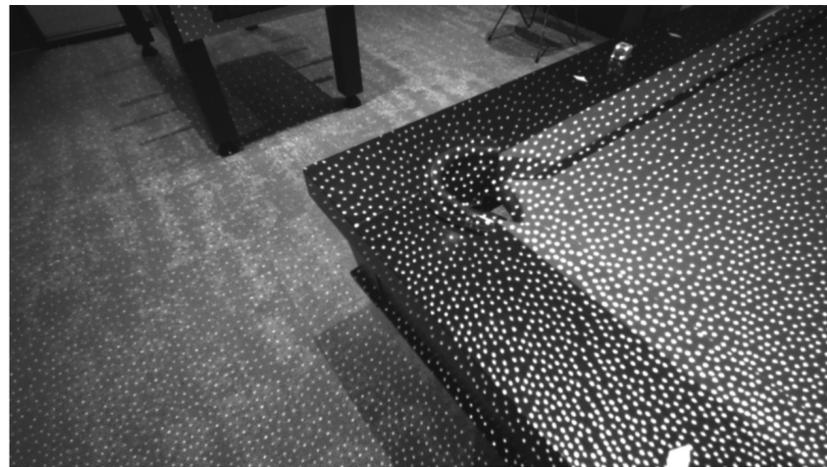
=



Photometric Loss

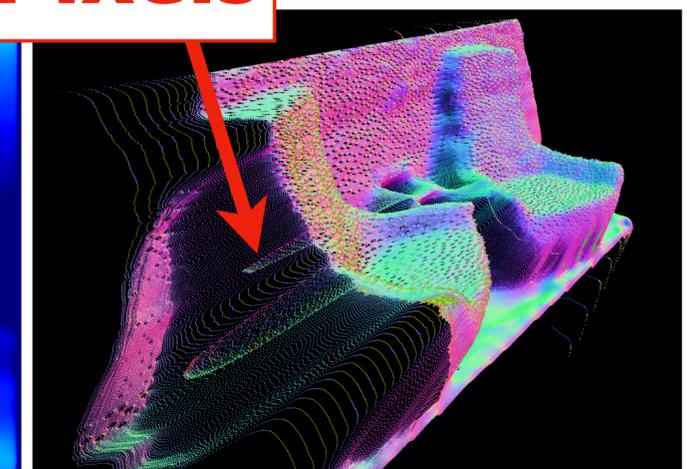
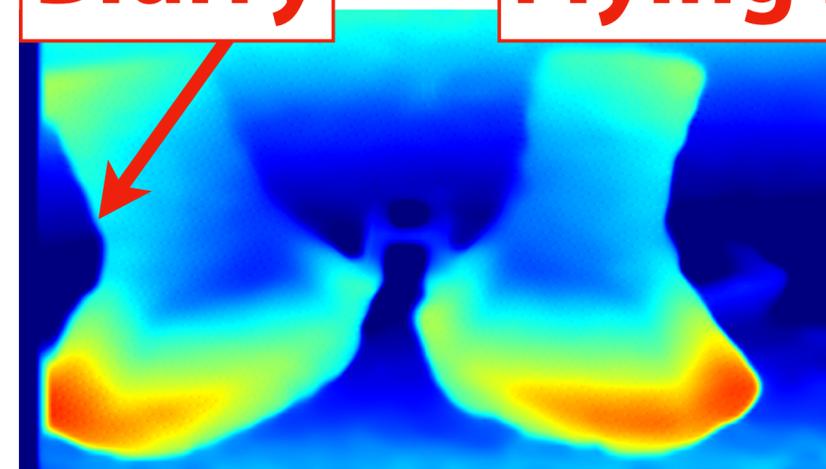
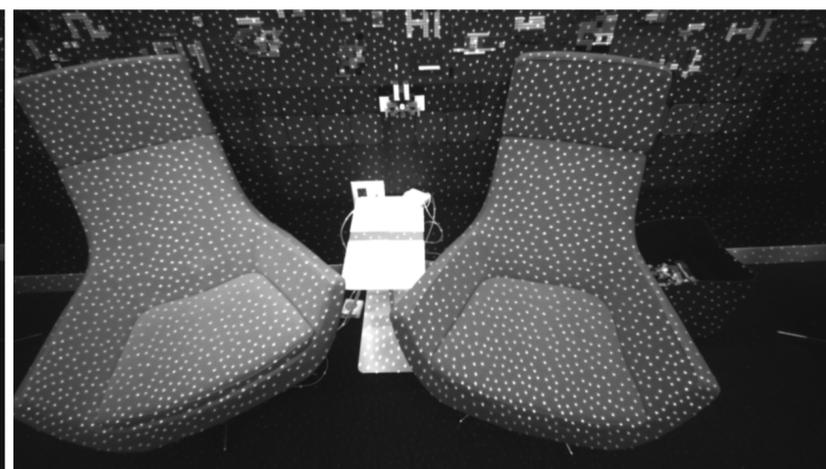
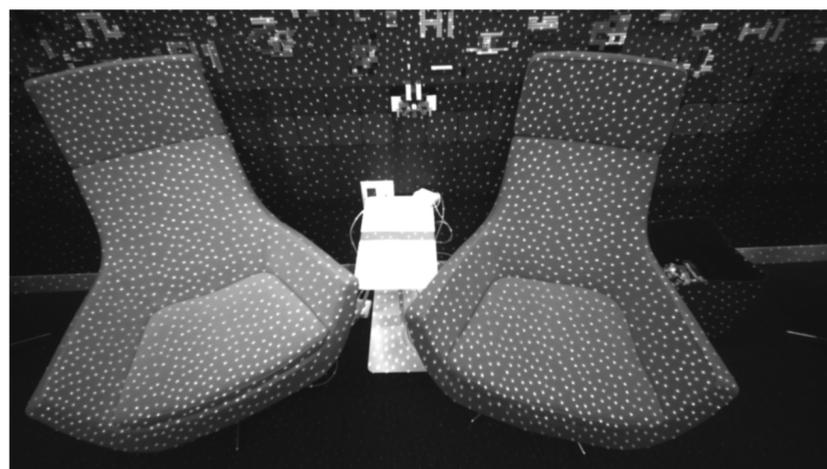
However...

$$\text{Photometric Loss} = | \text{Left View} - \text{Warping}(\text{Right View}, \text{Left Disparity}) |$$



Blurry

Flying Pixels



Left View

Right View

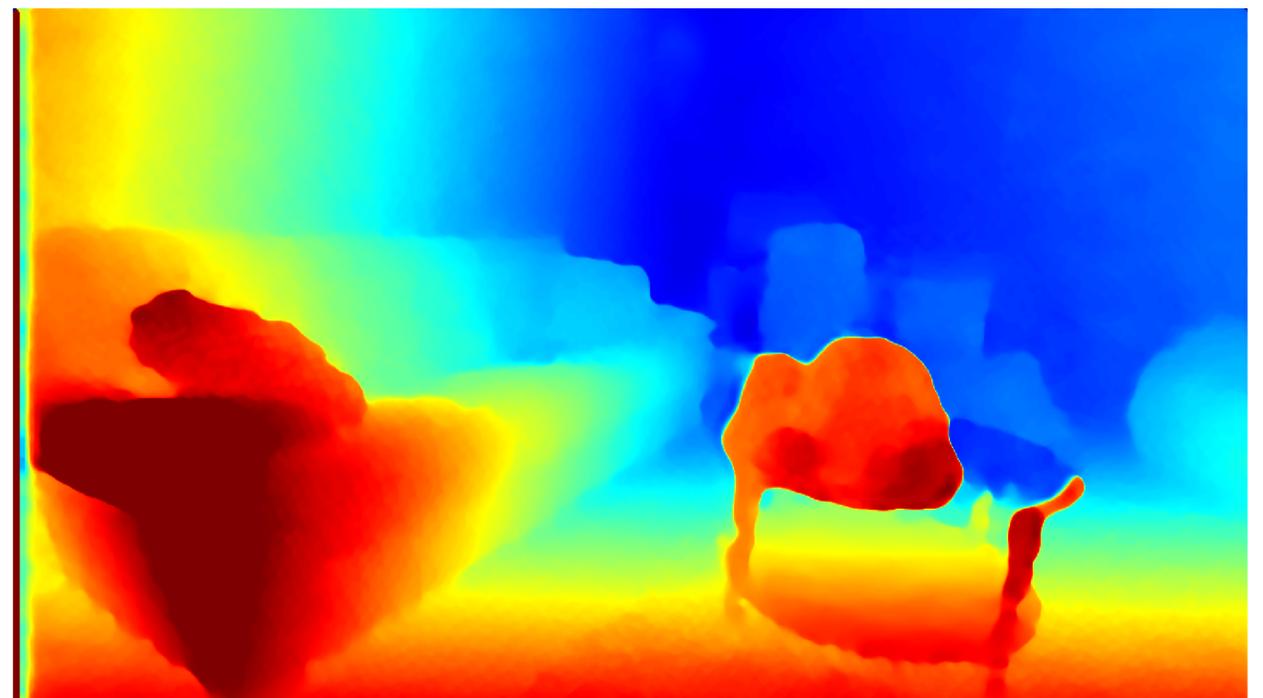
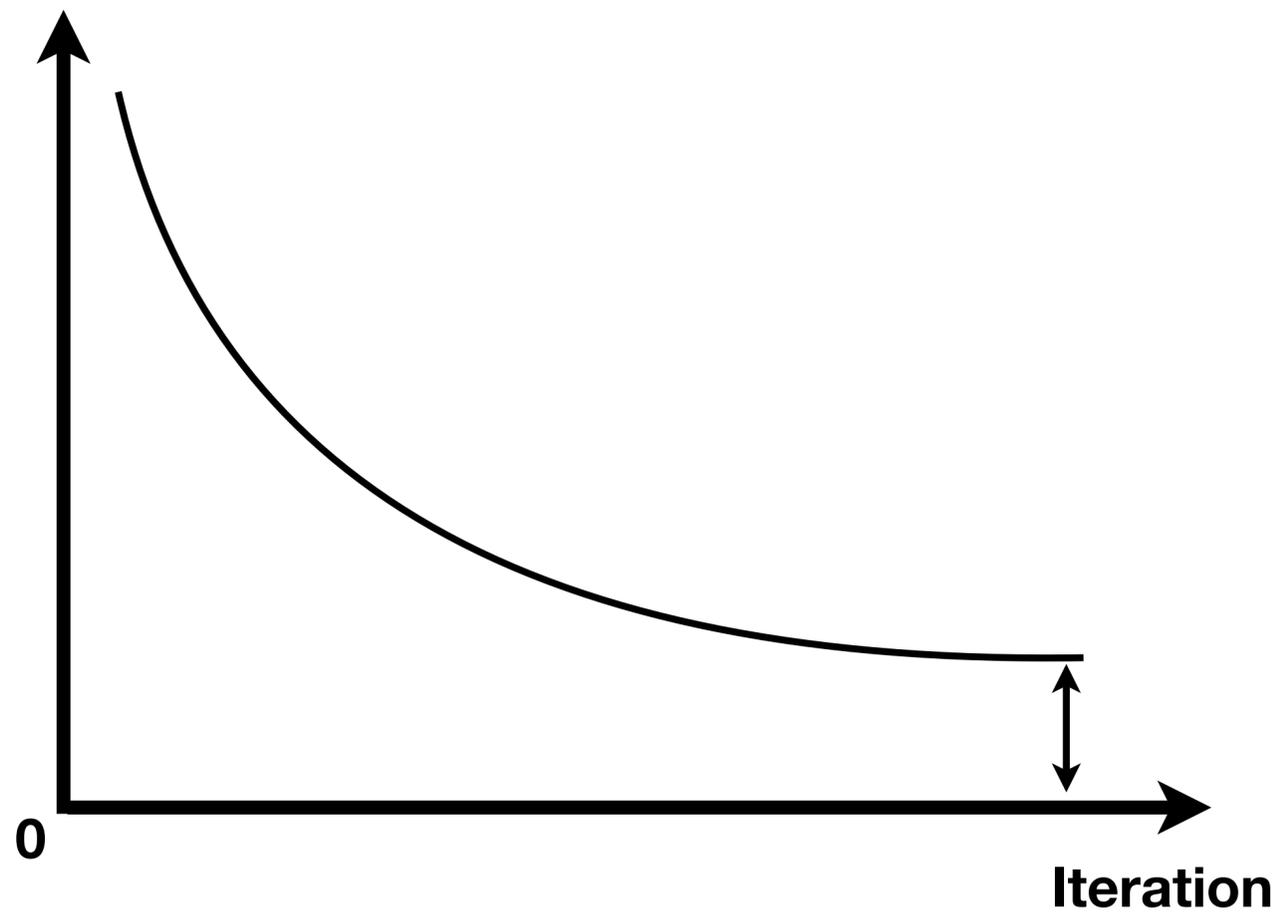
Estimated Disparity

Visualize in 3D

How to fix the problem?

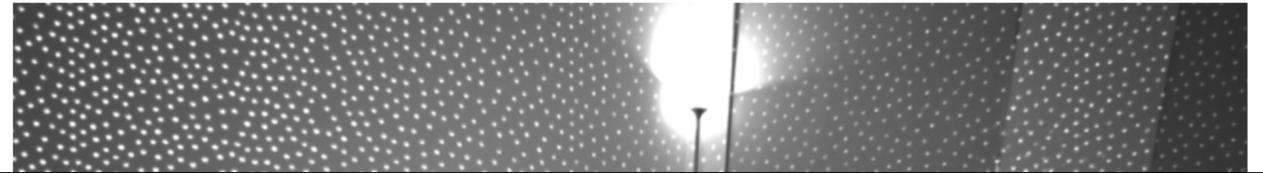
Disparity = argmin(Photometric Loss)

Photometric Loss



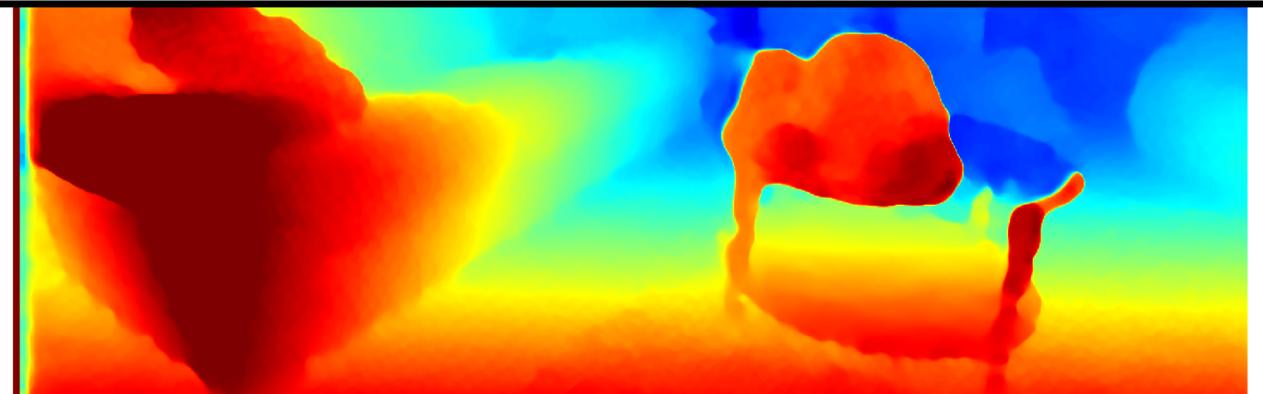
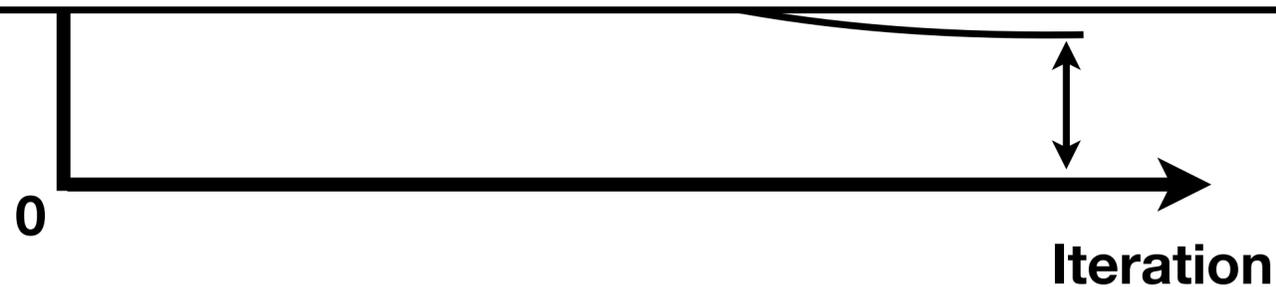
How to fix the problem?

Disparity = $\operatorname{argmin}(\text{Photometric Loss})$



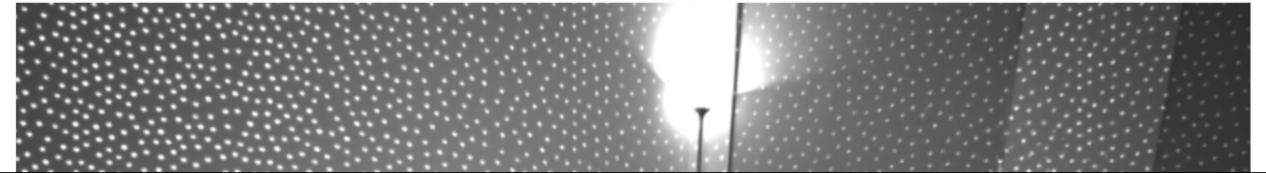
Lower Loss \rightarrow Better Solution?

- ✓ Entropy Loss (Image Classification, Semantic Segmentation)
- ✓ L1/L2 Loss (Depth estimation, Colorization)
- ✗ Photometric Loss for Stereo Matching



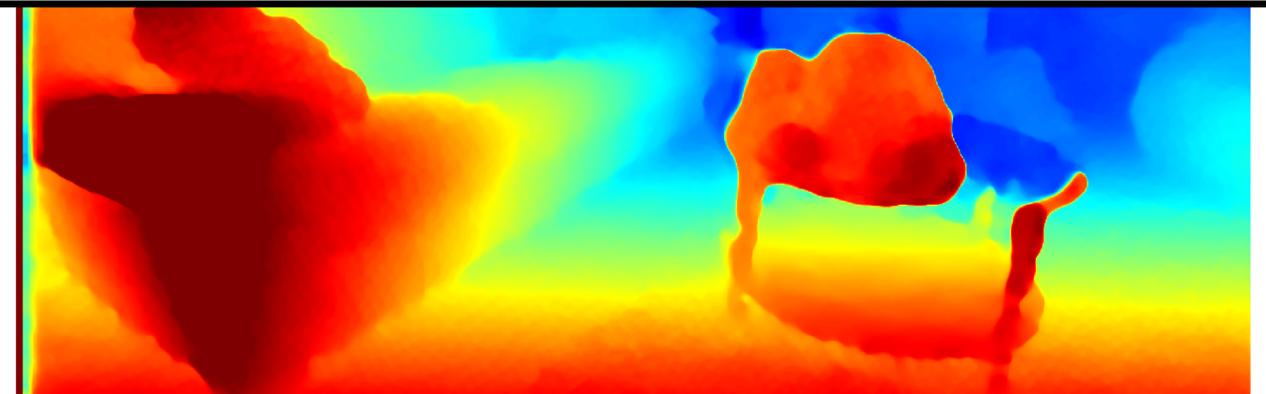
How to fix the problem?

Disparity = $\operatorname{argmin}(\text{Photometric Loss})$



Lower Loss \rightarrow Better Solution?

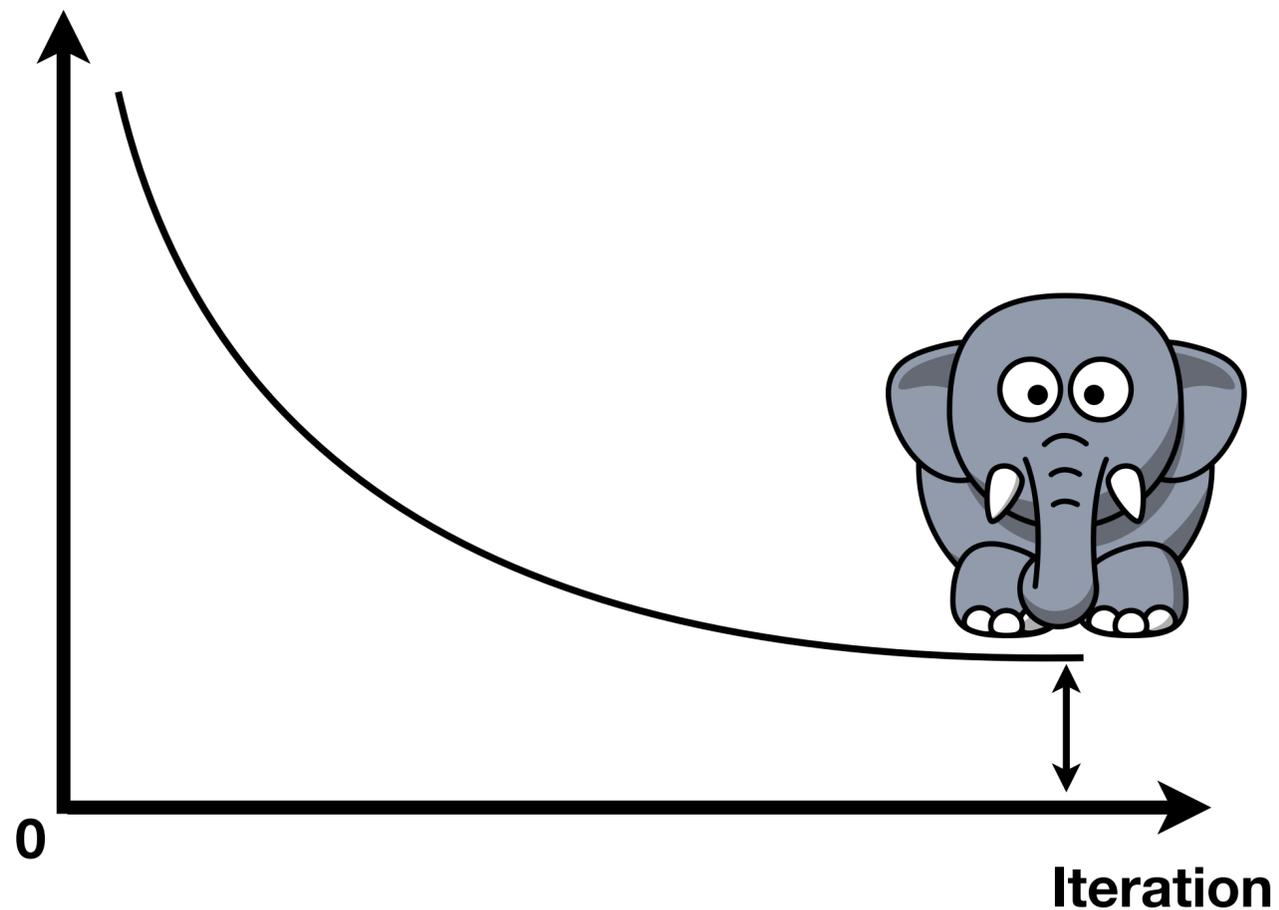
- ✘ Photometric Loss for Stereo Matching is not ideally zero!
 - Different Exposure
 - Occlusion
 - ...



How to fix the problem?

Disparity = $\operatorname{argmin}(\text{Photometric Loss})$

Photometric Loss



**Unnecessary over optimization
hurts performance!**

Fix the loss!

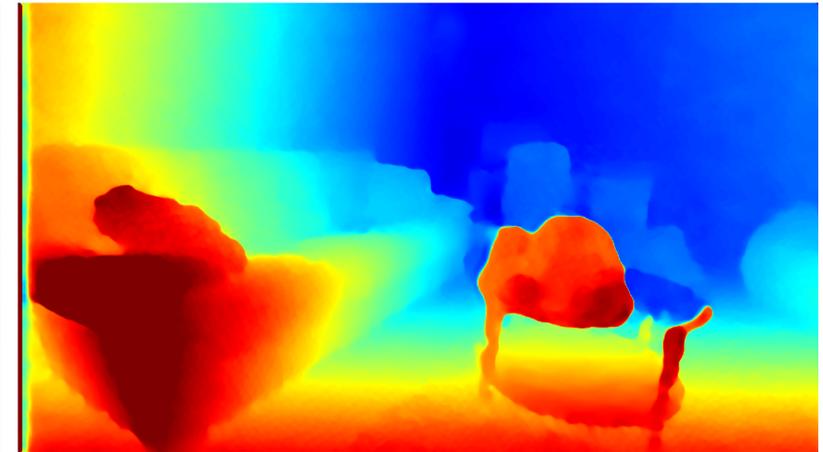
Improving Photometric Loss

1. Remove Unnecessary Dependence.
2. Remove Unexplainable Region.
3. Remove Bad Local Optima.

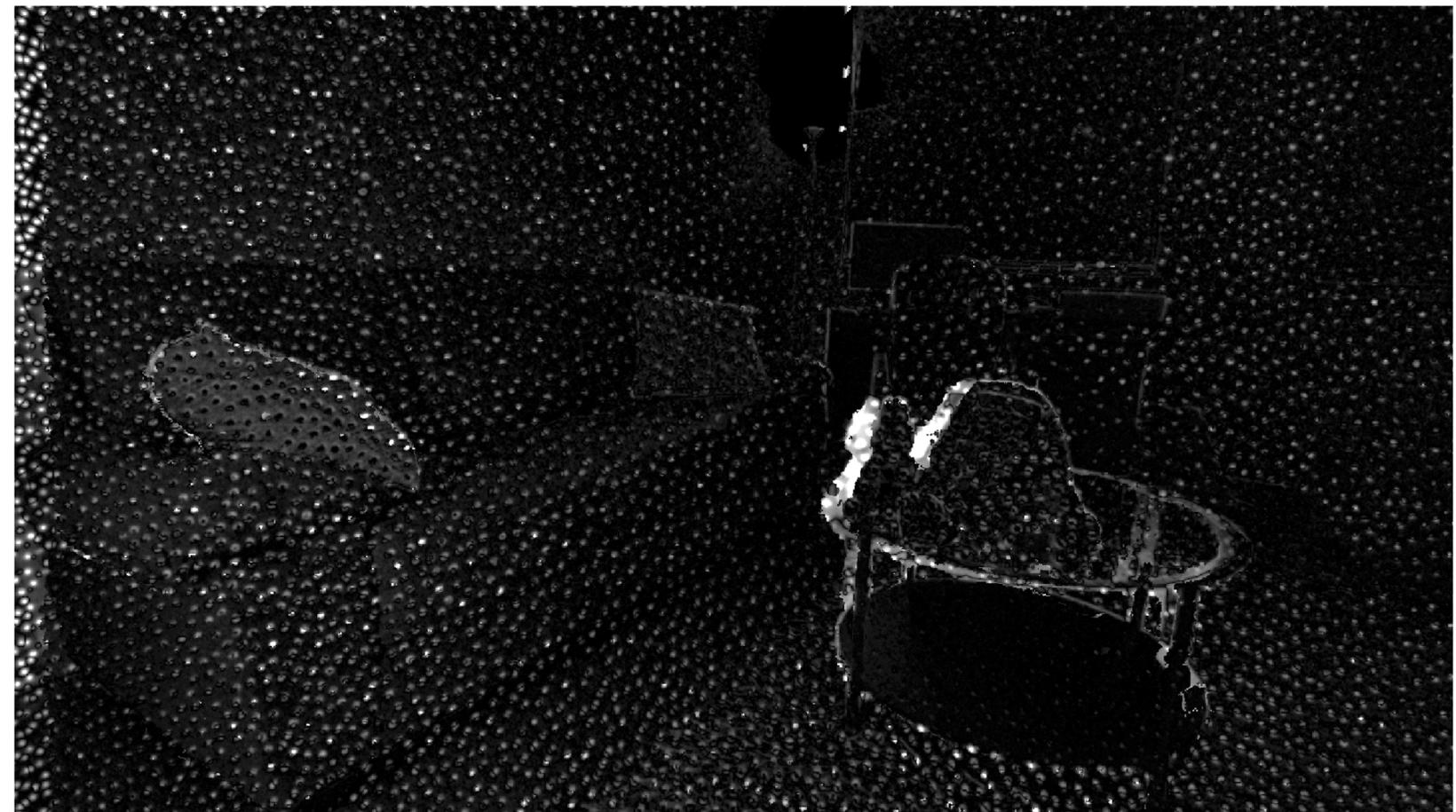
Left View



Disparity



Photometric Loss



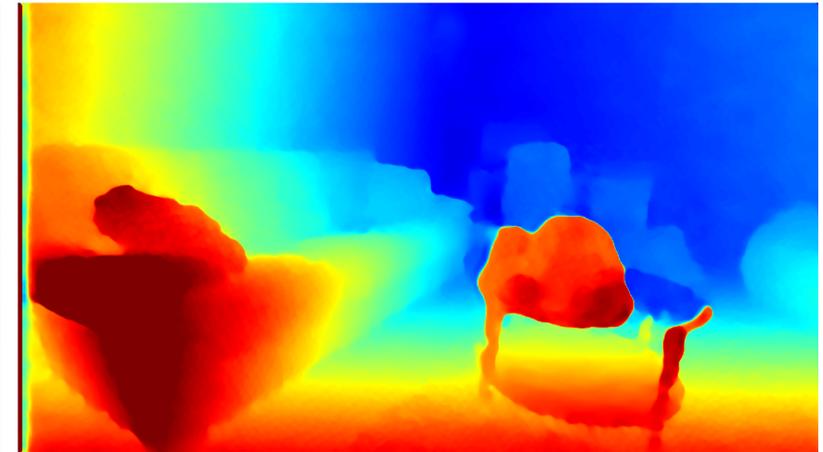
Improving Photometric Loss

1. Remove Unnecessary Dependence.

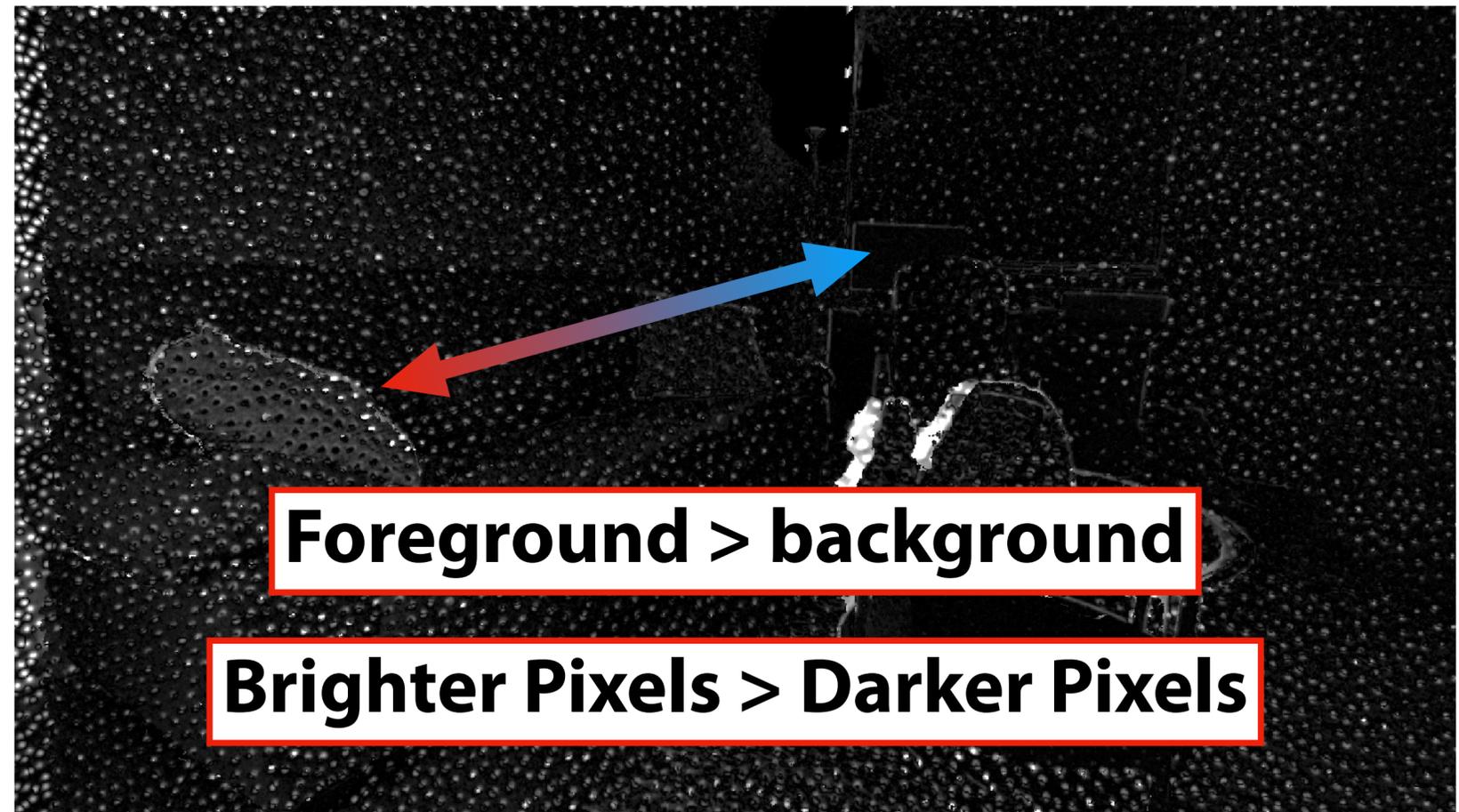
Left View



Disparity



Photometric Loss



Foreground > background

Brighter Pixels > Darker Pixels

Remove Dependence

1. Remove Unnecessary Dependence.

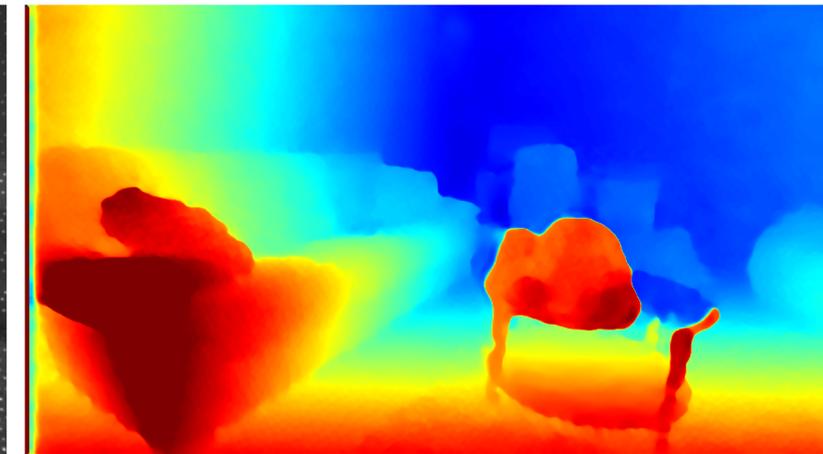
- Local Contrast Normalization

$$I_{LCN} = \frac{I - \mu}{\sigma + \eta}$$

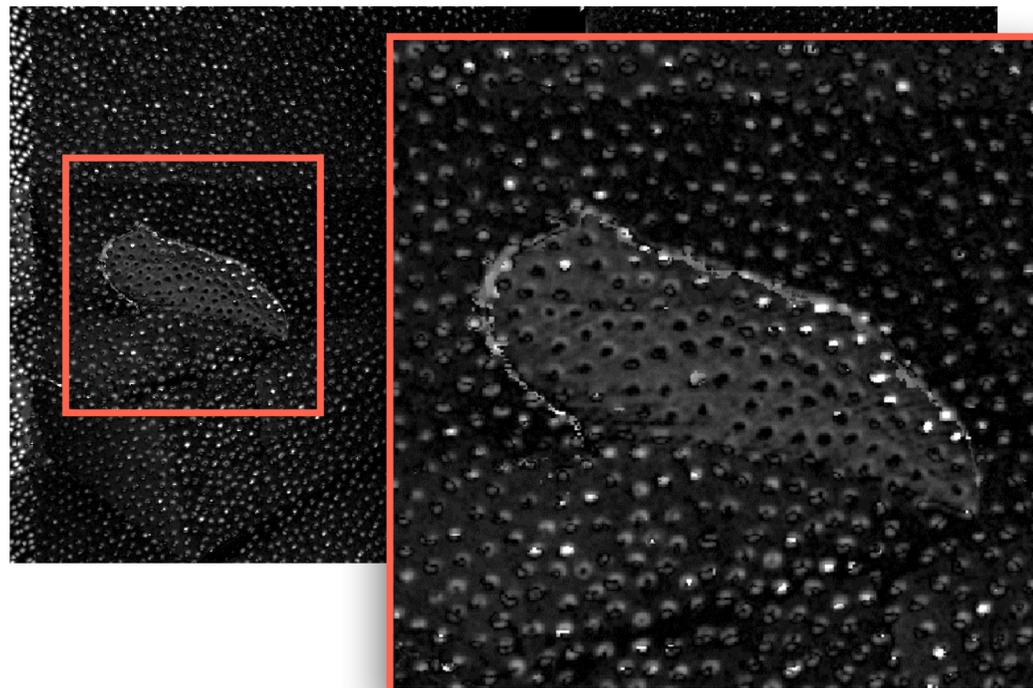
Left View



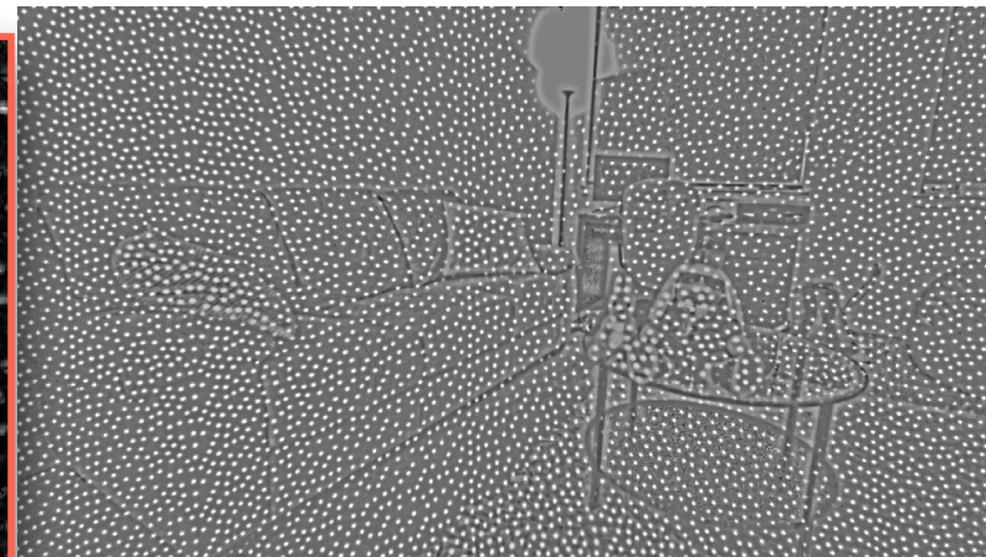
Disparity



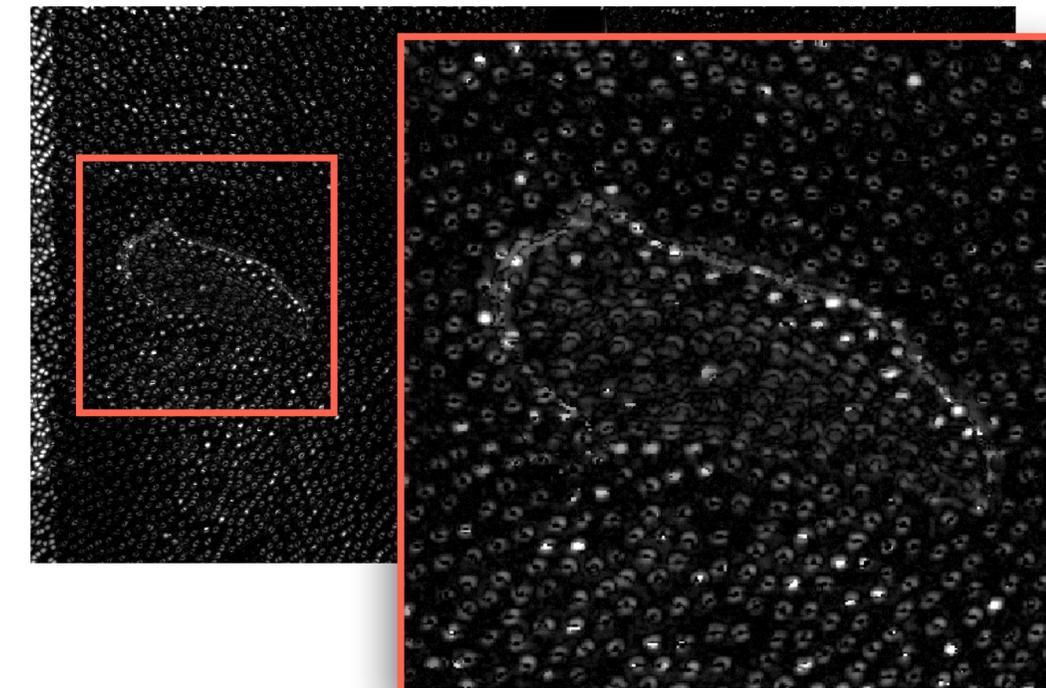
Loss on Raw IR



Local Contrast Normalization

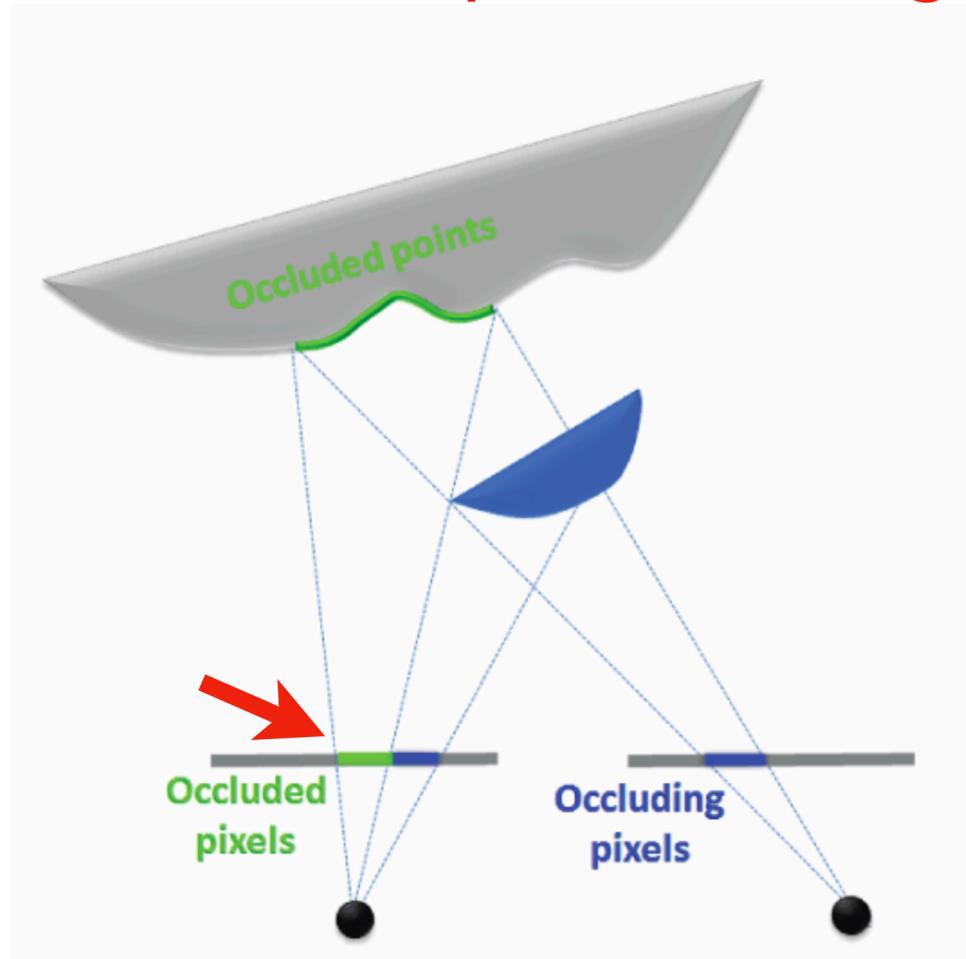


Loss on Normalized IR

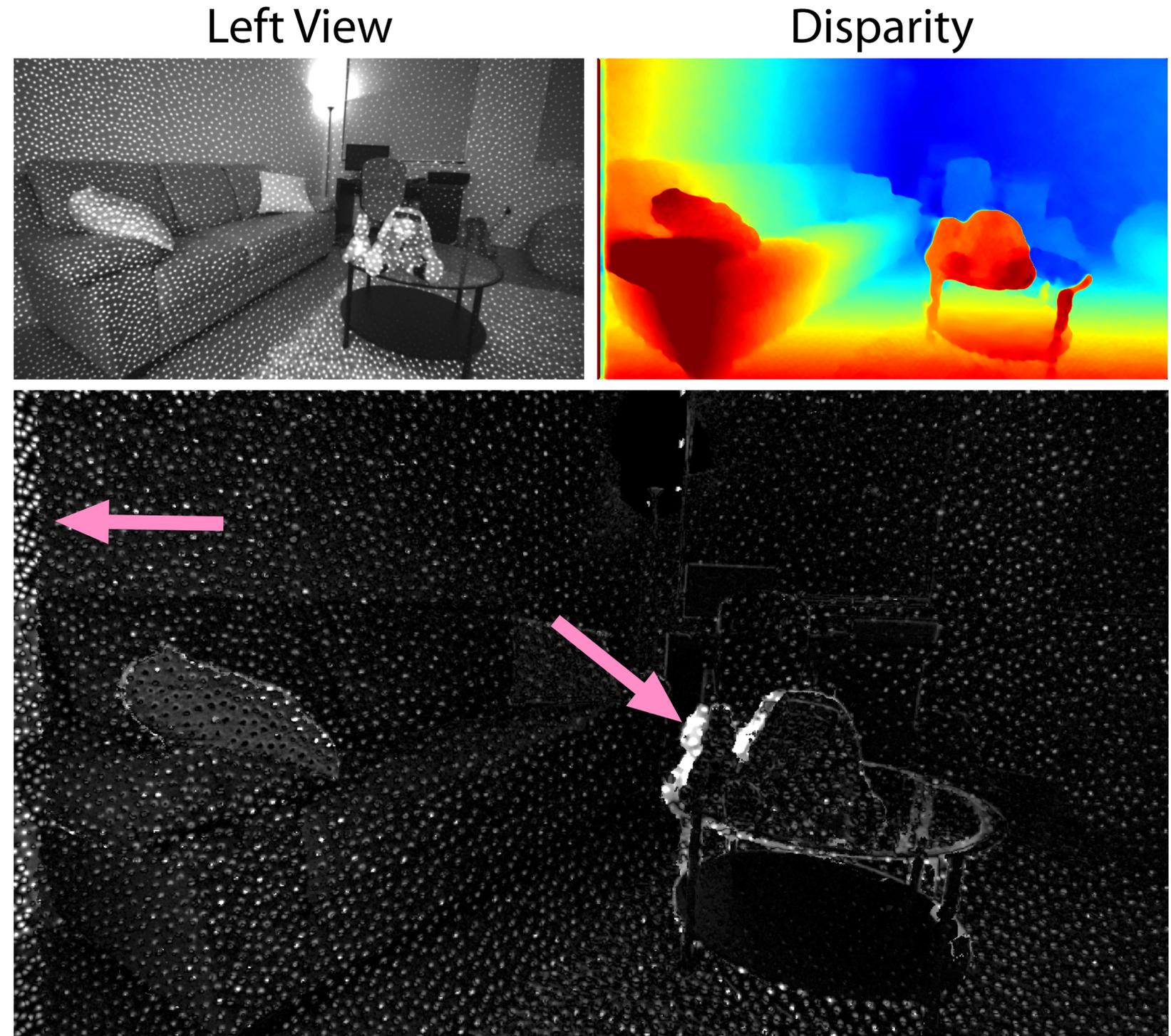


Improving Photometric Loss

1. Remove Unnecessary Dependence.
 - Local Contrast Normalization
2. Remove Unexplainable Region.



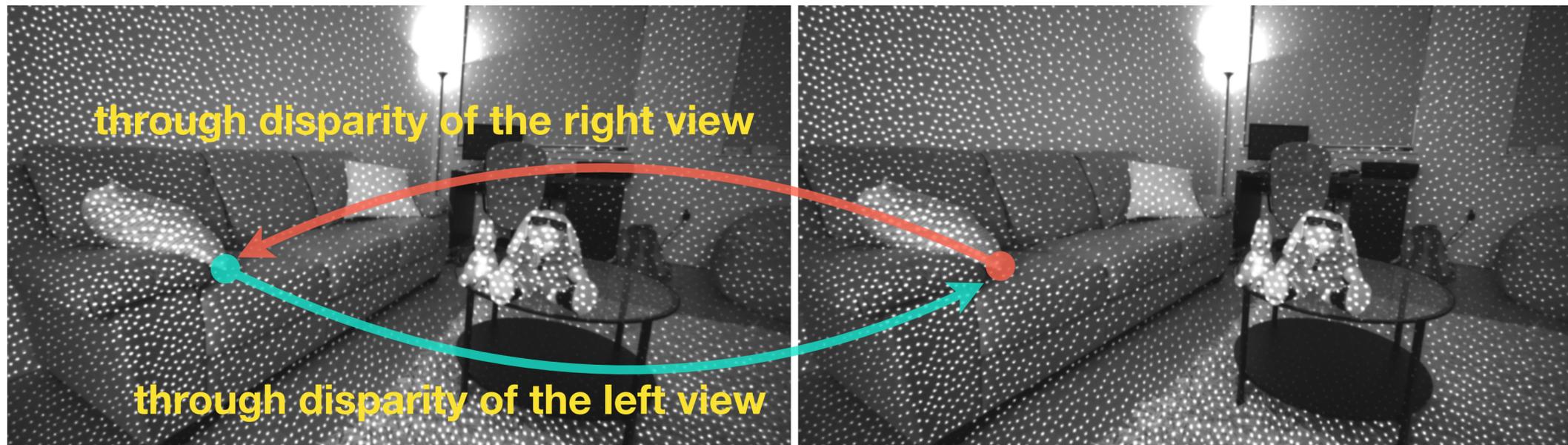
Photometric Loss



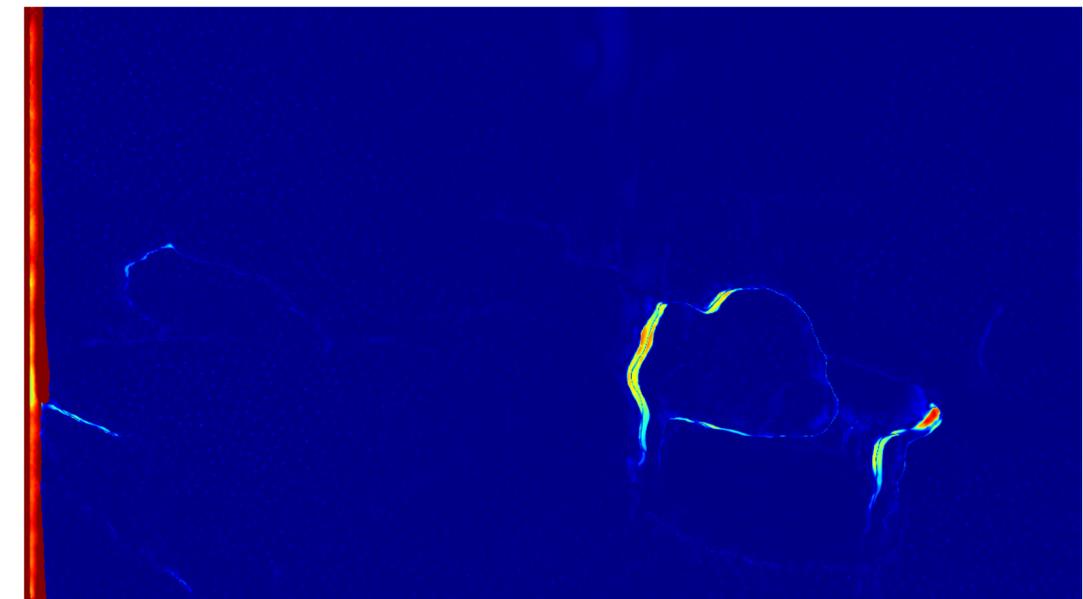
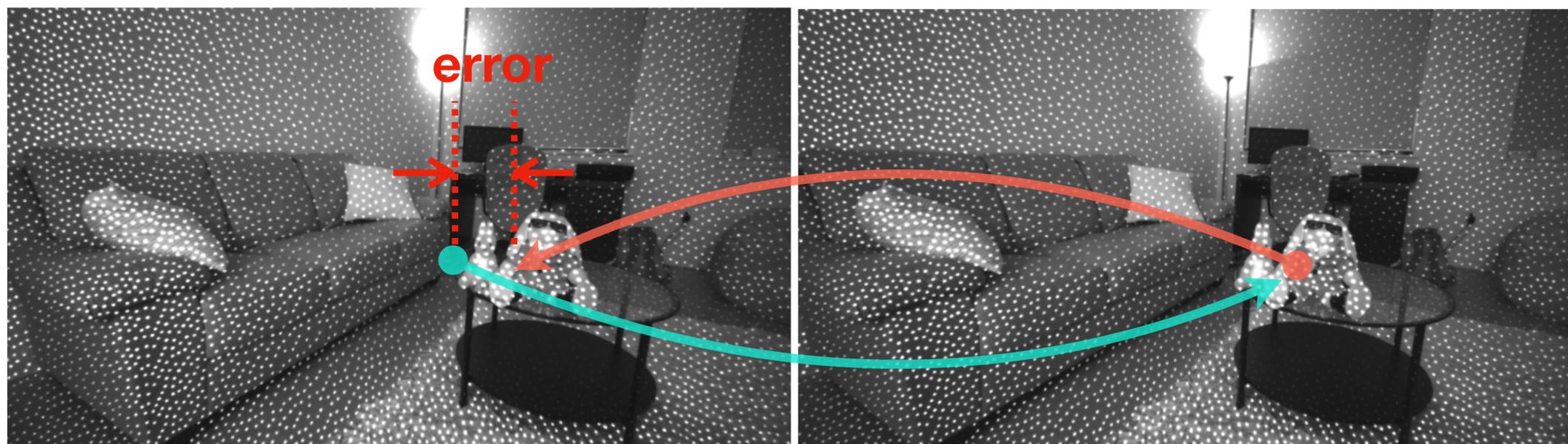
How to find occlusion?

Left View

Right View



Remove the pixel loss if error passes threshold.



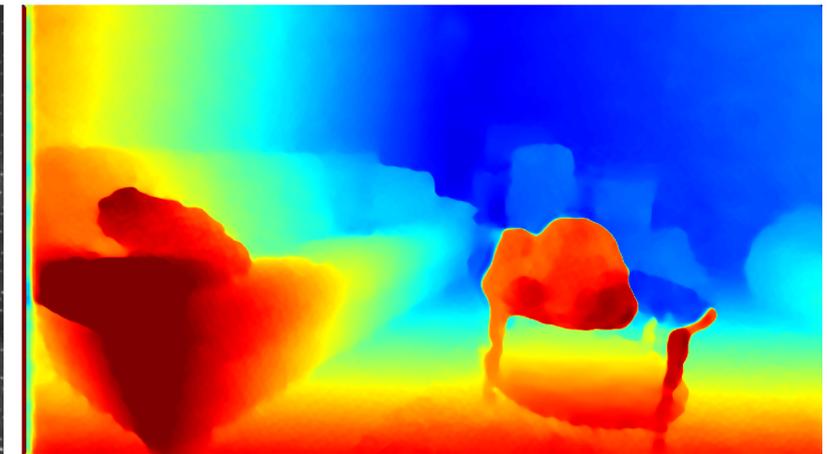
Improving Photometric Loss

1. Remove Unnecessary Dependence.
 - Local Contrast Normalization
2. Remove Unexplainable Region.
 - **Loop-based Check**

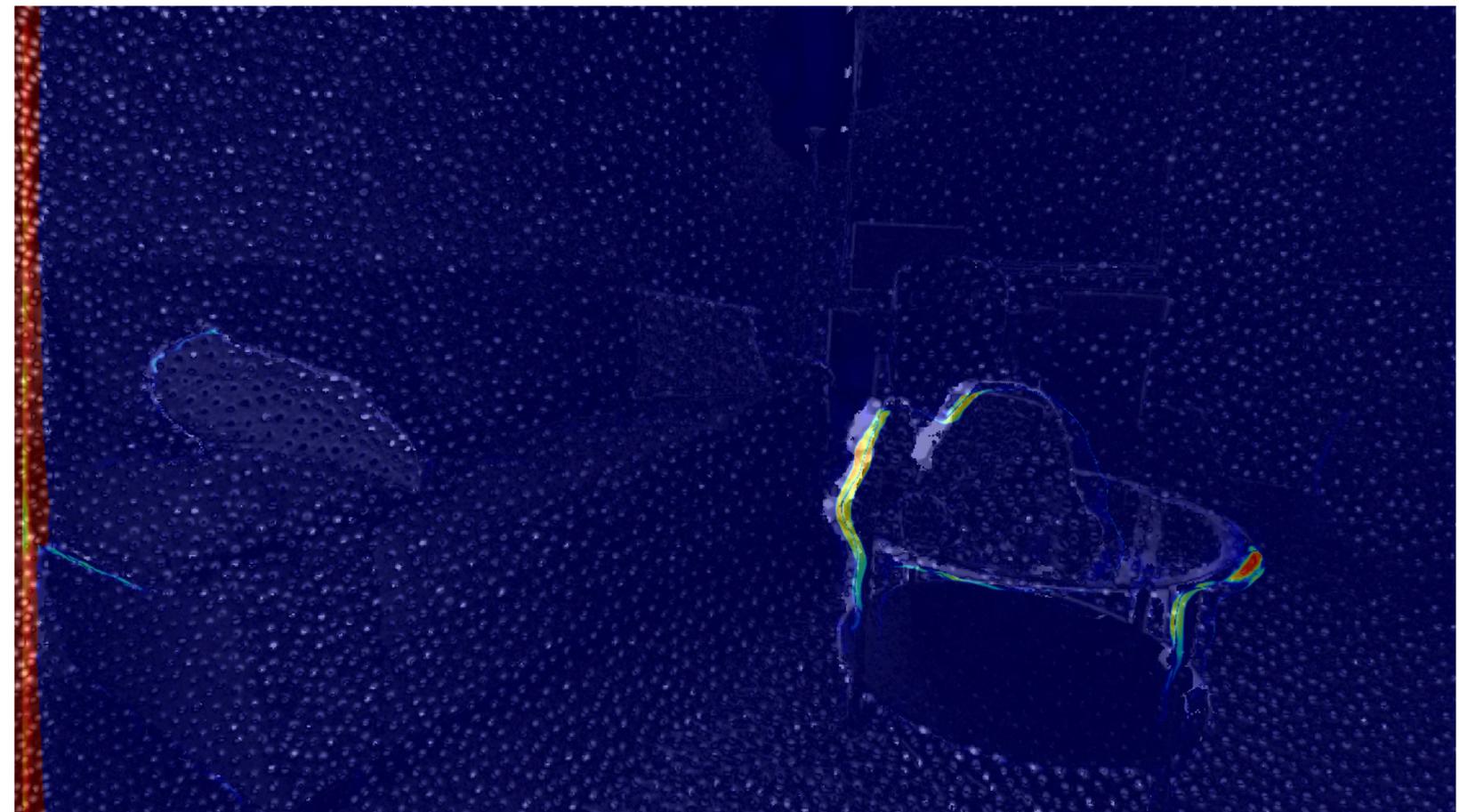
Left View



Disparity



Photometric Loss



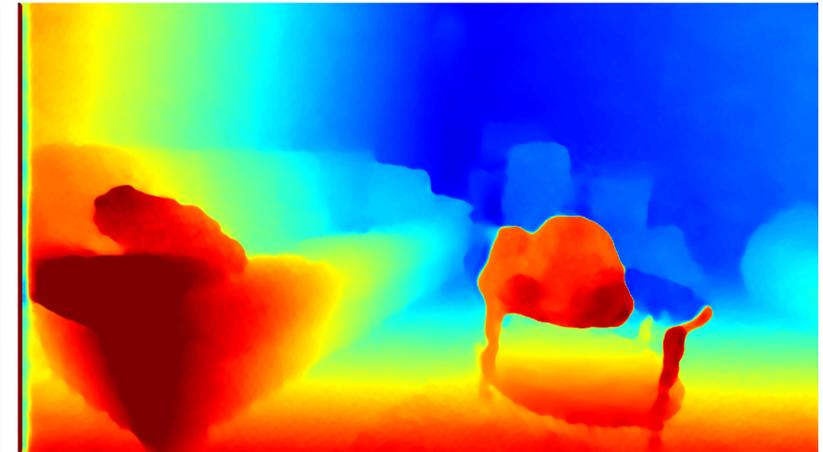
Improving Photometric Loss

1. Remove Unnecessary Dependence.
 - Local Contrast Normalization
2. Remove Unexplainable Region.
 - Loop-based Check
3. Remove Bad Local Optima.

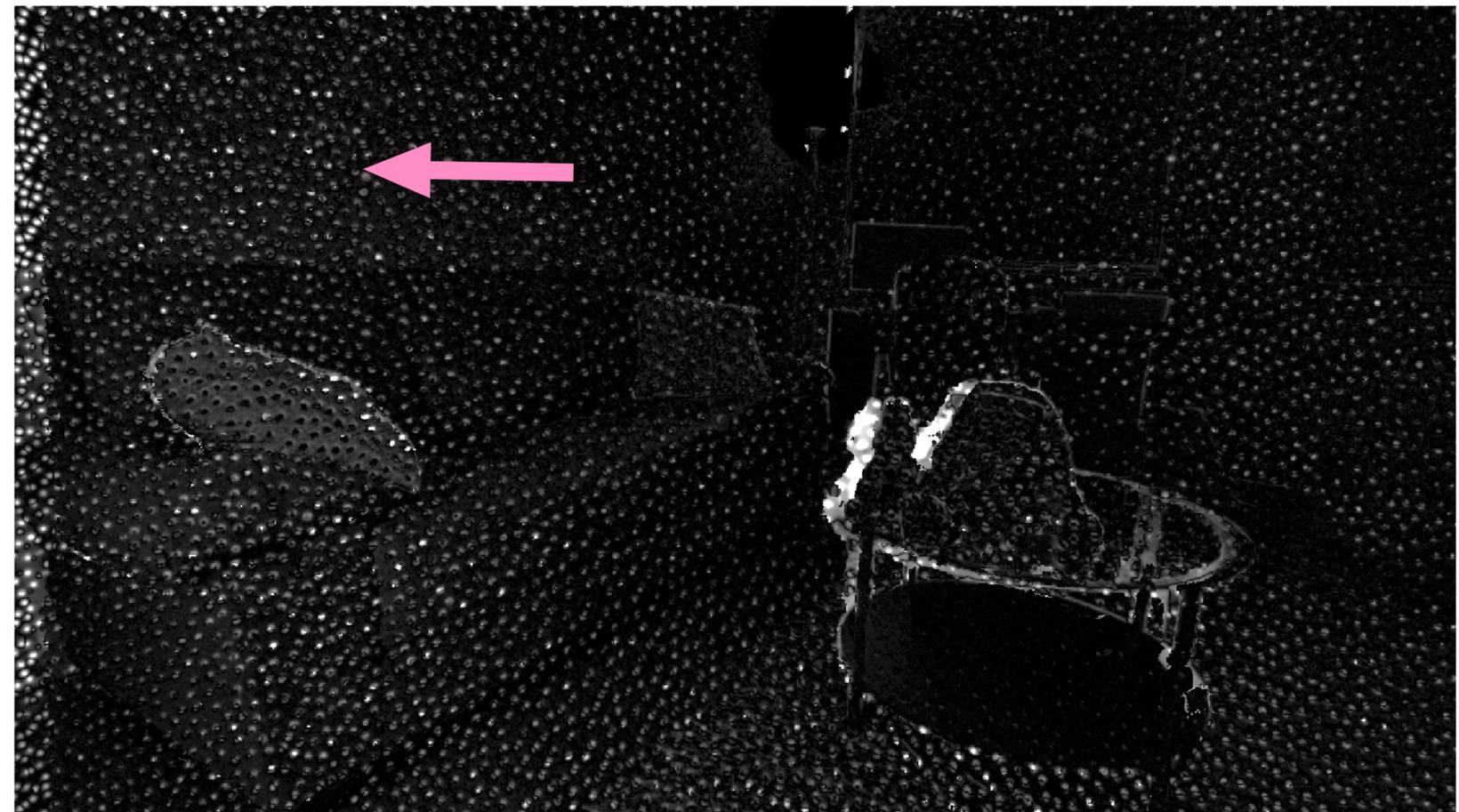
Left View



Disparity

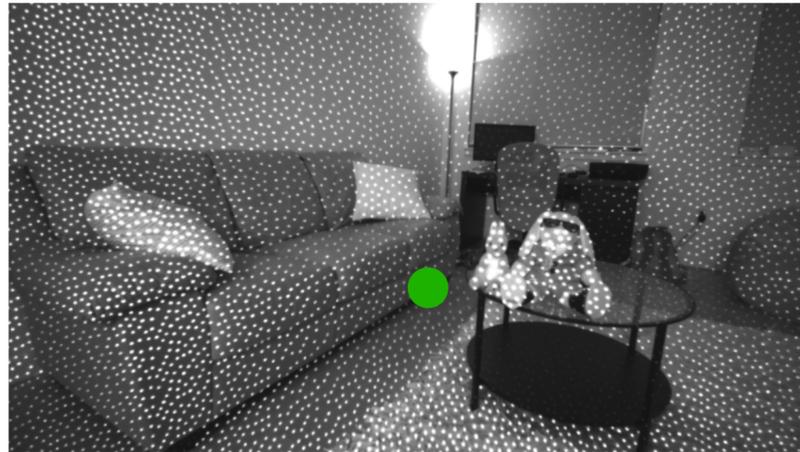


Photometric Loss

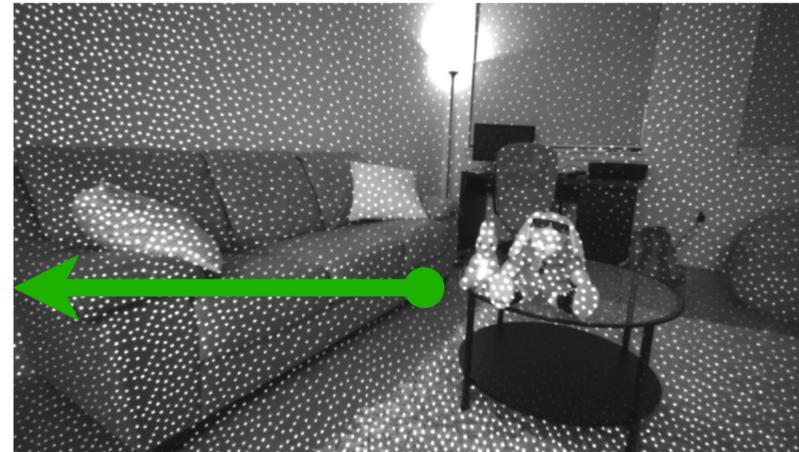


Remove Bad Optima

Left View



Right View



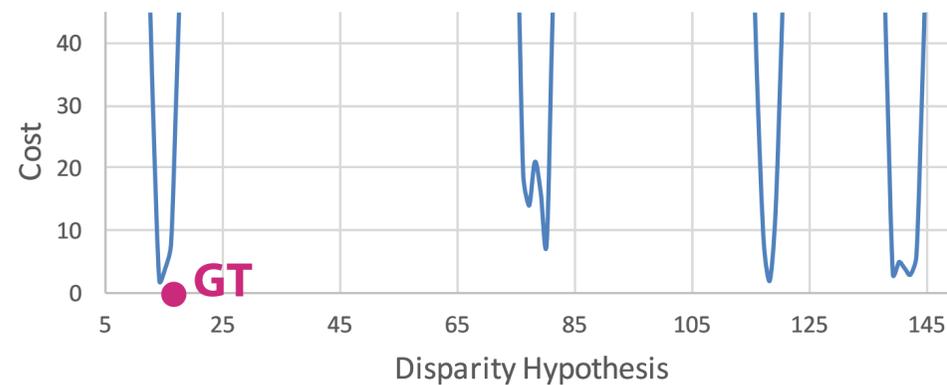
Adaptive Support Weight

$$C_{i,j} = \|\sigma_{i,j}(I_{LCN\ i,j} - \hat{I}_{LCN\ i,j})\|_1$$

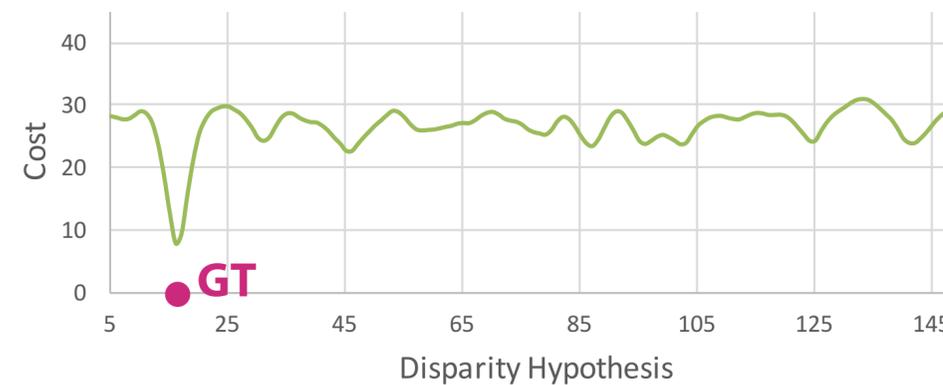
$$\hat{C}_{i,j} = \frac{\sum_{x=i-k}^{i+k-1} \sum_{y=j-k}^{j+k-1} \omega_{x,y} C_{x,y}}{\sum_{x=i-k}^{i+k-1} \sum_{y=j-k}^{j+k-1} \omega_{x,y}}$$

$$\omega_{x,y} = \exp\left(-\frac{|I_{i,j} - I_{x,y}|}{\sigma_\omega}\right)$$

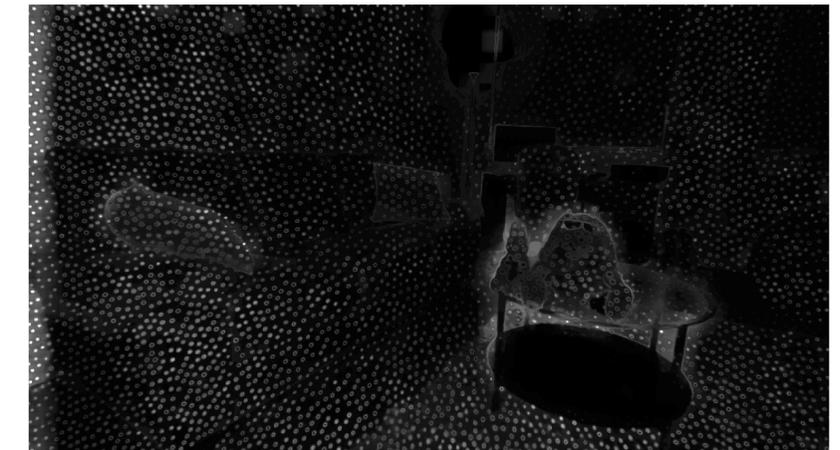
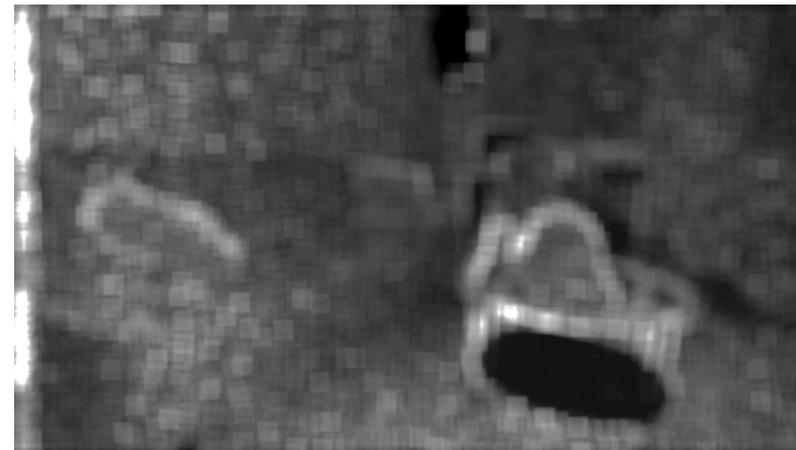
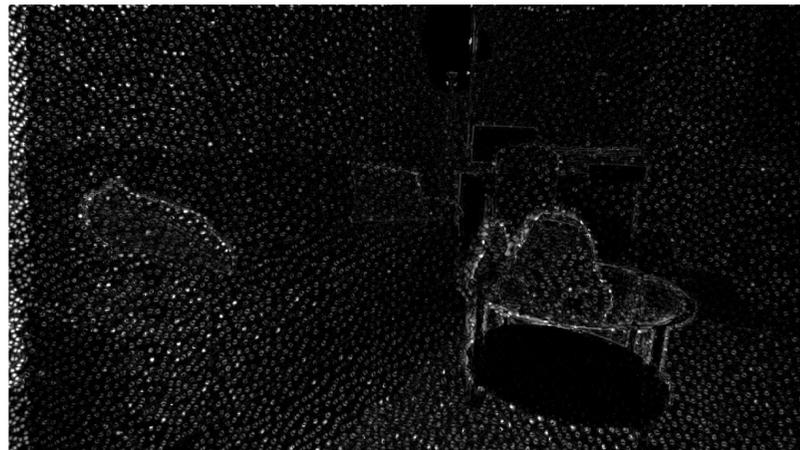
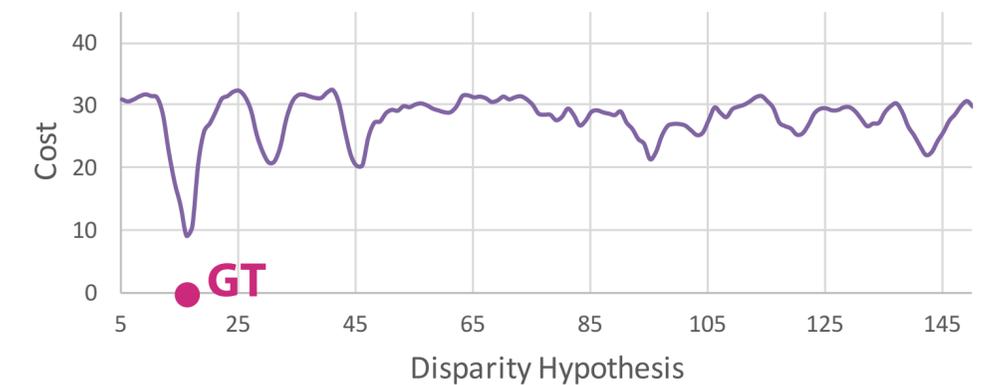
1x1 Window



32x32 Window



32x32 Window with ASW



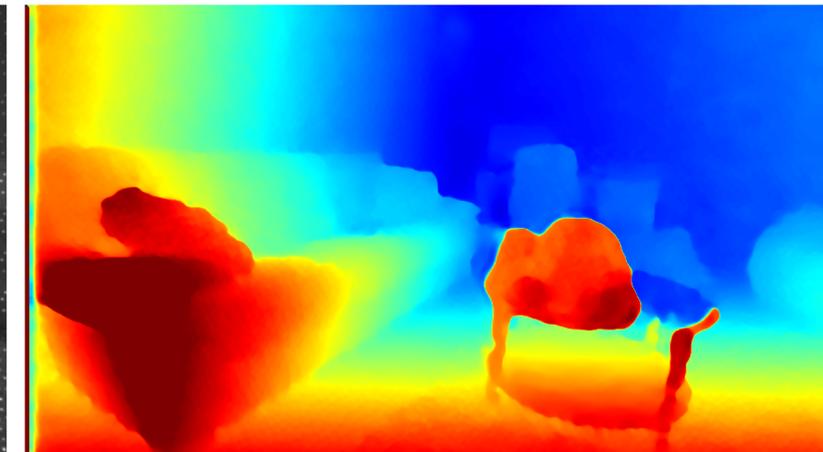
Improving Photometric Loss

1. Remove Unnecessary Dependence.
 - Local Contrast Normalization
2. Remove Unexplainable Region.
 - Loop-based Check
3. Remove Bad Local Optima.
 - **Window aggregation with ASW**

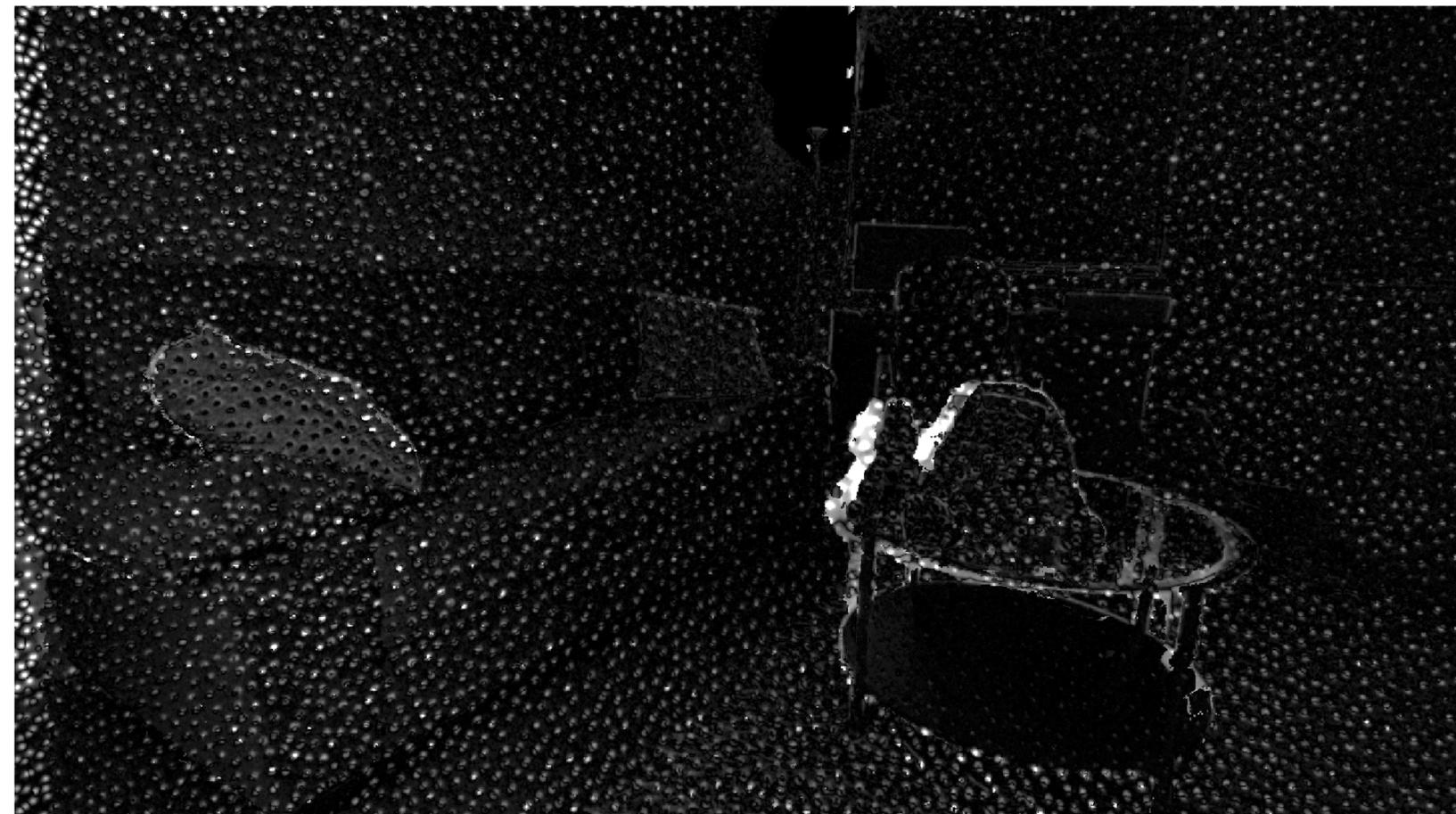
Left View



Disparity



Photometric Loss



Improving Photometric Loss

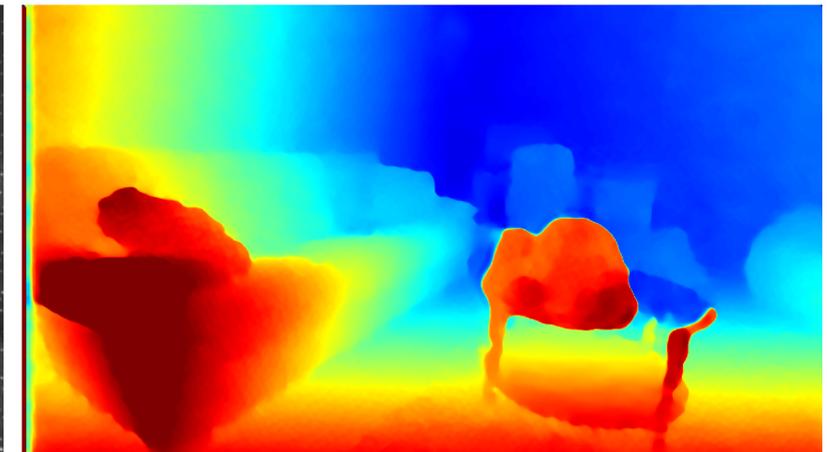
1. Remove Unnecessary Dependence.
 - Local Contrast Normalization
2. Remove Unexplainable Region.
 - Loop-based Check
3. Remove Bad Local Optima.
 - Window aggregation with ASW

$$L = \sum_{p \notin Inv} \hat{C}_p$$

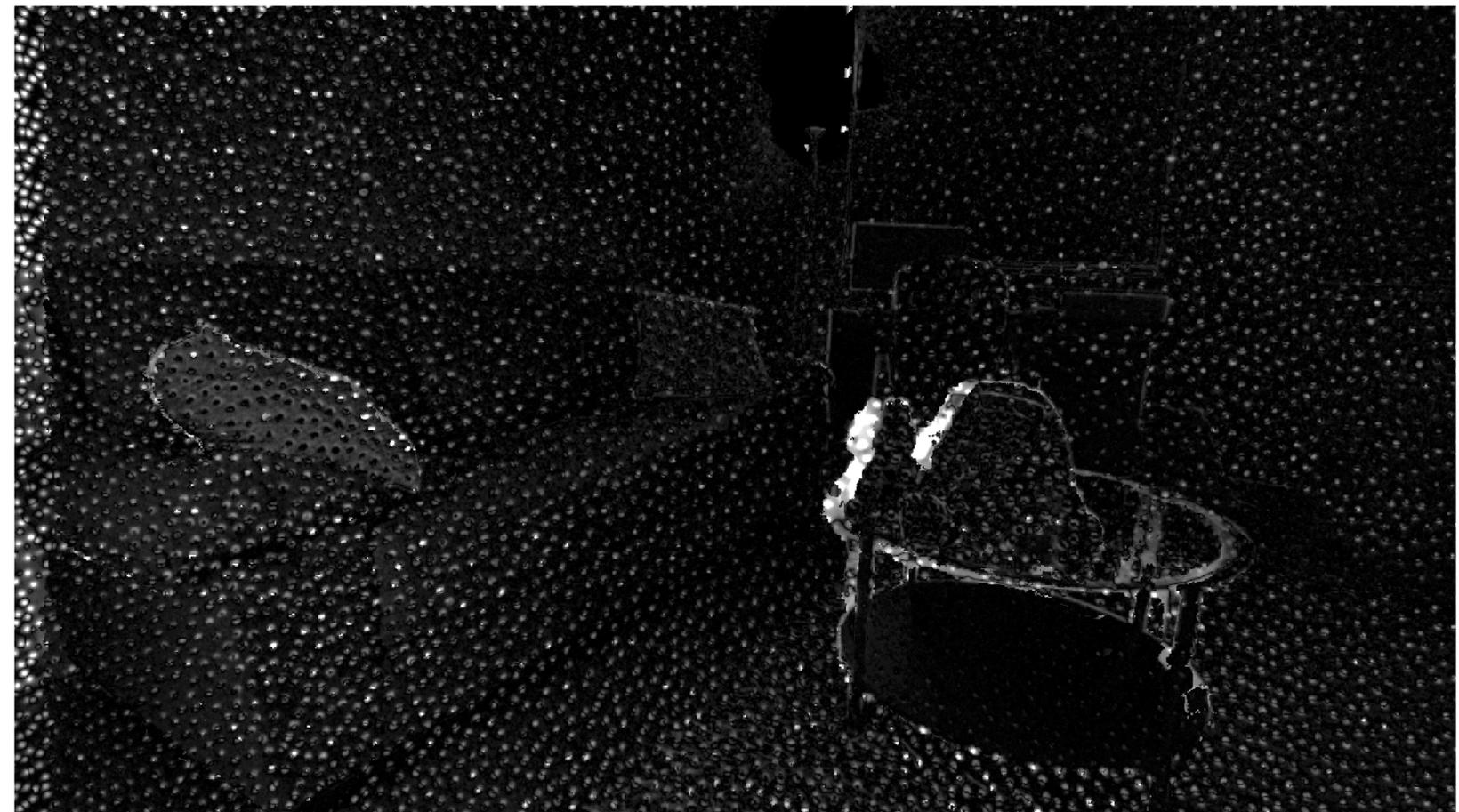
Left View



Disparity

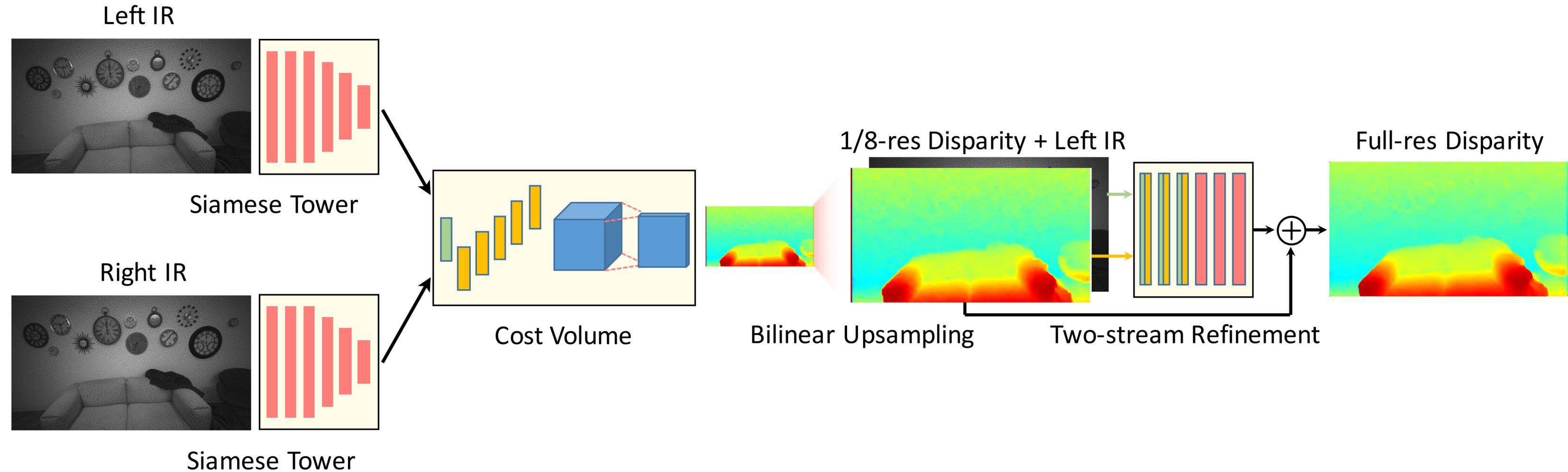


Photometric Loss

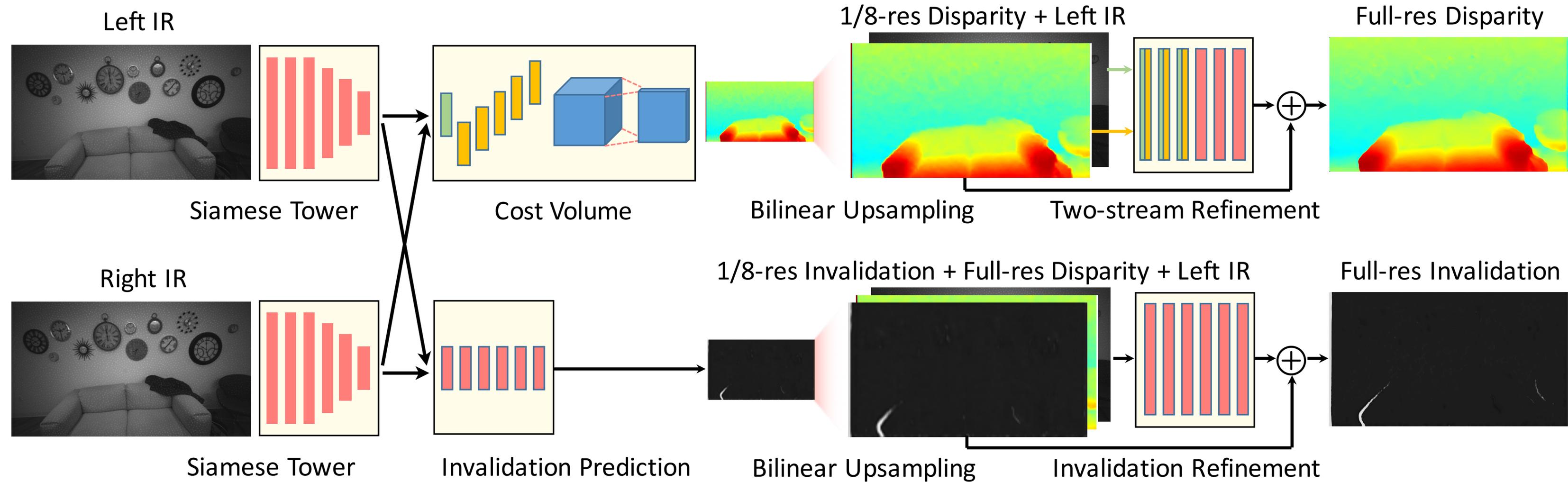


Network Architecture

Network Architecture

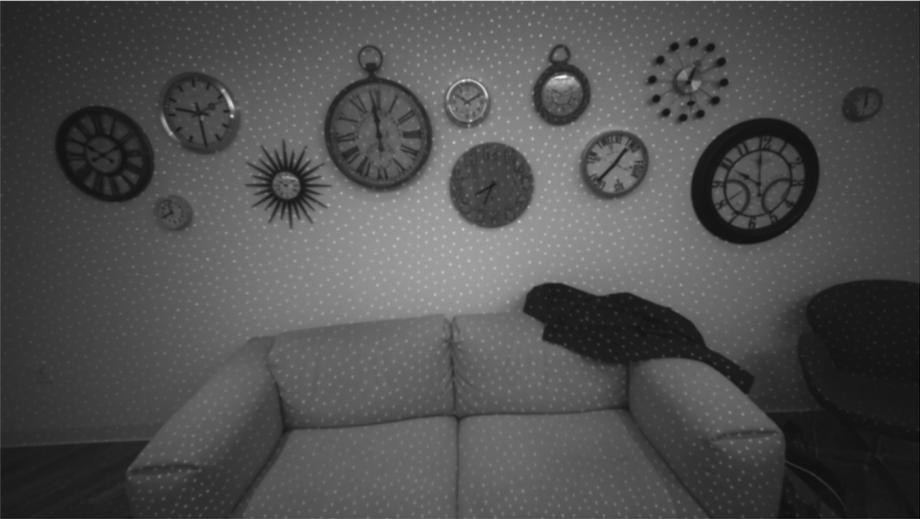


Network Architecture



Experiments

Experiments



Left IR Image

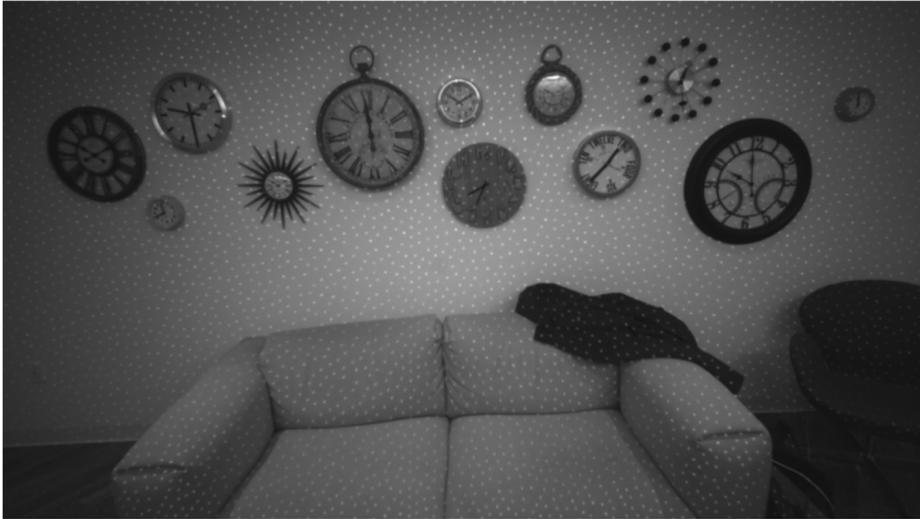


Intel RealSense D435

IR Stereo Camera

IR Projector

Color Camera



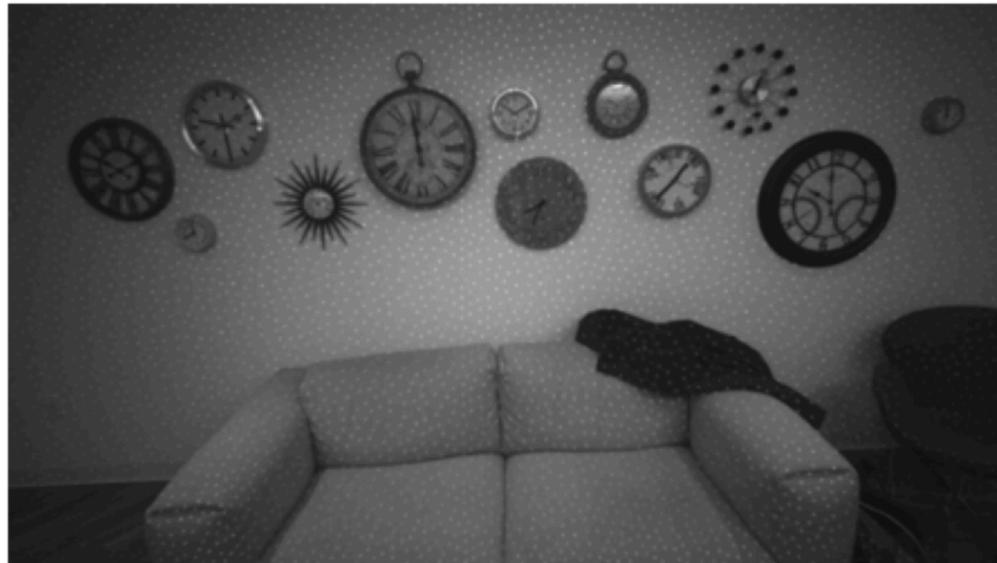
Right IR Image



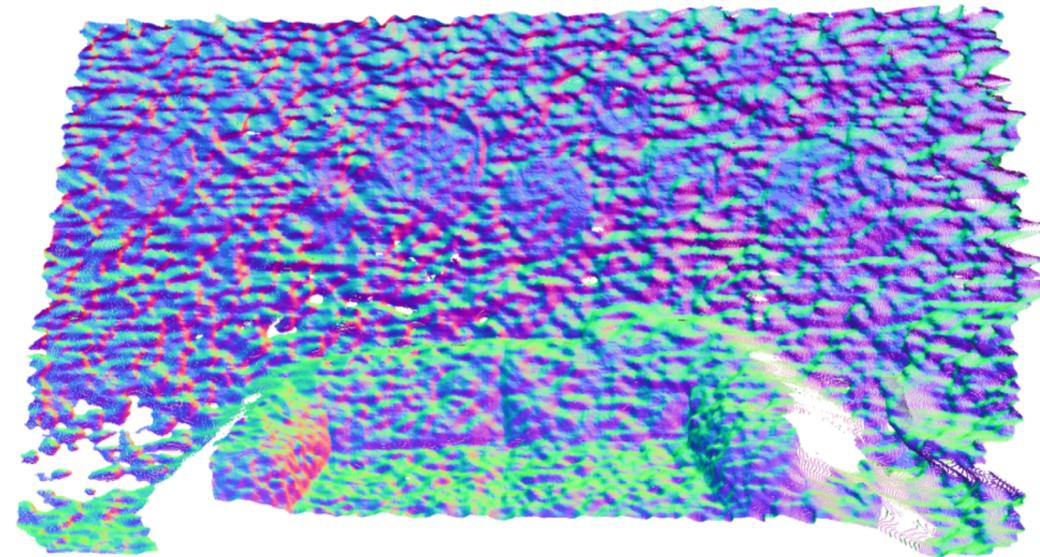
Color Image

Experiments

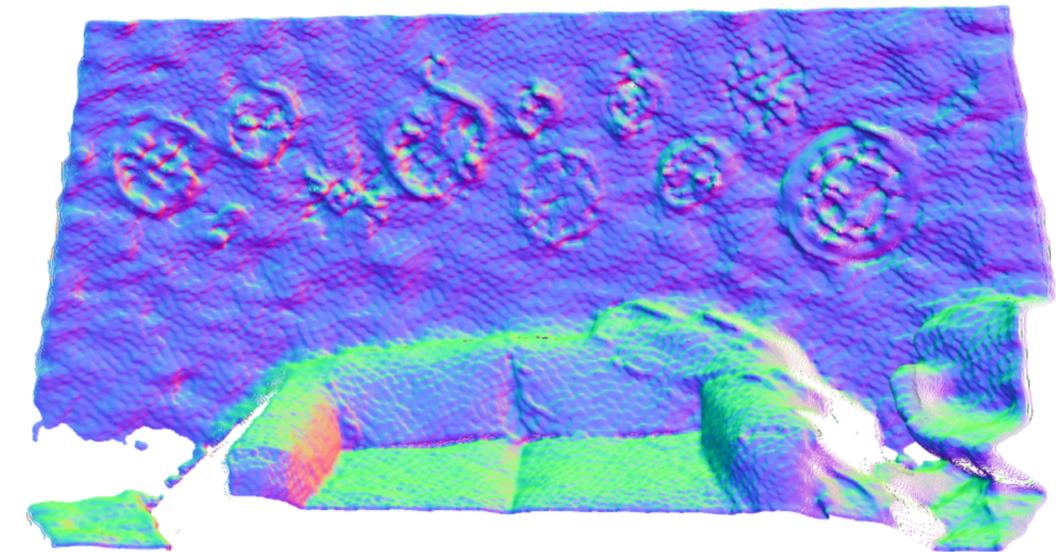
Intel RealSense D435



Left View



Intel RealSense D435



ActiveStereoNet



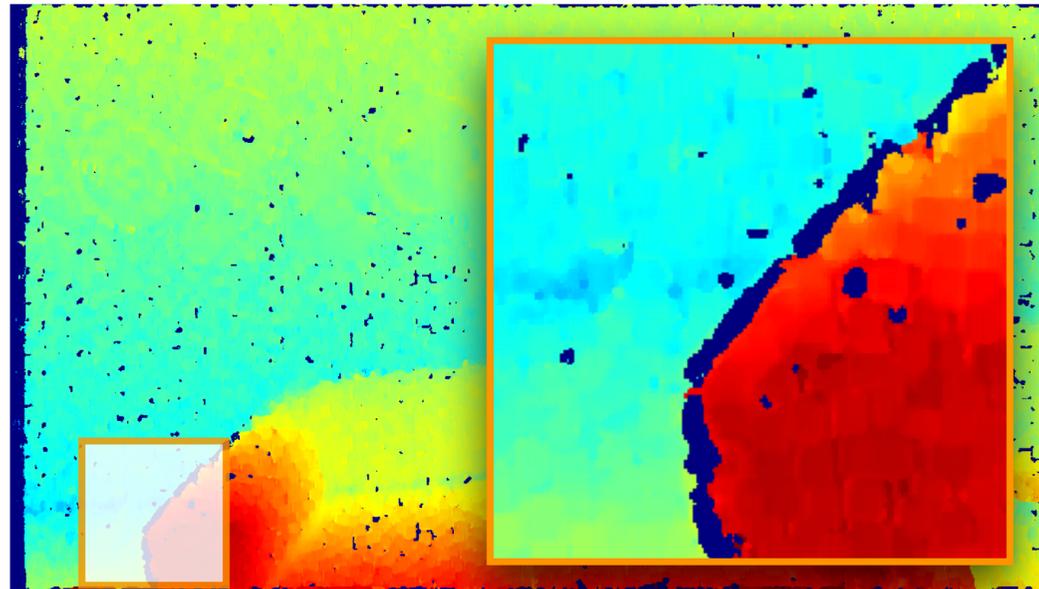
Color Image

Disparity Qualitative Result

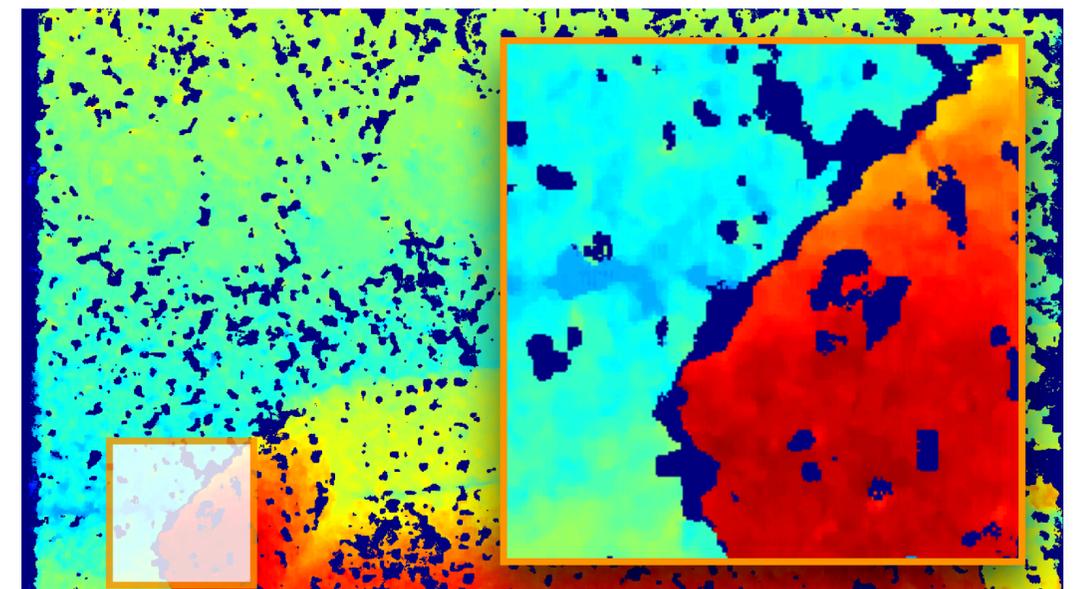
IR Left Input



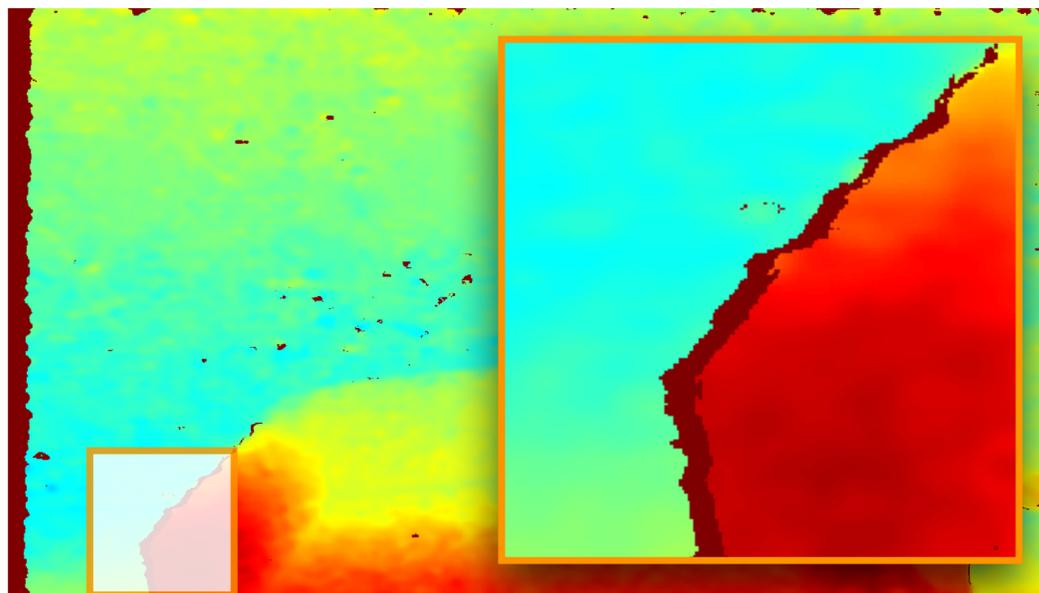
PatchMatch Stereo



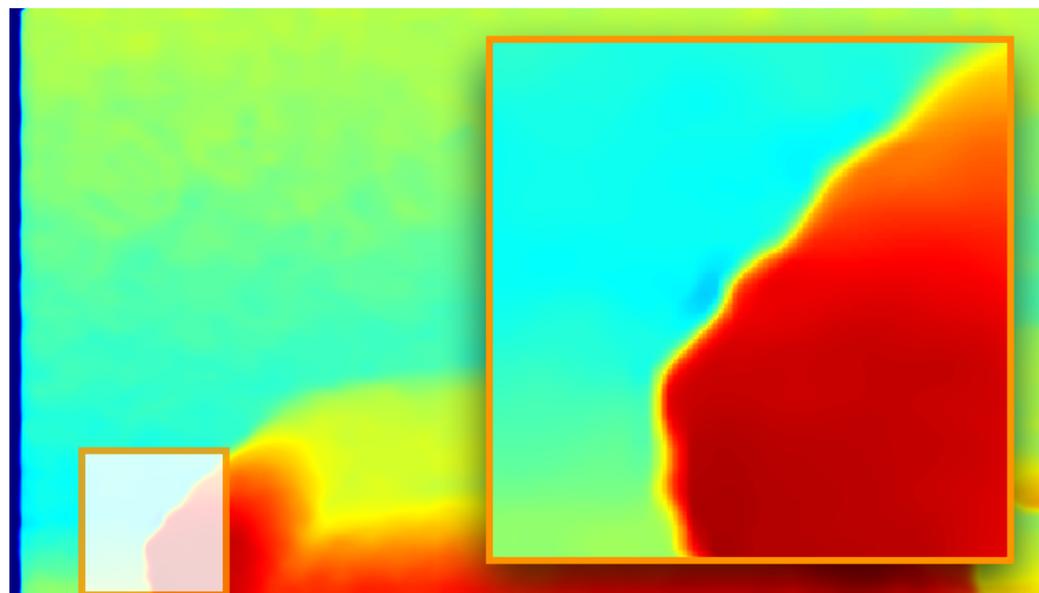
HashMatch Stereo



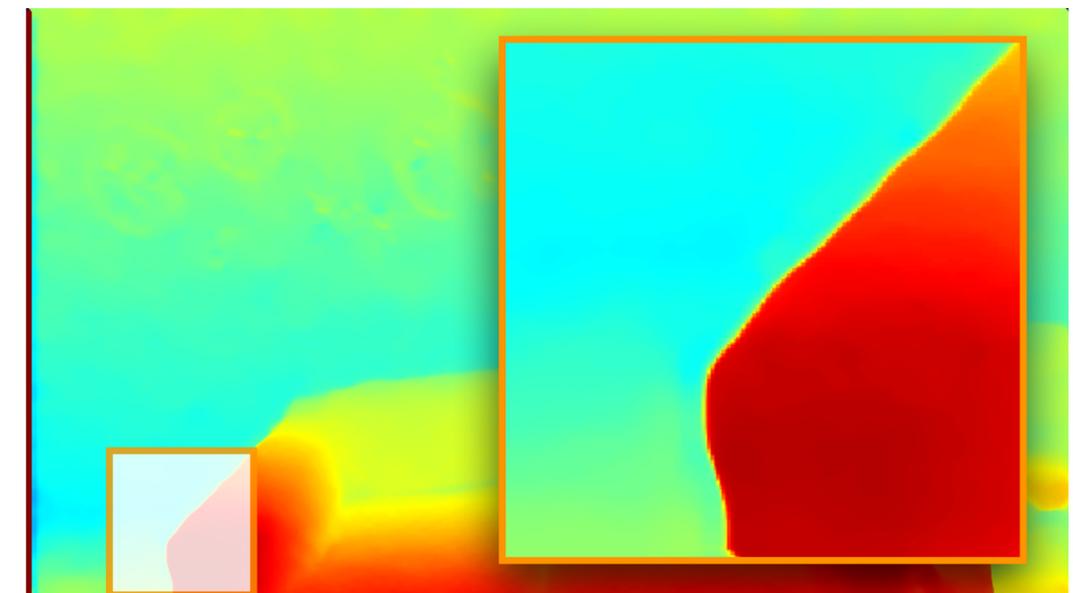
Sensor Output



ASN Semi Supervised (ours)

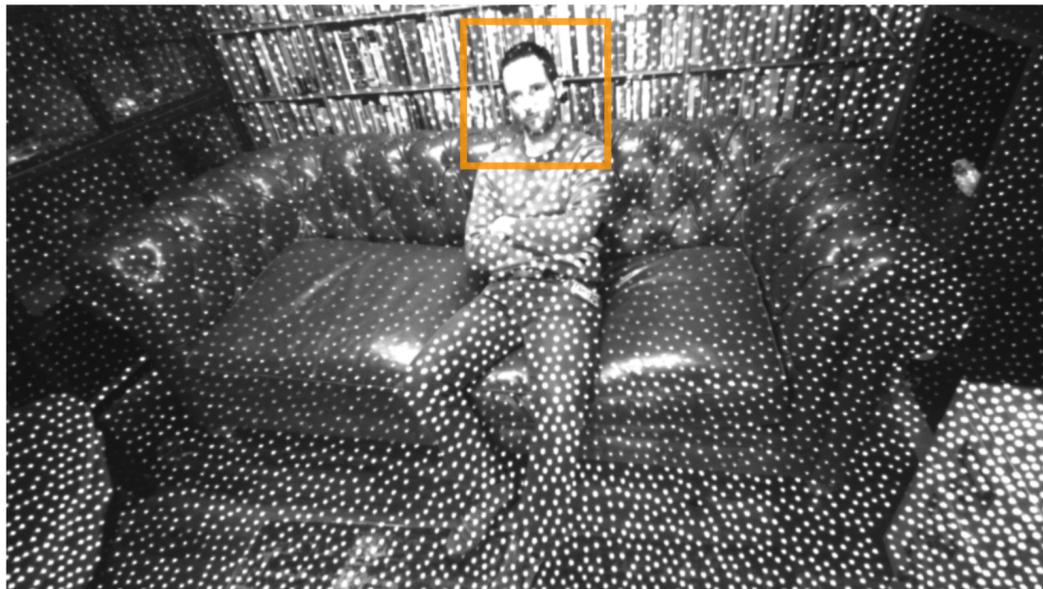


ASN Self-Supervised (ours)

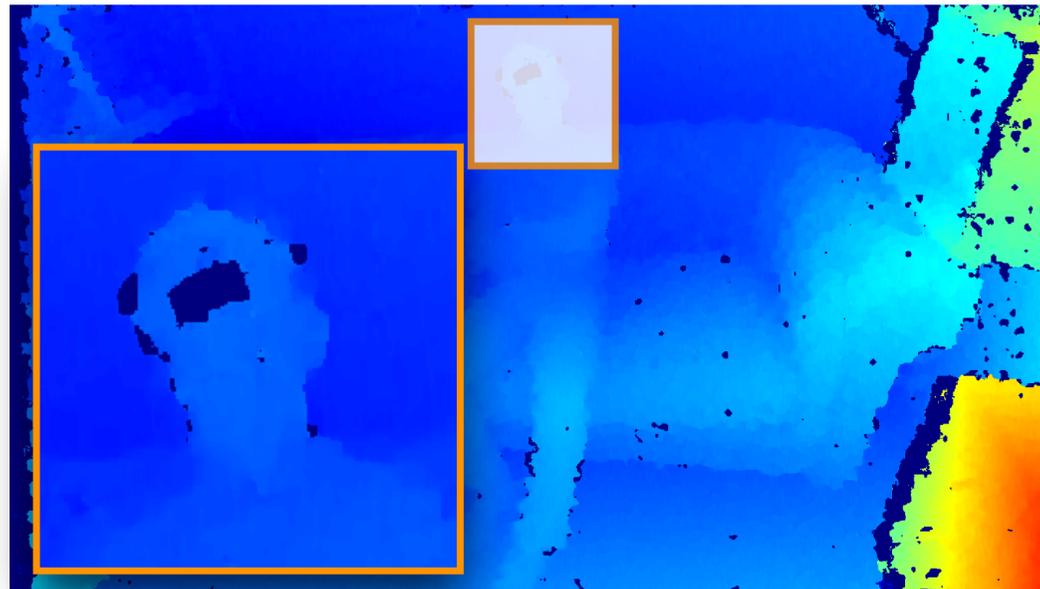


Disparity Quality Result

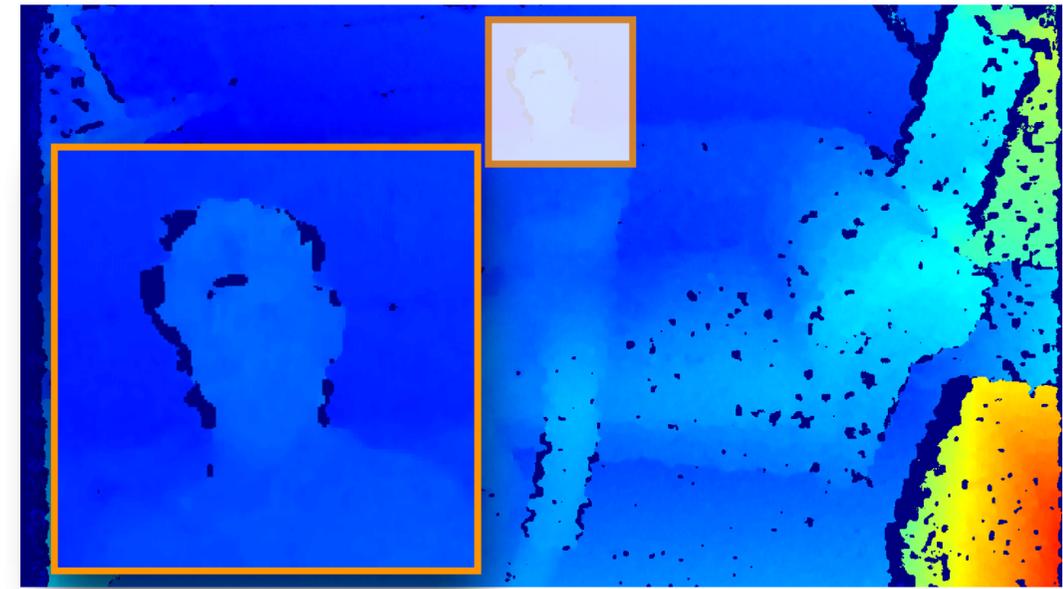
IR Left Input



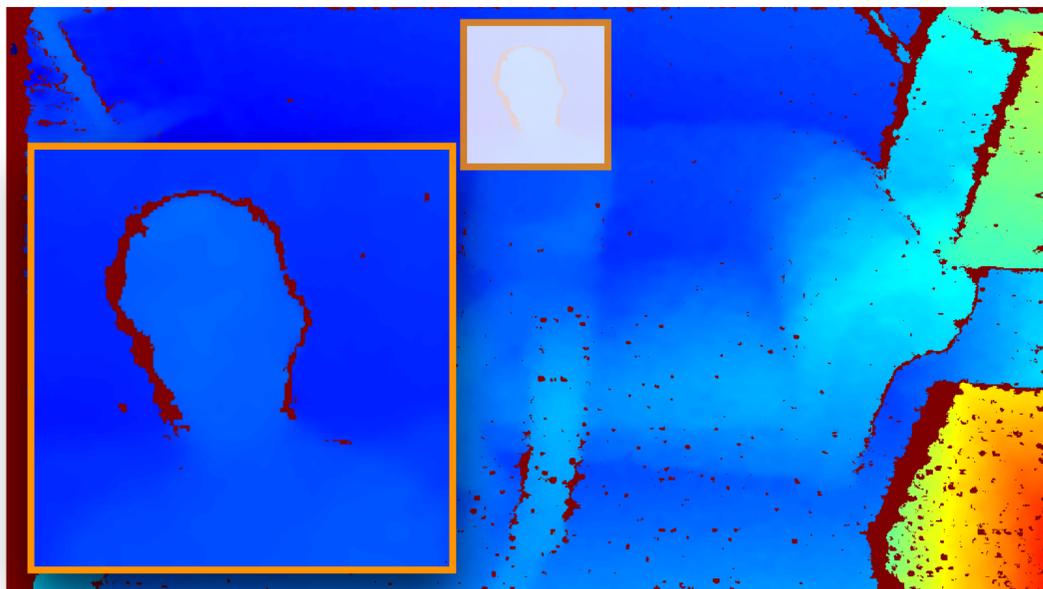
PatchMatch Stereo



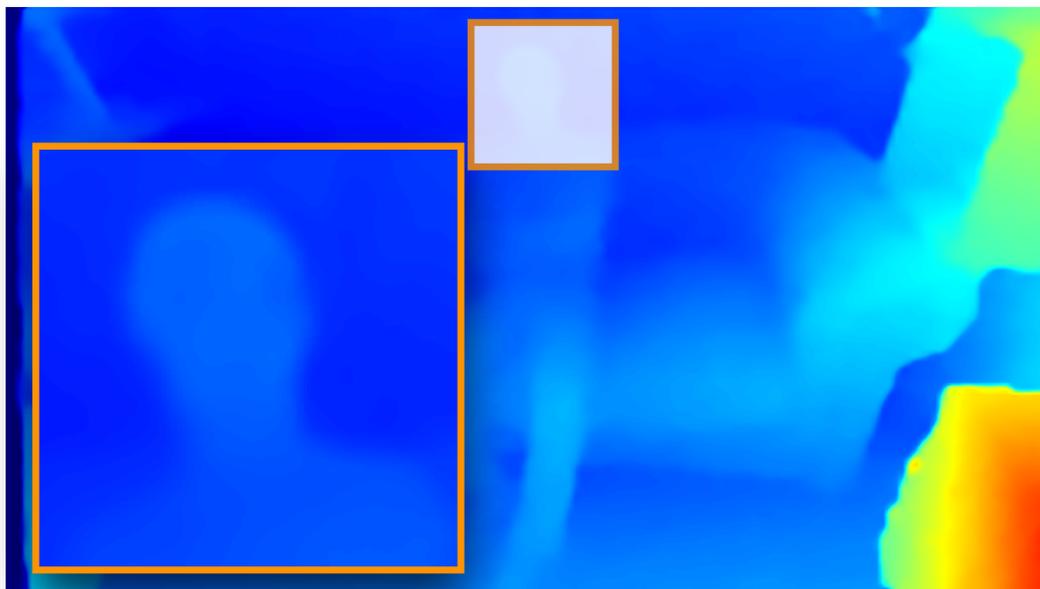
HashMatch Stereo



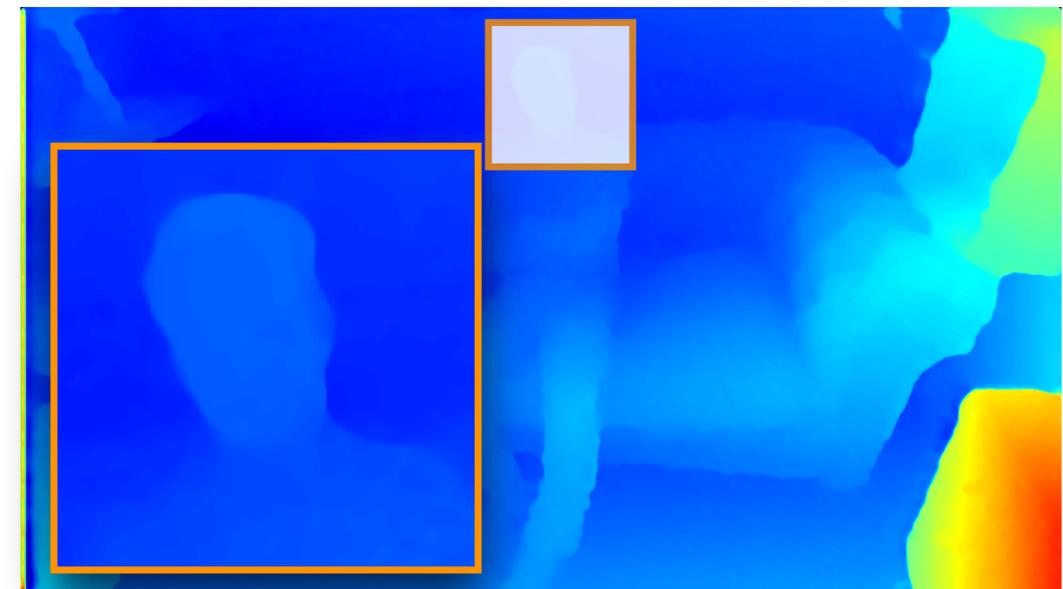
Sensor Output



ASN Semi Supervised (ours)

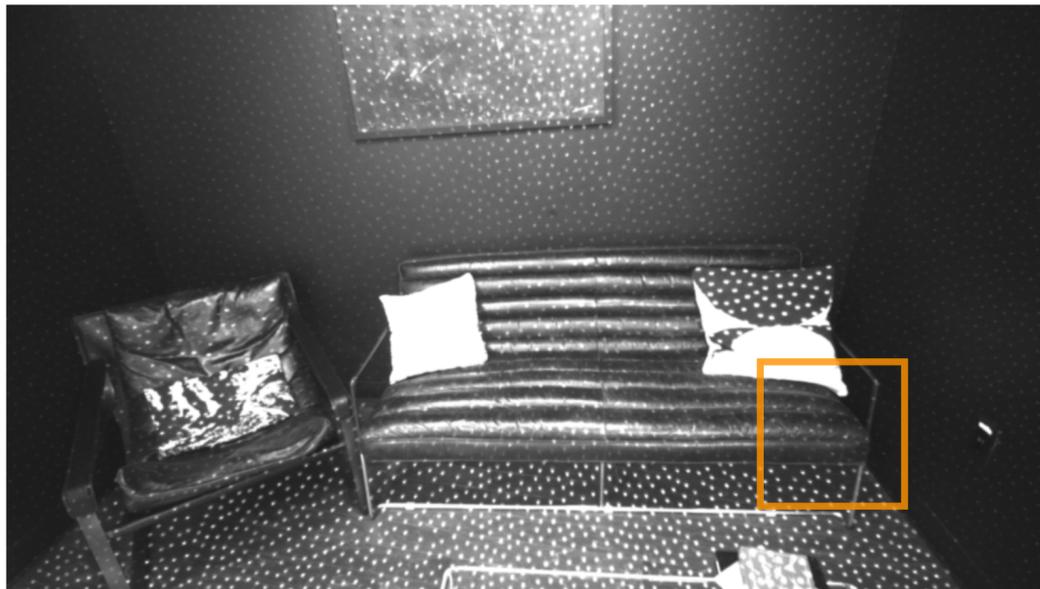


ASN Self-Supervised (ours)

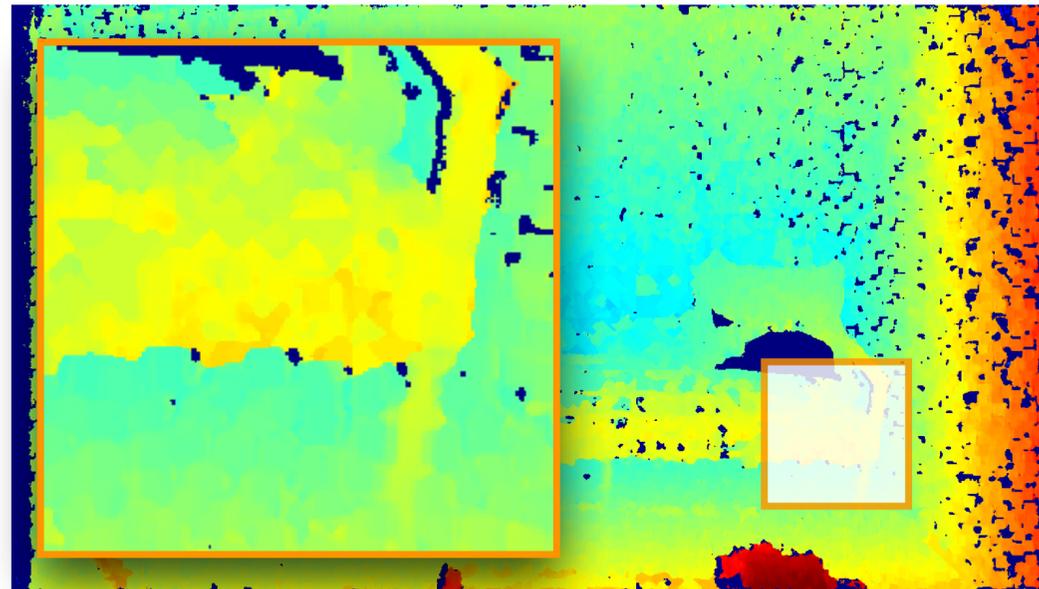


Disparity Qualitative Result

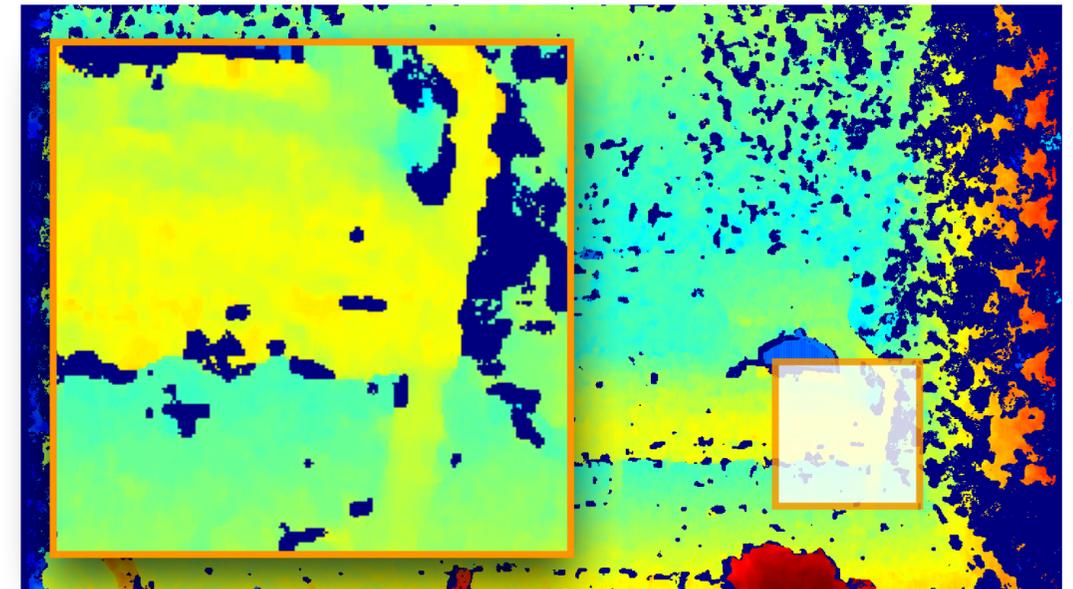
IR Left Input



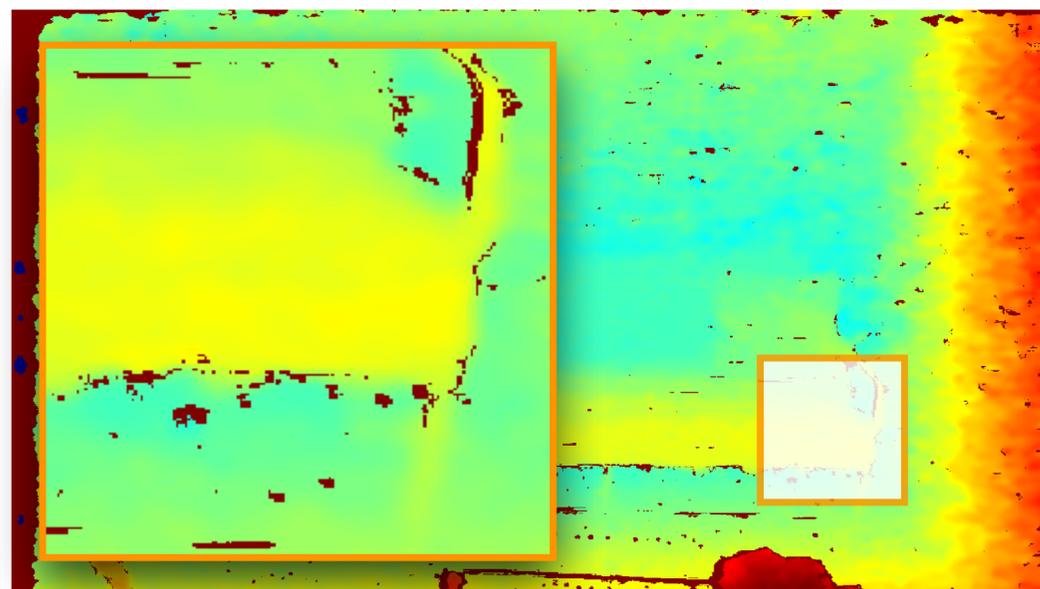
PatchMatch Stereo



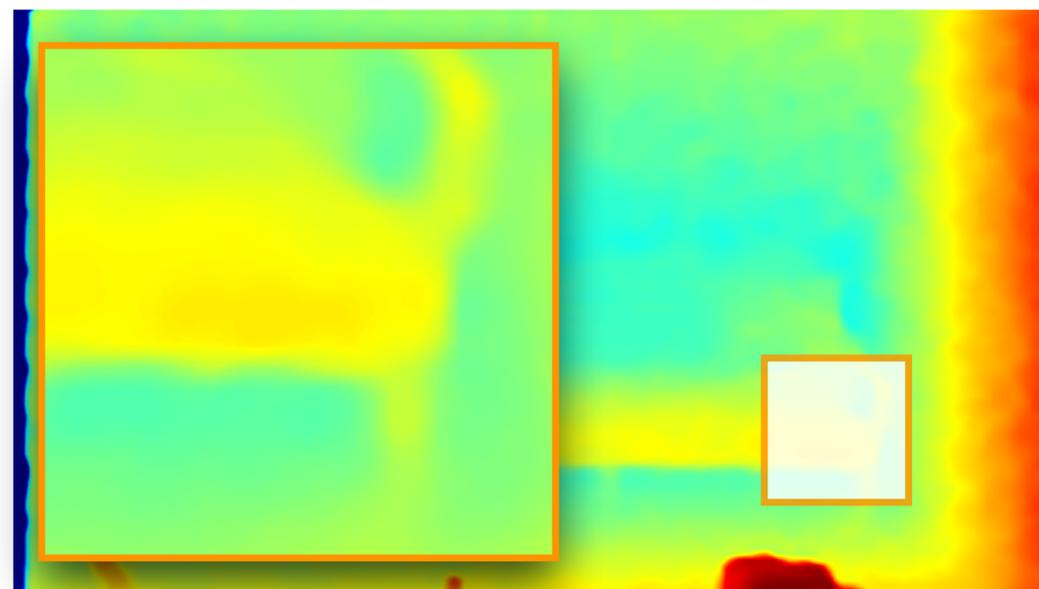
HashMatch Stereo



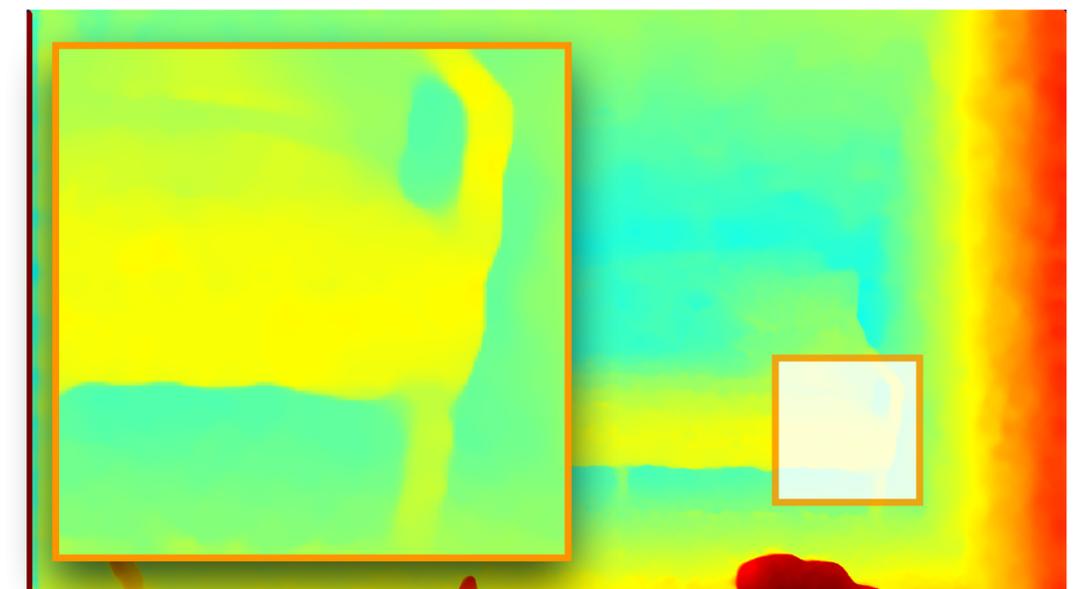
Sensor Output



ASN Semi Supervised (ours)



ASN Self-Supervised (ours)



Disparity Quantitative Result

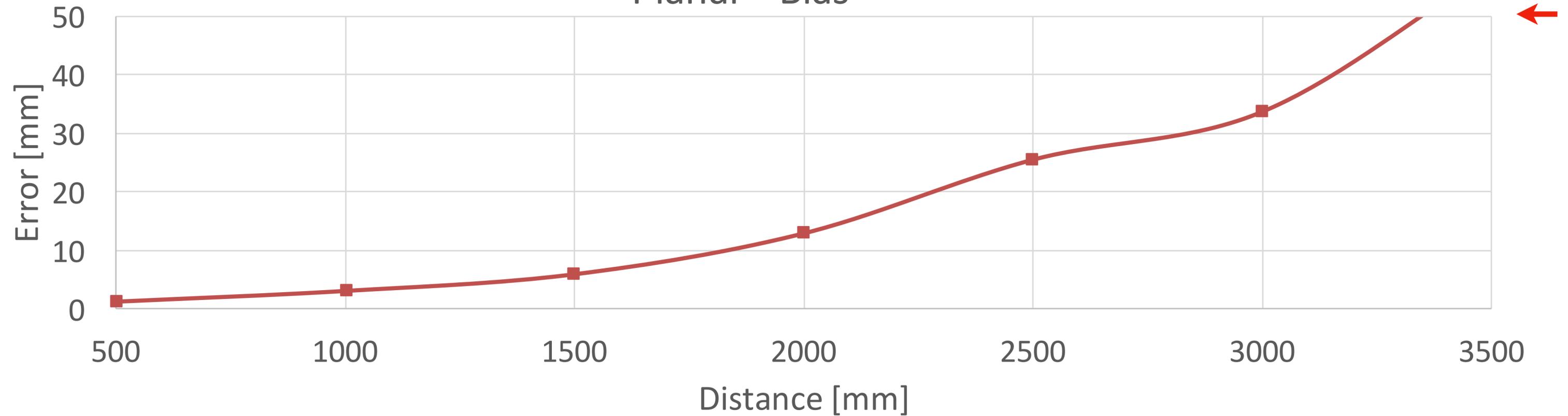


Fit plane on planar-wise scene as ground truth.

Disparity Quantitative Result

Sensor — Semi-Global Matching

Planar - Bias

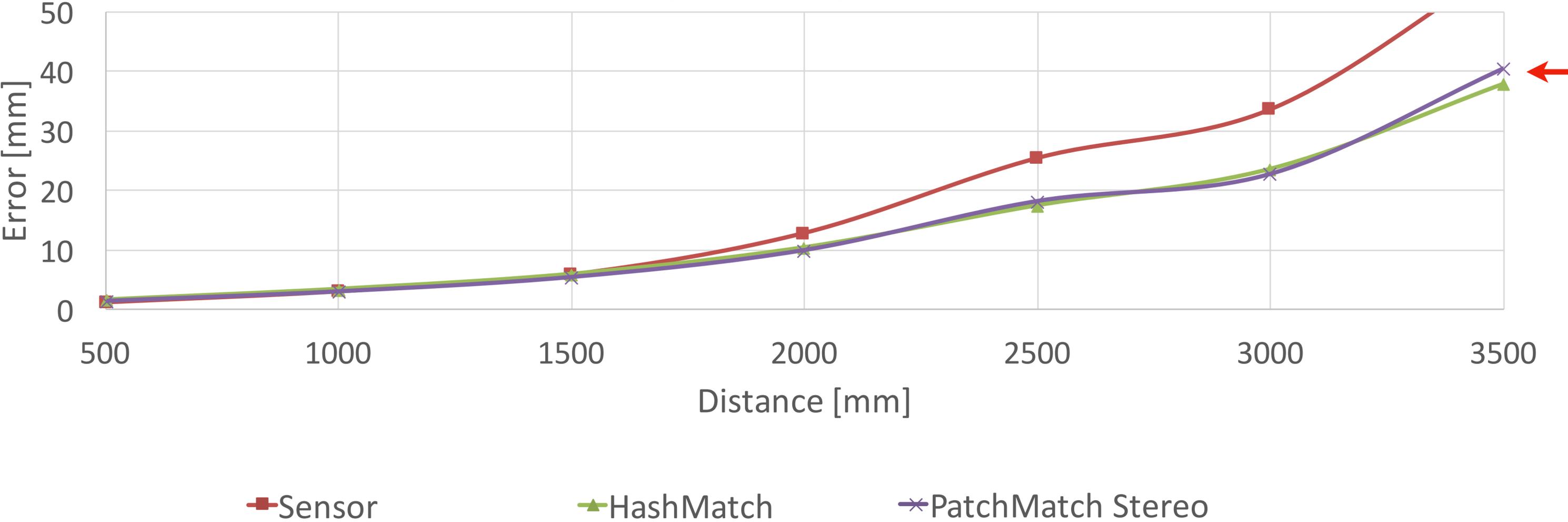


—■— Sensor

Disparity Quantitative Result

Traditional Methods

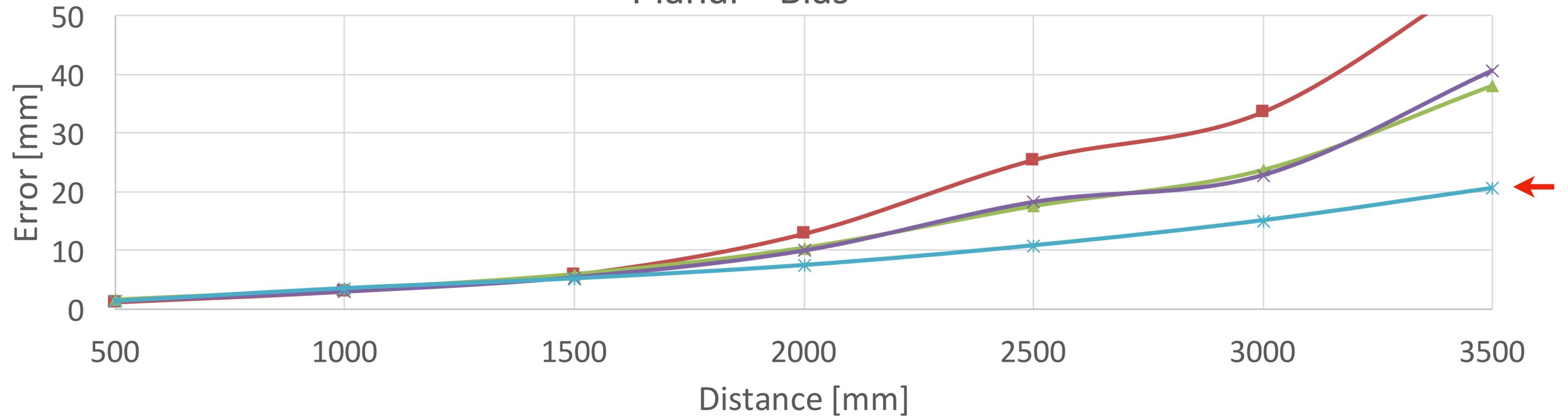
Planar - Bias



Disparity Quantitative Result

Previous Self-Supervised Method

Planar - Bias



■ Sensor

▲ HashMatch

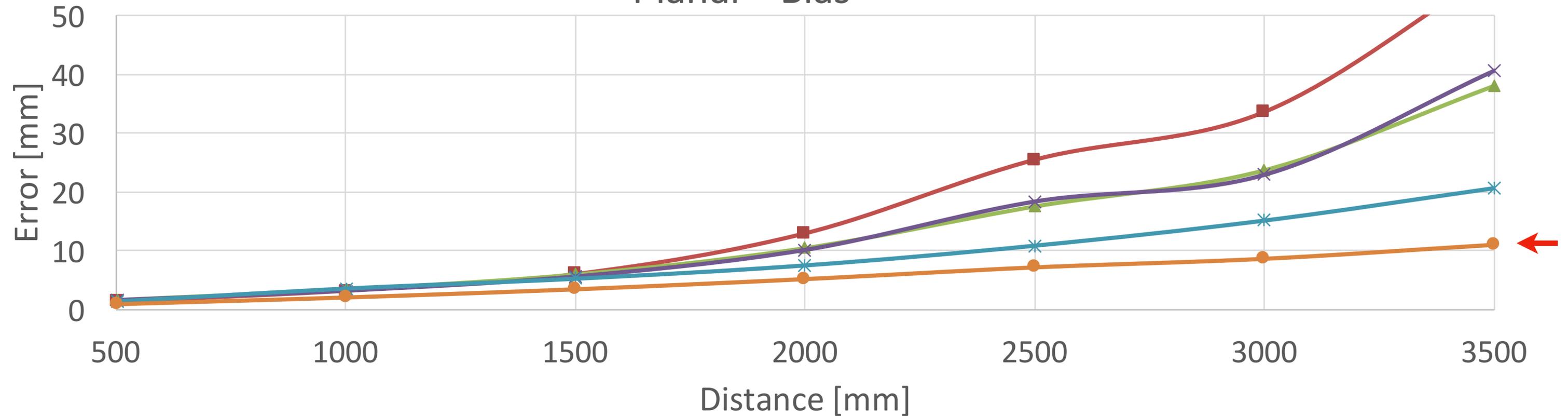
× PatchMatch Stereo

* Godard et al.

Disparity Quantitative Result

ActiveStereoNet

Planar - Bias



- Sensor
- ▲ HashMatch
- ✖ PatchMatch Stereo
- ✱ Godard et al.
- ASN (ours)

Disparity Quantitative Result

ActiveStereoNet

Planar - Bias

50

Disparity Error: 0.2 px \rightarrow 0.03 px

10

0

500

1000

1500

2000

2500

3000

3500

Distance [mm]

■ Sensor

▲ HashMatch

✱ PatchMatch Stereo

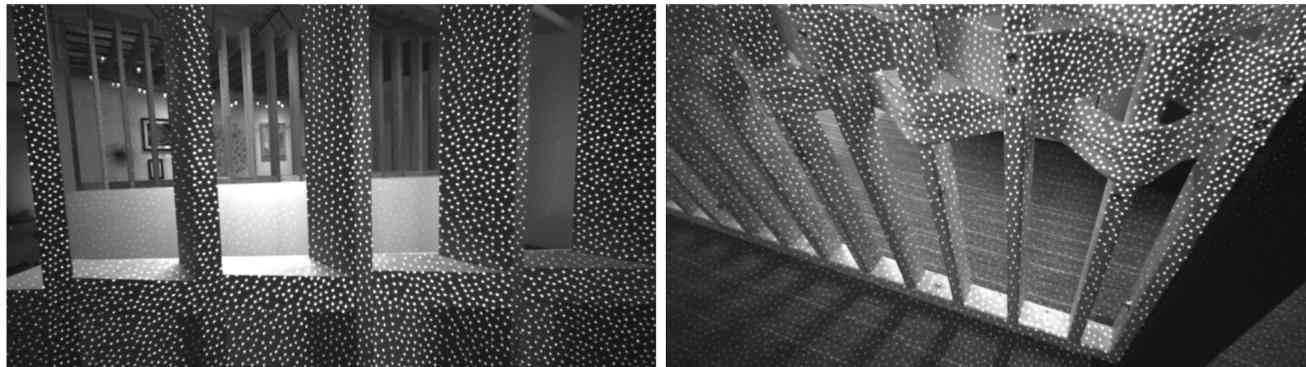
✱ Godard et al.

● ASN (ours)

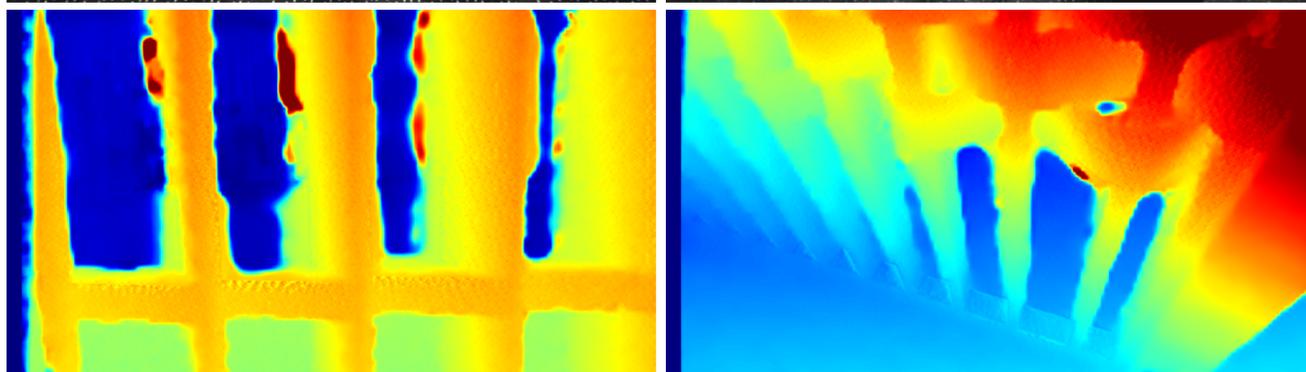
Ablation Study

Does LCN Matter?

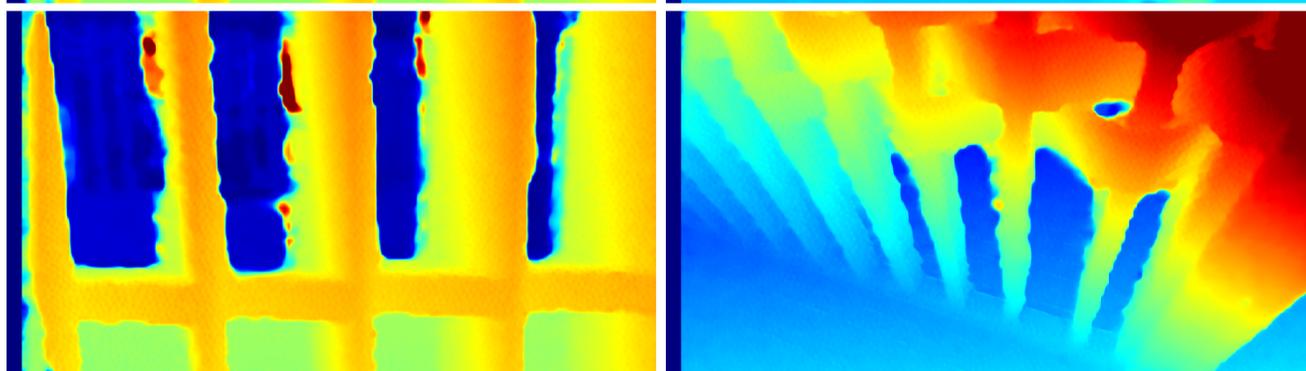
IR Left Input



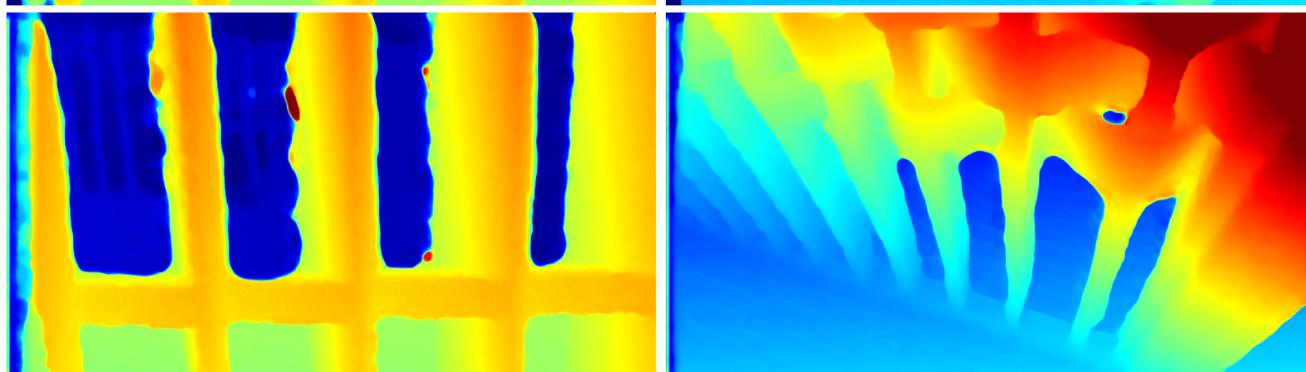
Traditional Photometric Loss



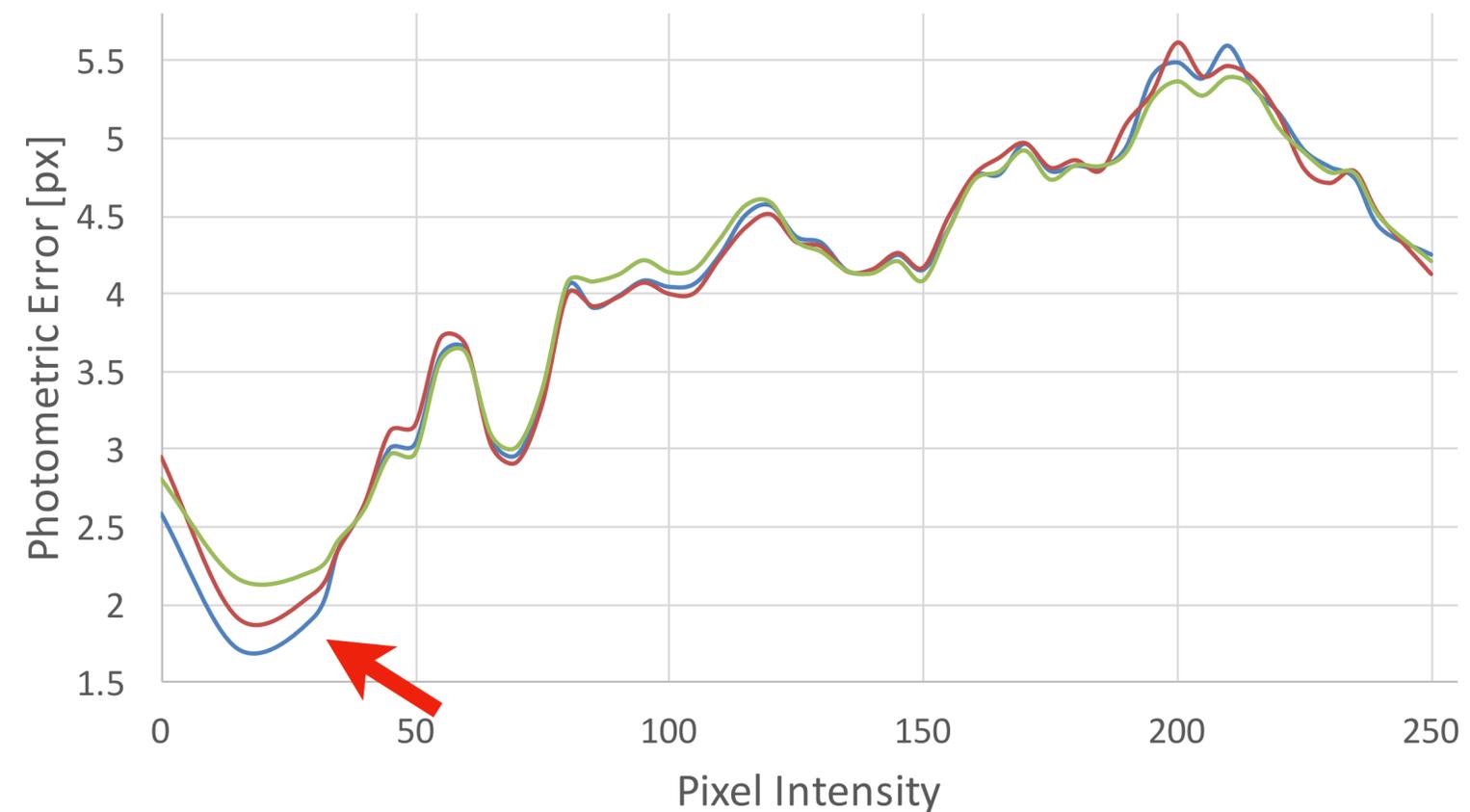
Perceptual Loss



LCN Loss



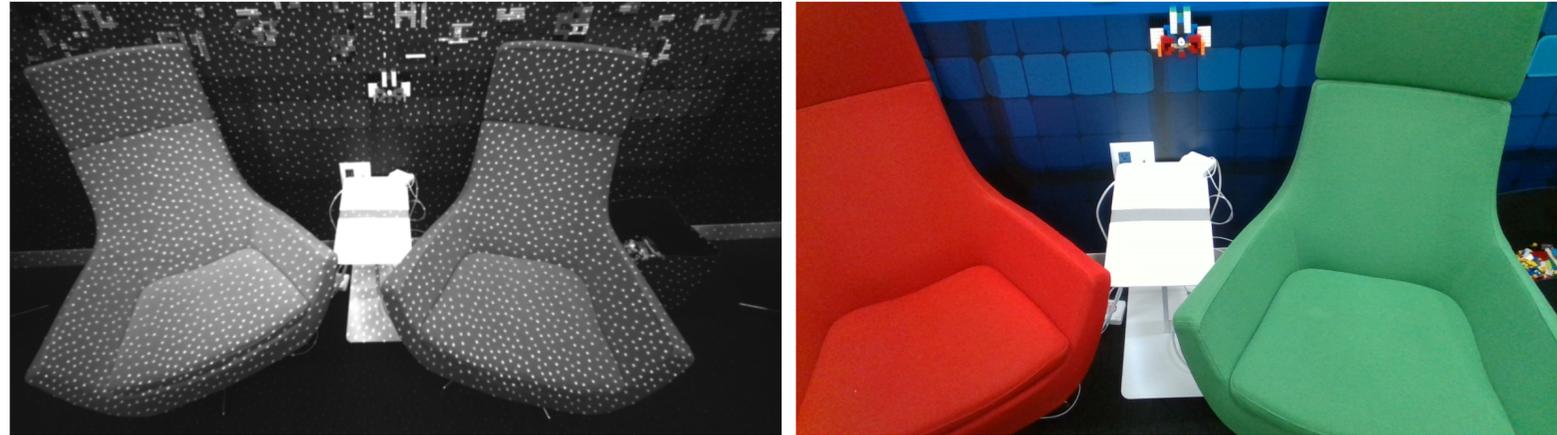
Reconstruction Loss Evaluation



— Proposed Loss — Perceptual Loss — Photometric Loss

LCN: Local Contrast Normalization

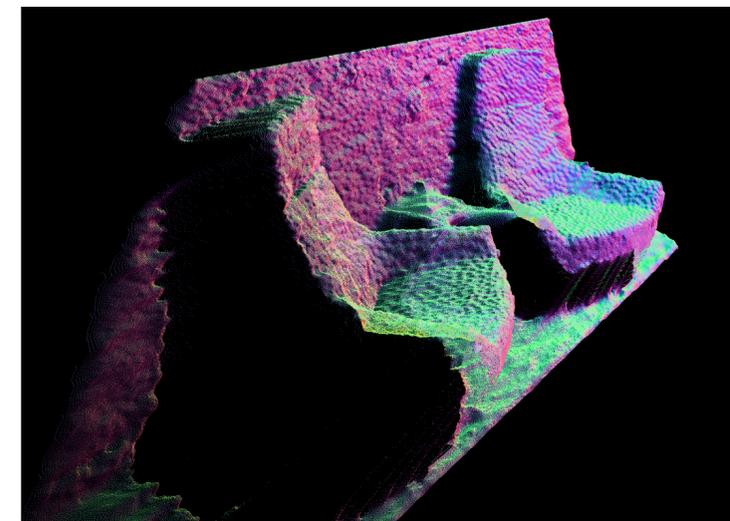
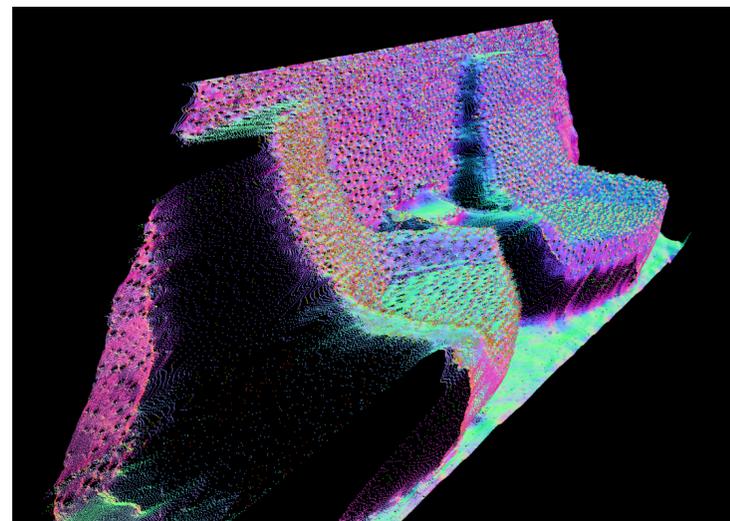
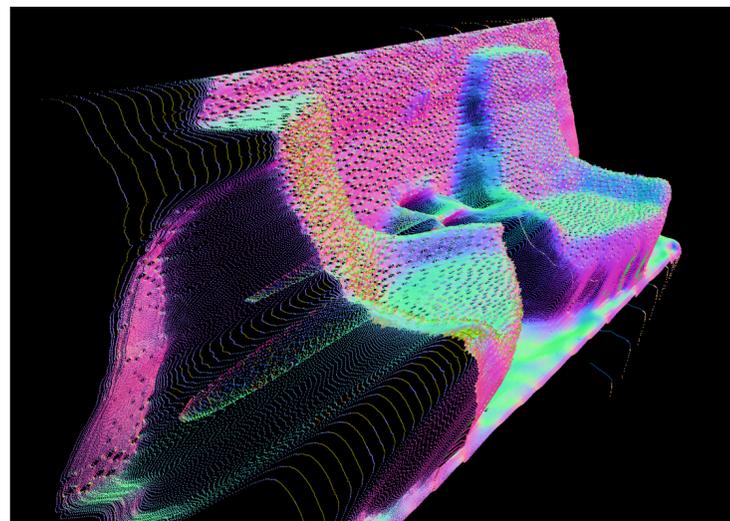
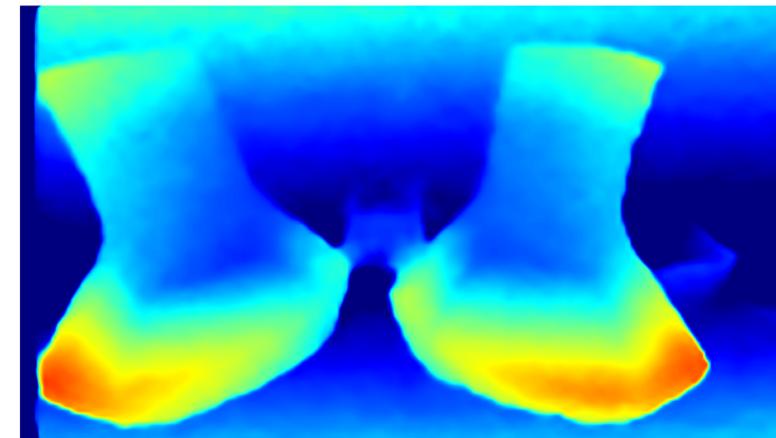
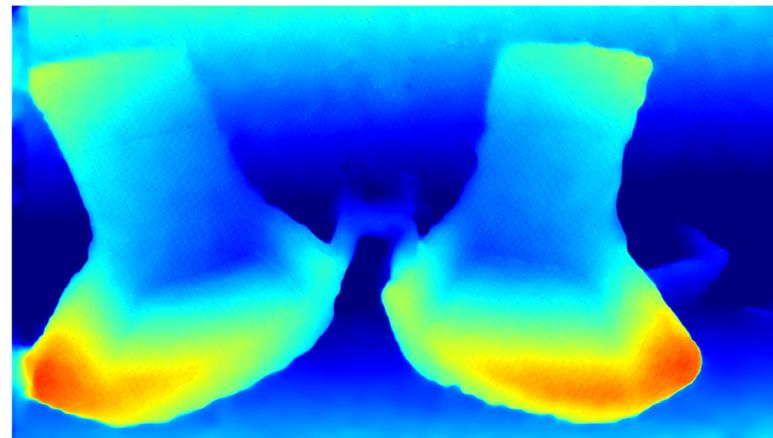
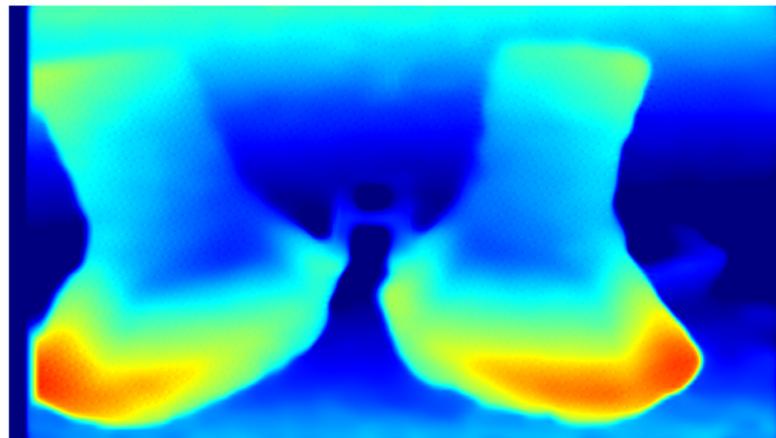
Does Invalidation Matter?



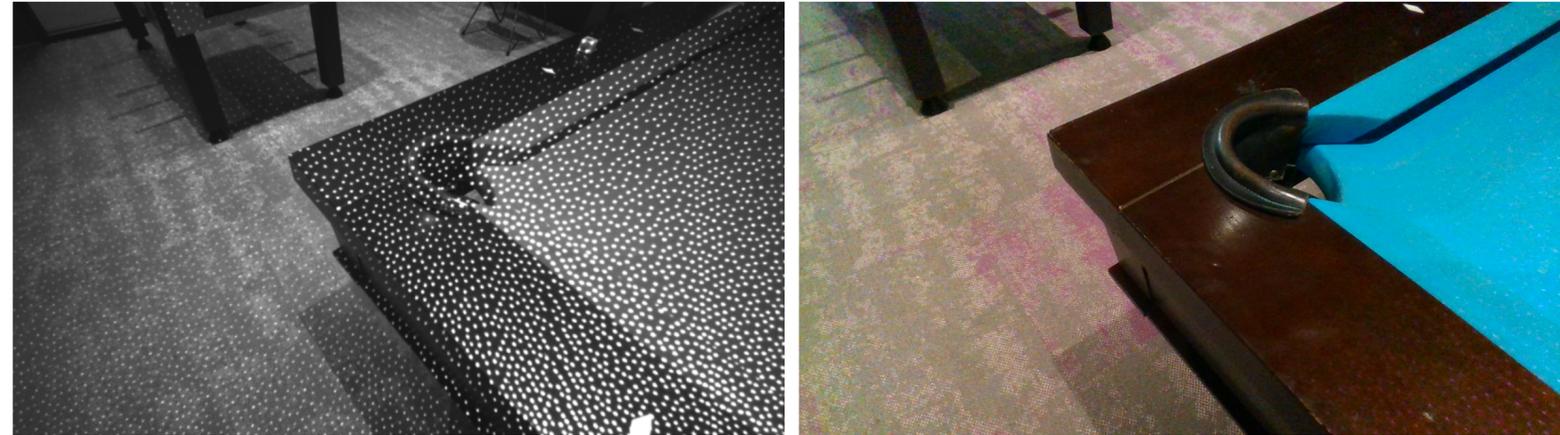
Godard et al.

ASN No Occlusion Mask

ASN Full



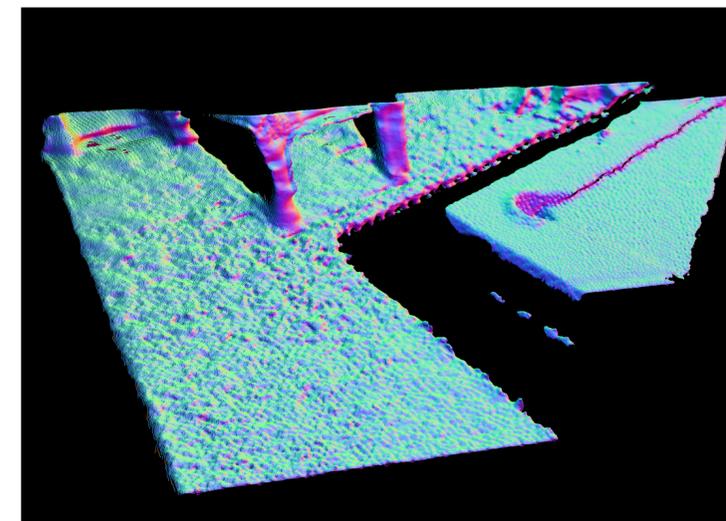
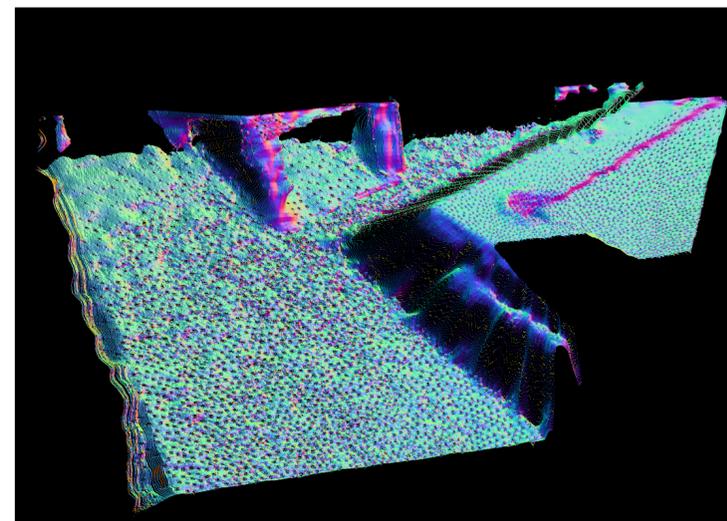
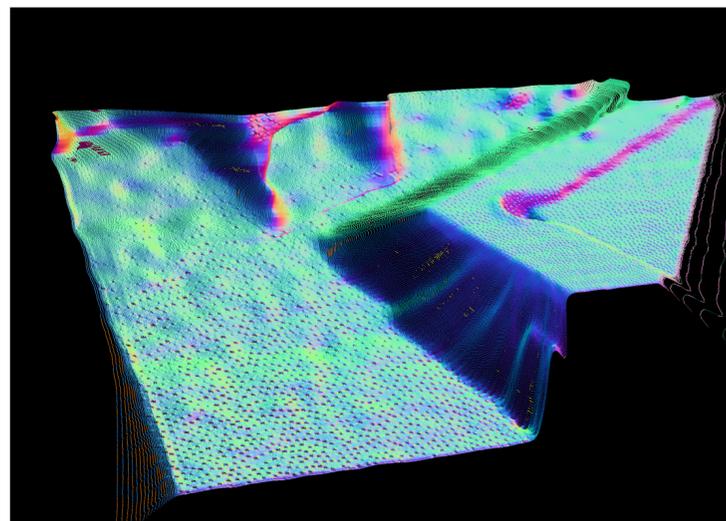
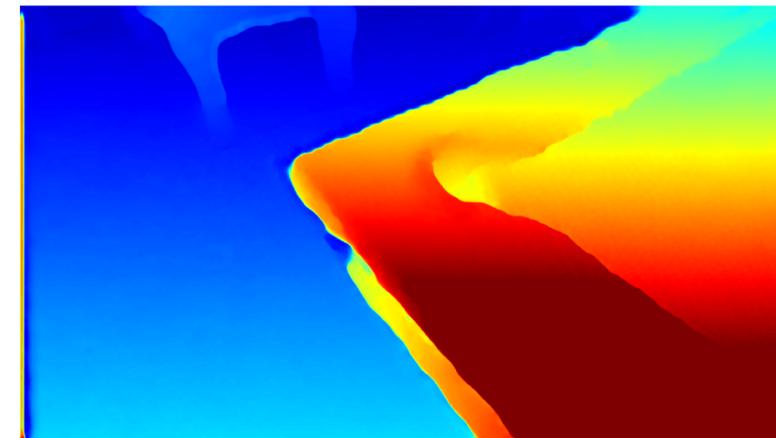
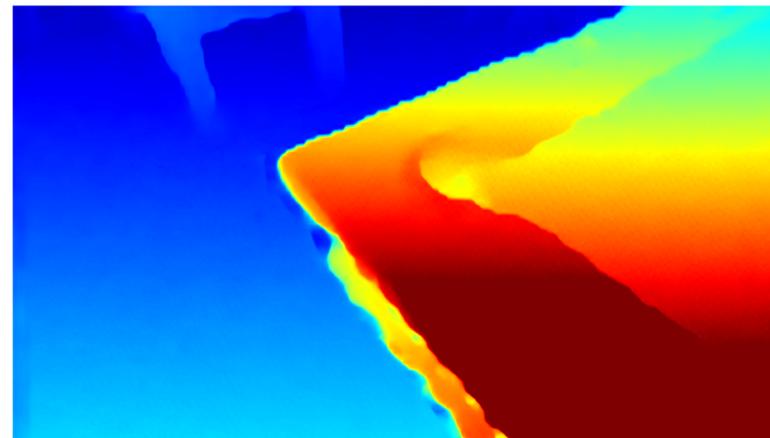
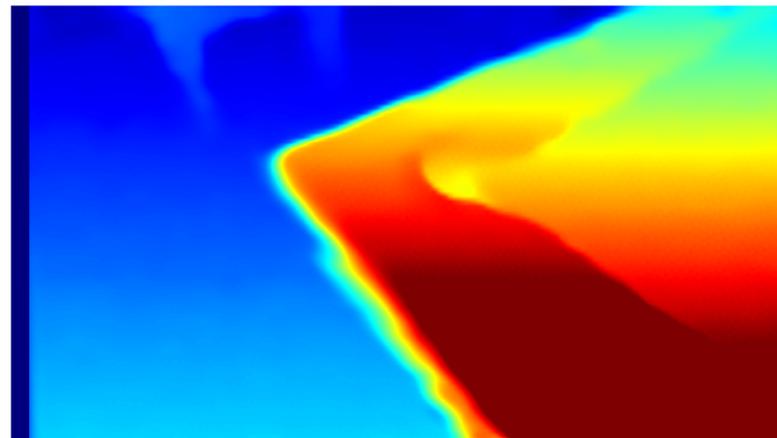
Does Invalidation Matter?



Godard et al.

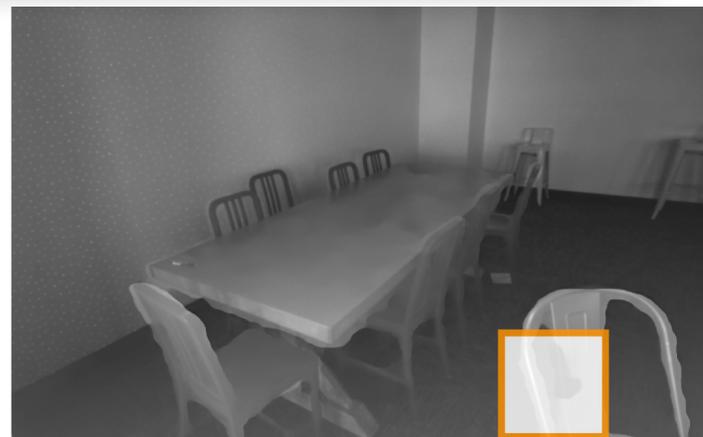
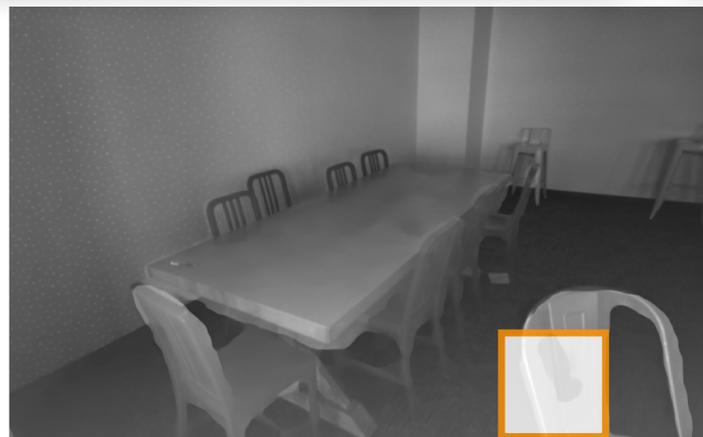
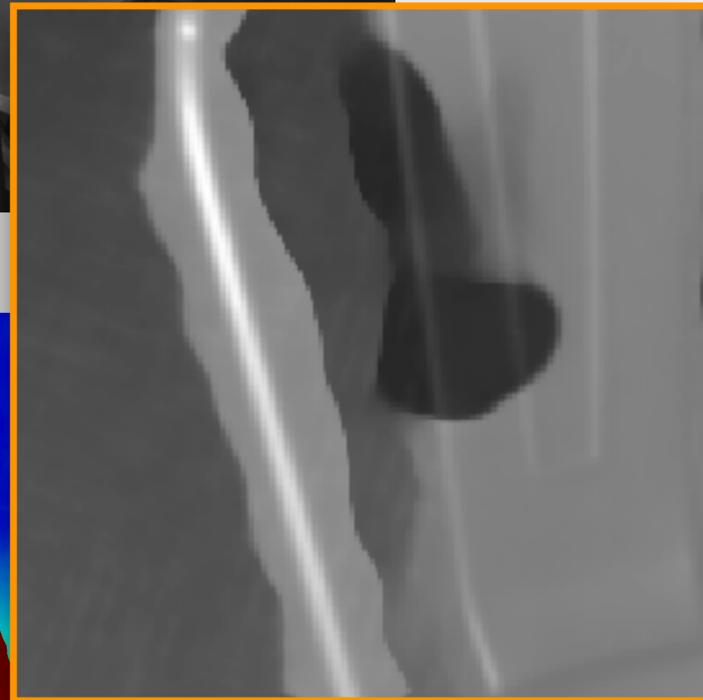
ASN No Occlusion Mask

ASN Full

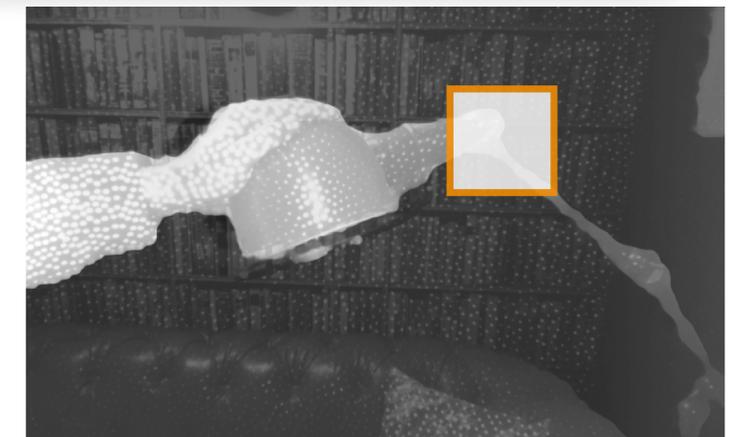
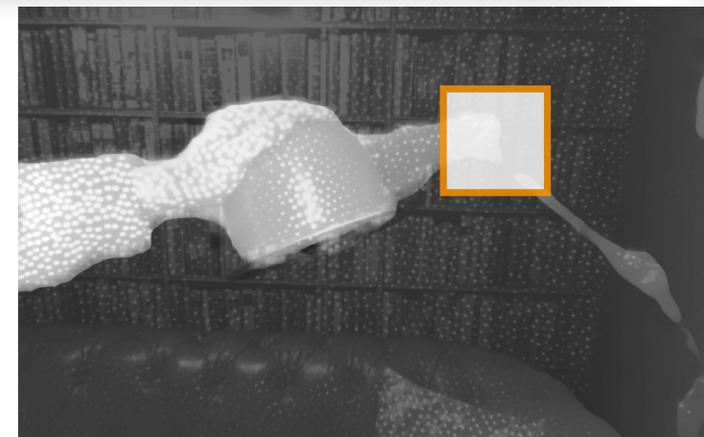


Does Window-based Matter?

Left IR Input



Left IR Input



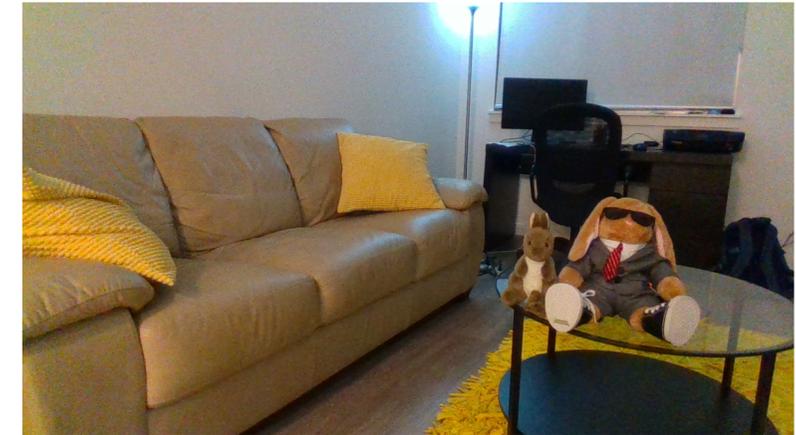
Why Active+Passive?

Only Passive

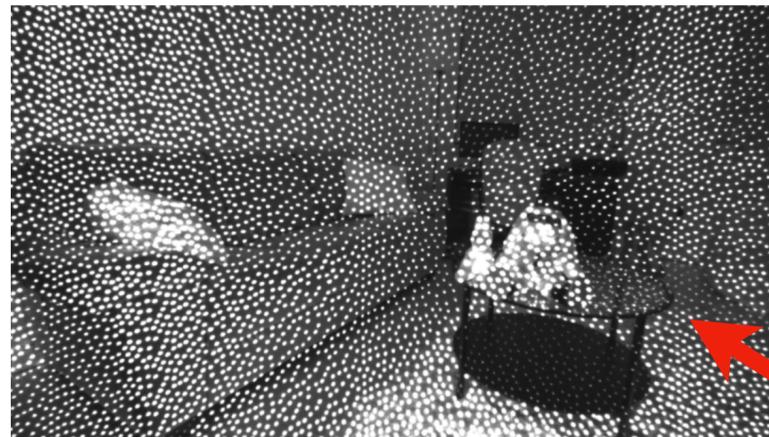
Only Active

Active+Passive

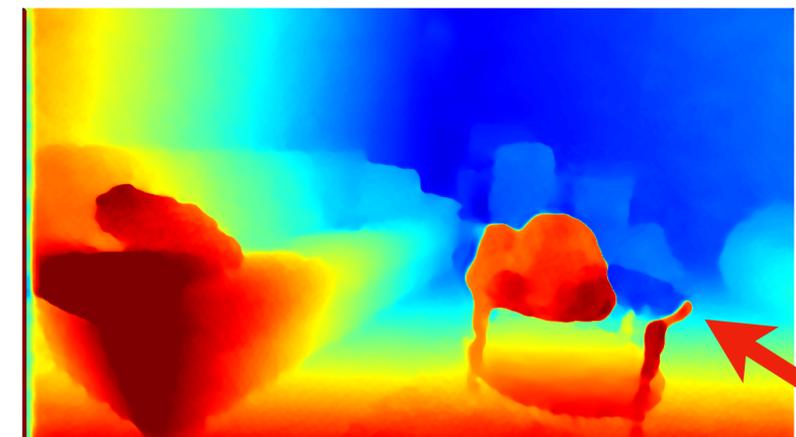
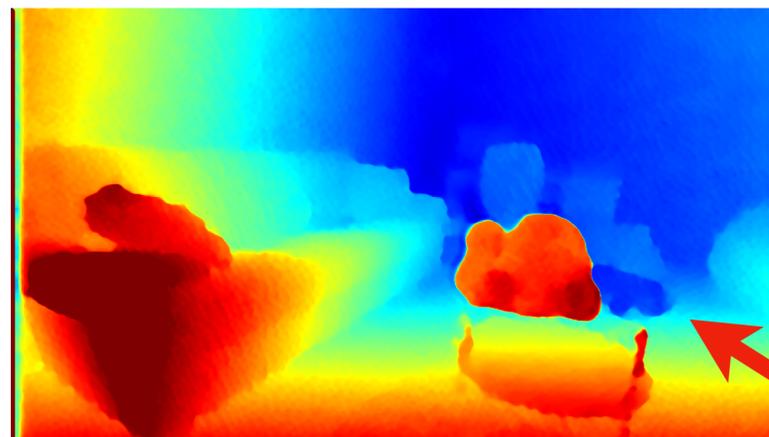
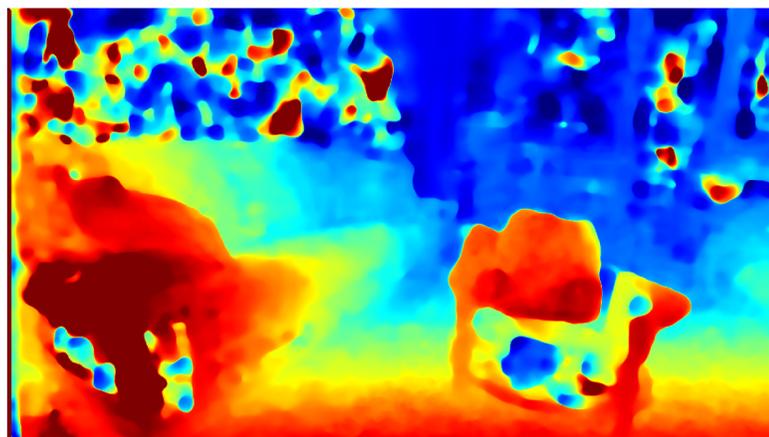
Color Image



IR Image

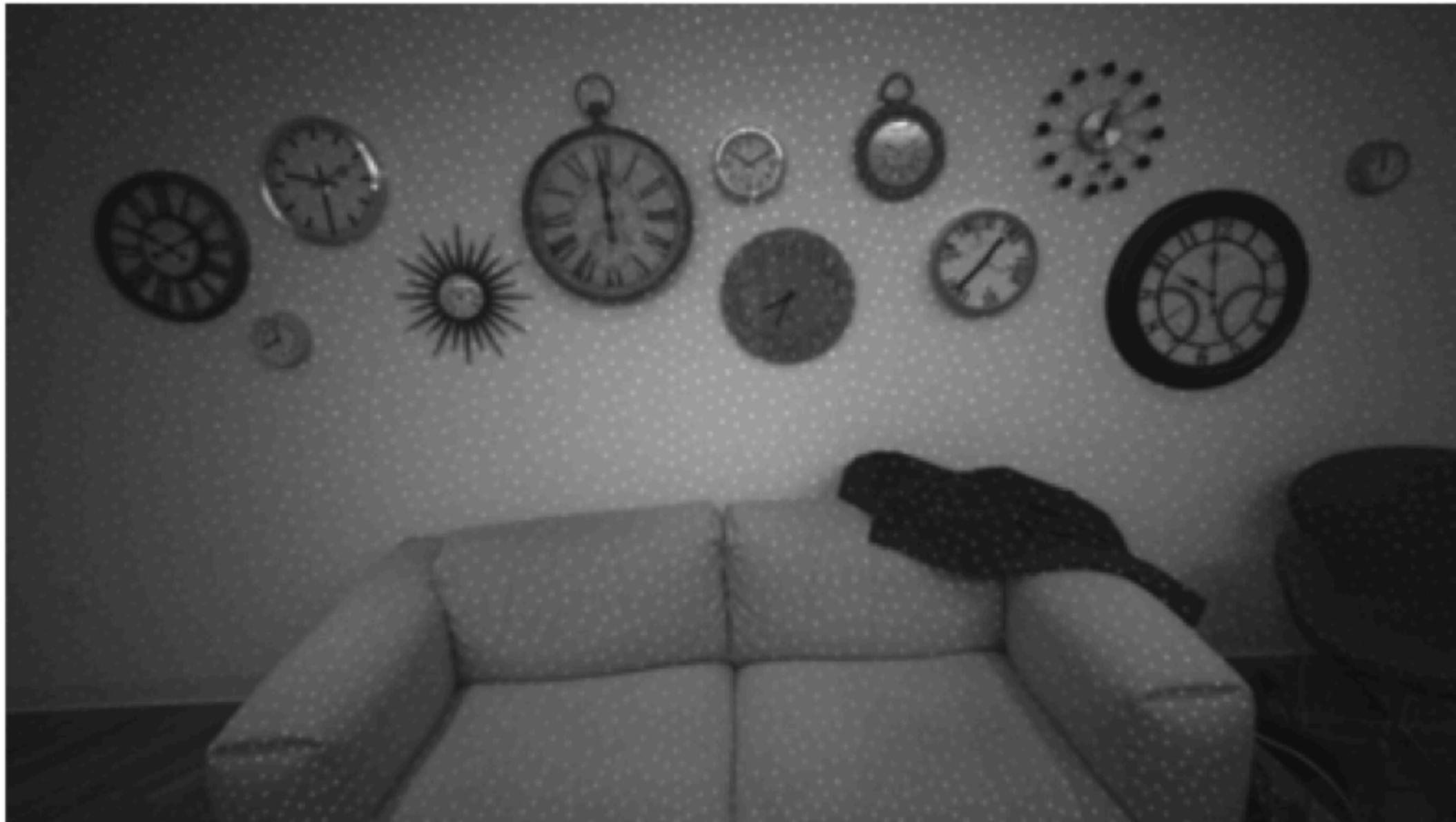


Disparity



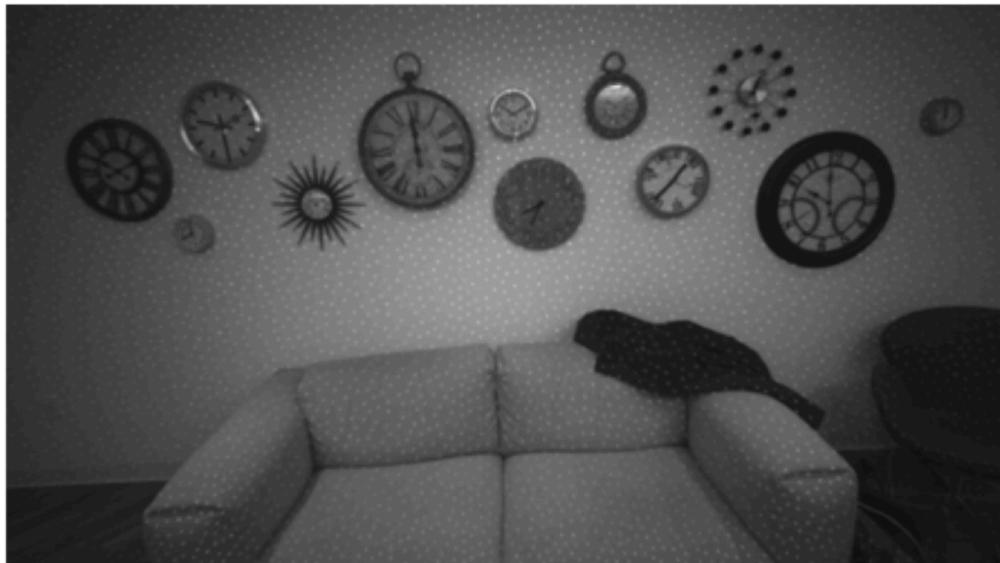
Conclusion

- We investigate active stereo with deep learning.

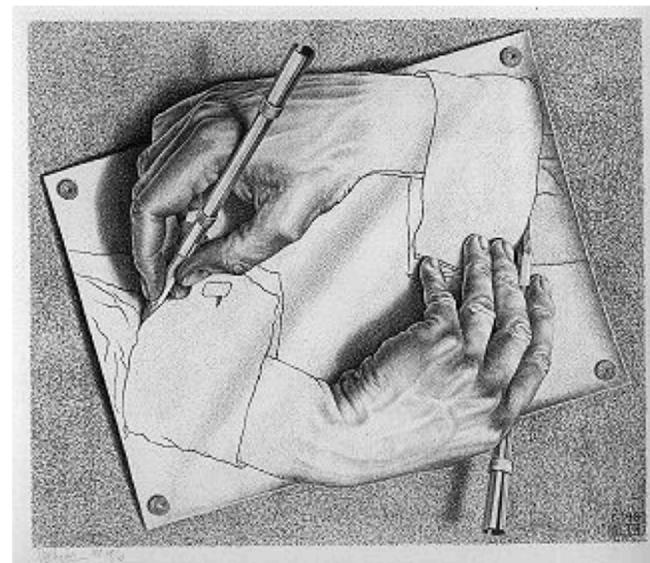


Conclusion

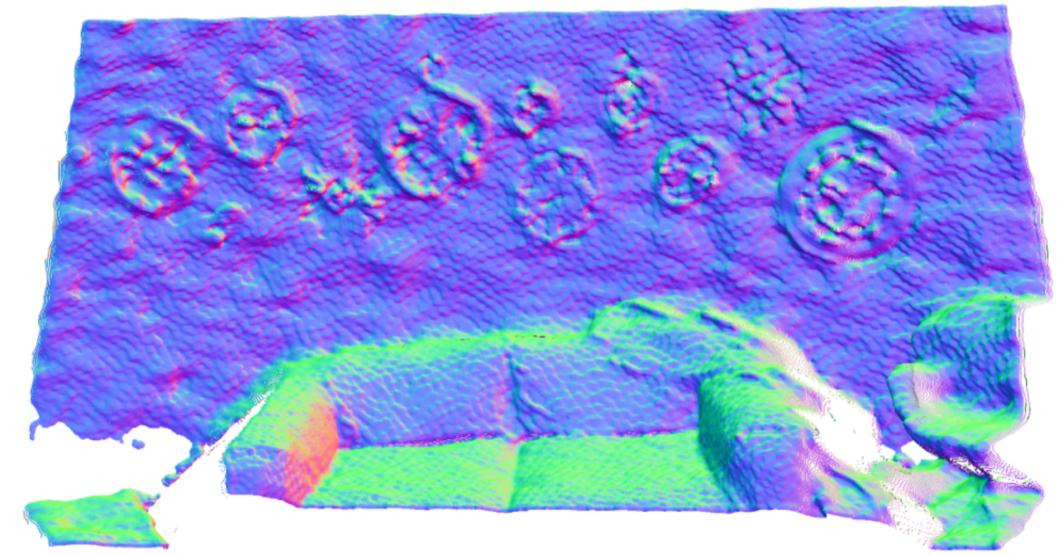
- We investigate active stereo with deep learning.
- We propose self-supervised learning for active stereo system using the improved photometric loss.
- We reduce the disparity error from 0.2 px to 0.03 px.



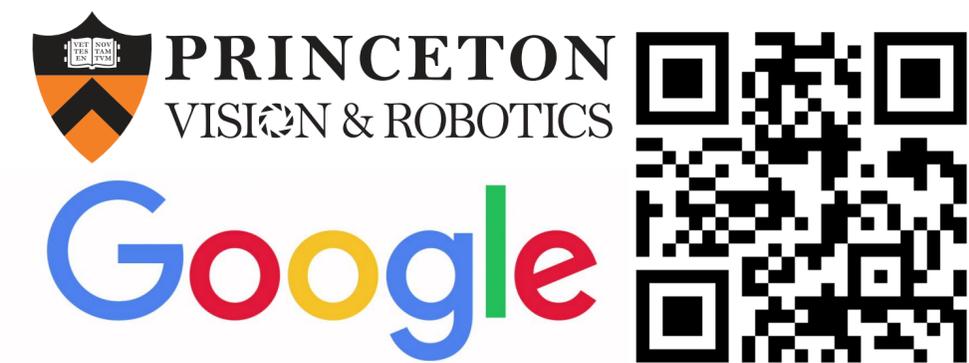
Active Stereo



Self-supervised Learning



Disparity Error: 0.03px



ActiveStereoNet

End-to-End Self-Supervised Learning for Active Stereo Systems

Presenter: Yinda Zhang

Yinda Zhang^{1,2}, Sameh Khamis¹, Christoph Rhemann¹, Julien Valentin¹, Adarsh Kowdle¹, Vladimir Tankovich¹, Michael Schoenberg¹, Shahram Izadi¹, Thomas Funkhouser^{1,2}, Sean Fanello¹

Google Inc.¹ Princeton University²

Project Webpage: <http://asn.cs.princeton.edu/>