



清華大學



**NUS**  
National University  
of Singapore



# A Two-Stage Masked Autoencoder Based Network for Indoor Depth Completion

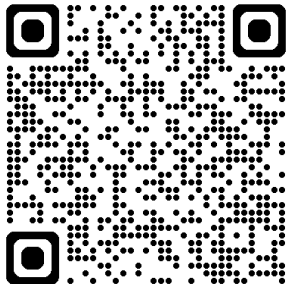
---

Kailai Sun, Zhou Yang, Qianchuan Zhao

Email: [skl23@nus.edu.sg](mailto:skl23@nus.edu.sg)

2024.6.18

Project:



WeChat:



# Opening

- **Scan to BIM**: a workflow or process that translates scanned, point-cloud digital models into building information modeling (BIM) platforms
- **Indoor 3D reconstruction** [1]: create a 3D digital spatial information representation of the interior of a building

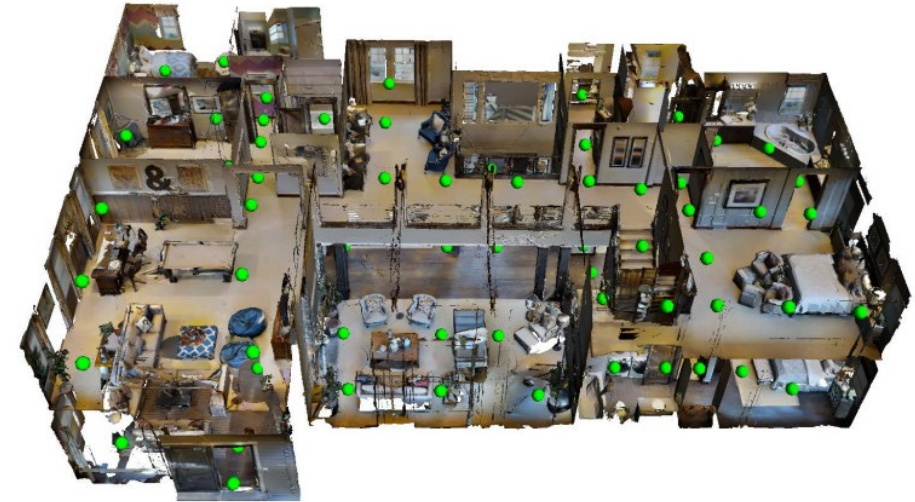


A point cloud is displayed in Autodesk ReCap

[1] Shayan Nikooohemat, et al. Indoor 3D reconstruction from point clouds for optimal routing in complex buildings to support disaster management. Automation in Construction, 113, 2020.

# Challenge

- **Depth completion**: an **important** task focuses on using part of the depth data measured in the real scene to obtain more **dense and complete** depth data.
- Cause: illumination or the materials of the scene objects, limited distance
- However, the latest methods often suffer from sensitivity to **dynamic environmental lighting** conditions.

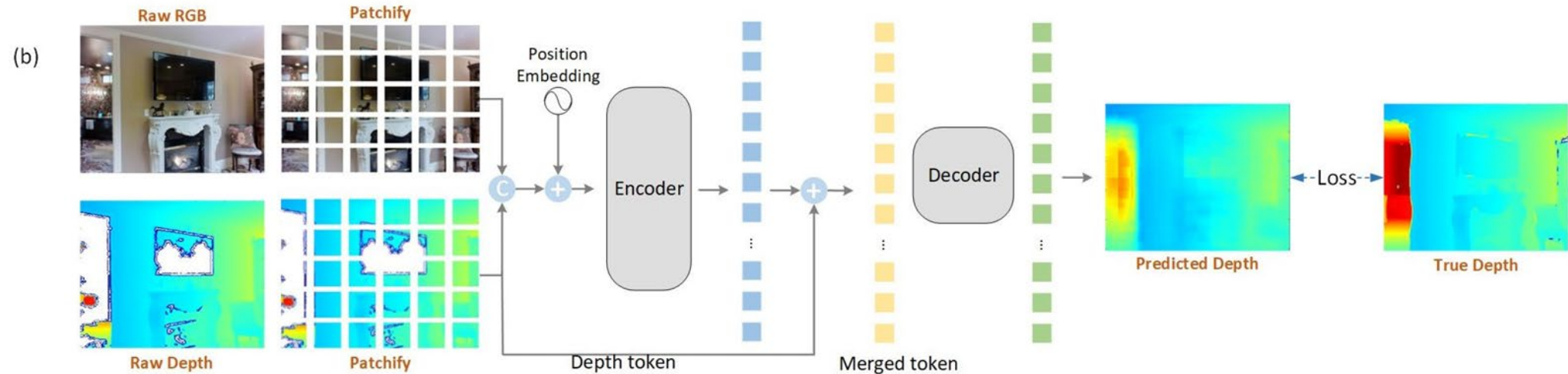
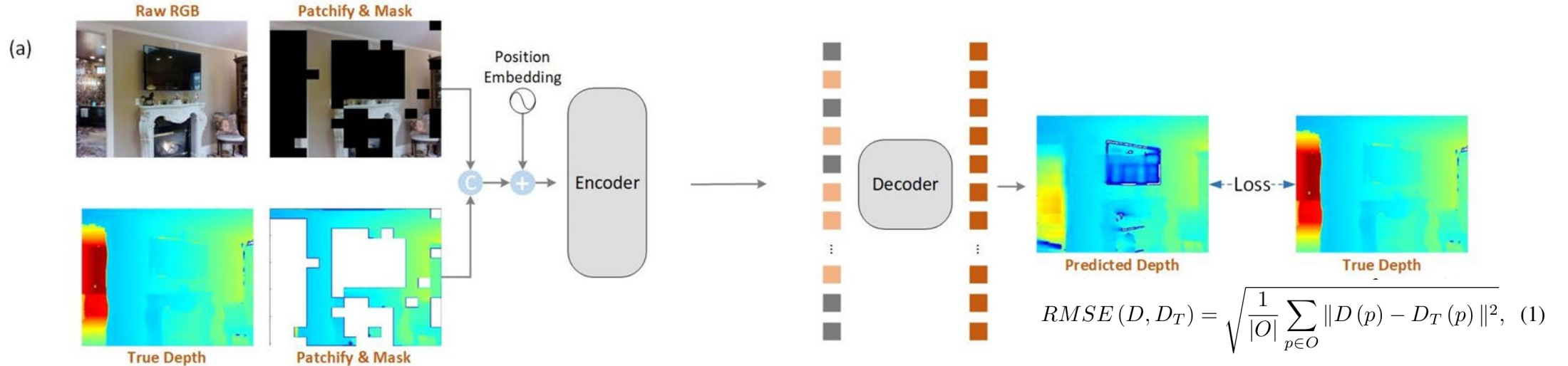


Matterport3D dataset

**15%** depth values are missing

- Masked Autoencoder only apply partial observation to reconstruct the entire image, learning **robust** features and improving the generalization ability.
- We consider: Missing depth patches <sup>simulate</sup> ← ? Masks.
- We propose a Vision Transformer-based two-stage network for indoor depth completion:
  - (1) an MAE-based self-supervision pre-training **encoder** to learn an effective latent representation from the jointly masked RGB and depth images;
  - (2) a **decoder** based on token fusion to complete (reconstruct) the full depth from an incomplete depth image.

# Method



# Result

Methods	RMSE↓	ME↓	SSIM↑	$\delta_{1.25} \uparrow$	$\delta_{1.25^2} \uparrow$
Joint Bilateral Filter	1.978	0.774	0.507	0.613	0.689
MRF[1]	1.675	0.618	0.692	0.651	0.780
AD[2]	1.653	0.610	0.696	0.663	0.792
FCN	1.262	0.517	0.605	0.681	0.808
Zhang[3]	1.316	0.461	0.762	0.781	0.851
Huang[4]	1.092	0.342	<b>0.799</b>	0.850	0.911
Struct-MDC[5]	1.060	0.503	0.534	0.656	0.713
Pre-training	1.216	0.675	0.642	0.705	0.800
Fine-tuning w/o Pre-training	<b>0.660</b>	0.243	0.654	0.794	0.904
Fine-tuning w/ Pre-training	<b>0.690</b>	<b>0.206</b>	<b>0.765</b>	<b>0.852</b>	<b>0.912</b>

$$ME(D, D_T) = \frac{1}{|O|} \sum_{p \in O} \|D(p) - D_T(p)\|$$

$$SSIM(D, D_T) = \frac{(2\mu_{D_T}\mu_D + c_1)(2\sigma_{D_T D} + c_2)}{(\mu_{D_T}^2 + \mu_D^2 + c_1)(\sigma_{D_T}^2 + \sigma_D^2 + c_2)}$$

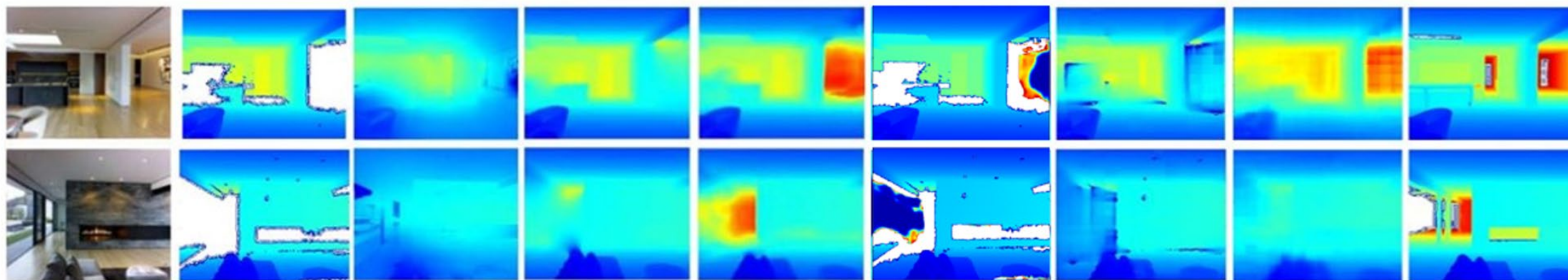
$$Max(\frac{D(p)}{D_T(p)}, \frac{D_T(p)}{D(p)}) < t,$$

The pre-training model is better than the traditional methods (e.g., Joint bilateral filter and MRF) and the method of Zhang on RMSE.

Our fine-tuning model achieves superior performance on the Matterport3D dataset, and performs best on  $\delta$  and ME.



# Result



RGB image

Raw depth

Bilateral

Zhang

Huang

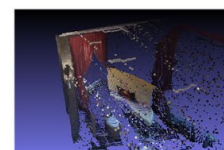
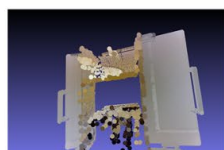
Struct-MDC

Pre-training  
(Ours)

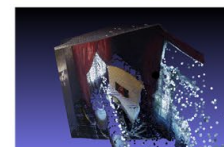
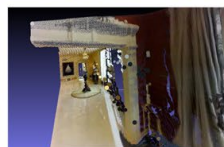
Fine-tuning  
(Ours)

True depth

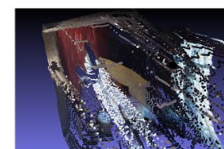
Raw  
PCD



GT  
PCD



Pre-training  
PCD



Fine-tuning  
PCD



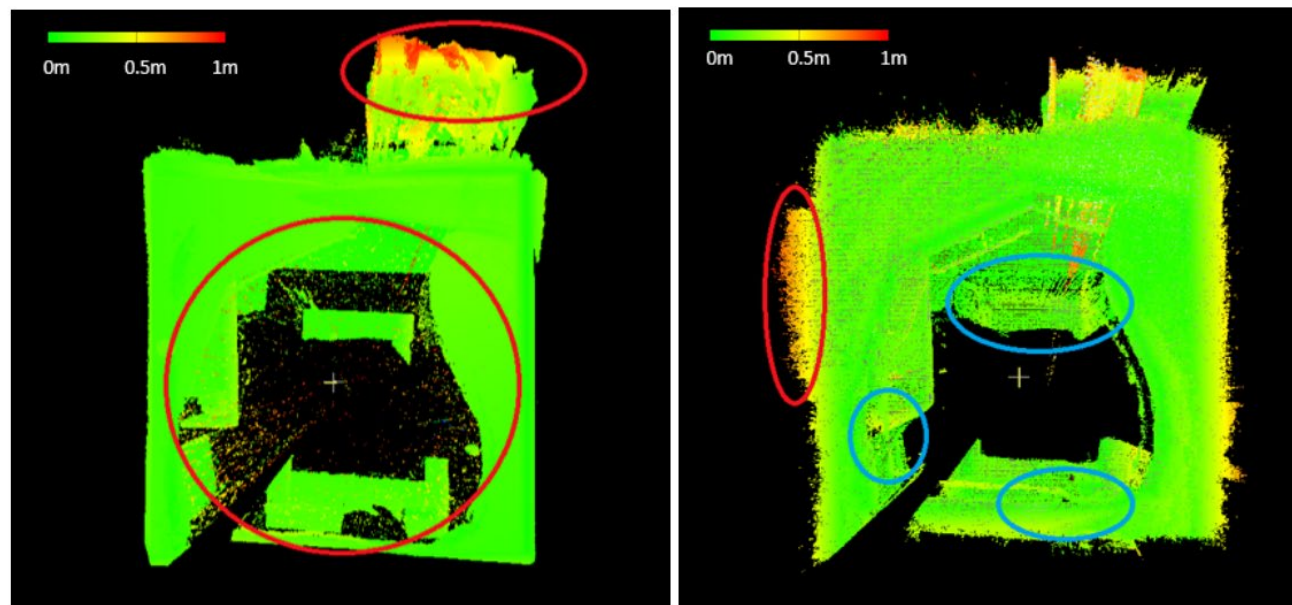
# Application: Indoor 3D Reconstruction



ICL-NUIM dataset

Depth completion

ORB-SLAM3



Reconstruction errors before/after depth completion

Methods	Mean (m)↓	Median (m)↓	Standard Deviation (m)↓	Minimum (m)	Maximum (m)↓
Depth incompletion	0.138	0.053	0.200	0.0	1.106
Depth completion	<b>0.086</b>	0.057	<b>0.101</b>	<b>0.0</b>	<b>1.100</b>



清華大學



**NUS**  
National University  
of Singapore



# Thank you.

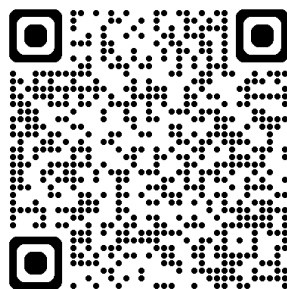
A Two-Stage Masked Autoencoder Based Network for Indoor Depth Completion

Kailai Sun, Zhou Yang, Qianchuan Zhao

Email: [skl23@nus.edu.sg](mailto:skl23@nus.edu.sg)

2024.6.18

Project:



WeChat:

