

# NeRF-MAE:

## Masked AutoEncoders for Self-Supervised 3D Representation Learning for Neural Radiance Fields

European Conference on Computer Vision, ECCV 2024

also appeared at CVPR Neural Rendering Intelligence Workshop, CVPR 2024



Zubair Irshad



Sergey Zakharov



Vitor Guizilini



Adrien Gaidon



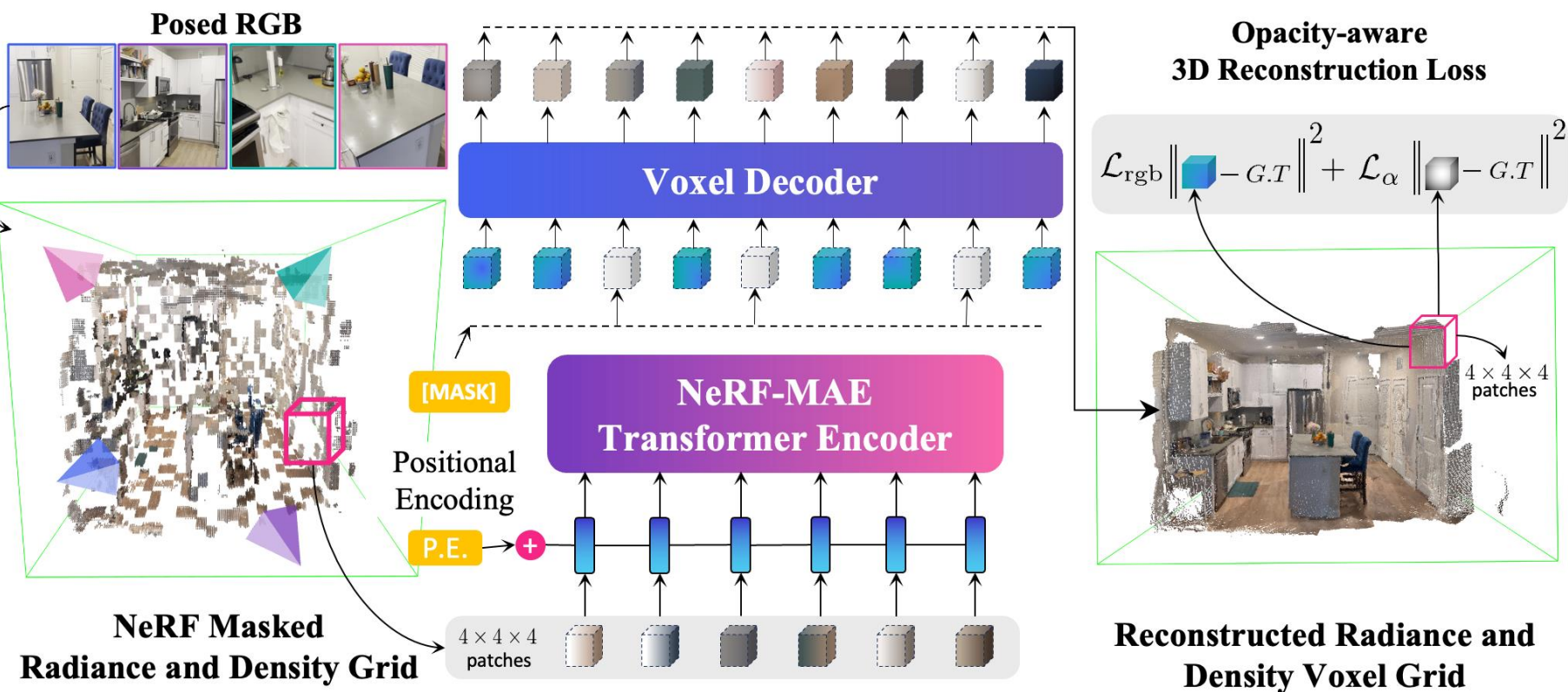
Zsolt Kira



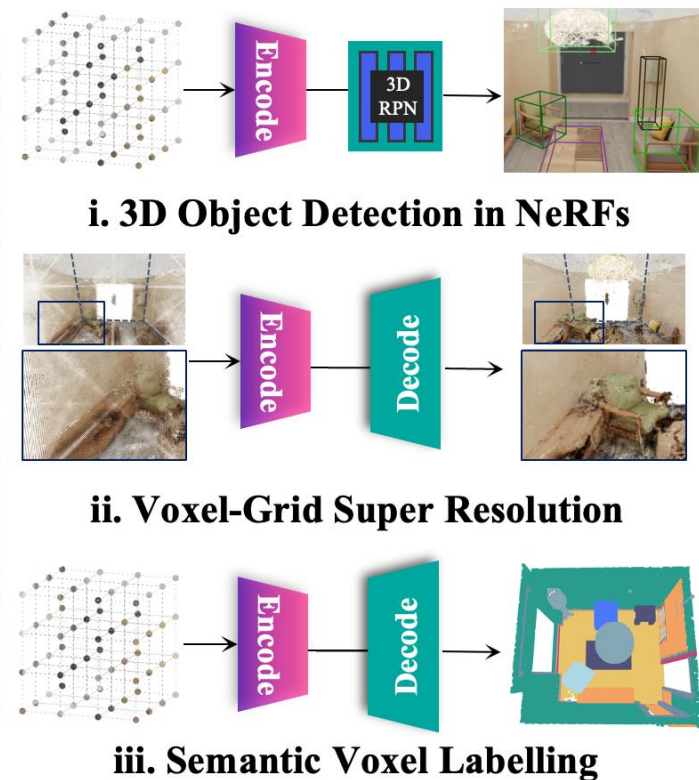
Rares Ambrus

# We propose **NeRF-MAE**, a framework for self-supervised **3D representation learning** from **NeRFs**

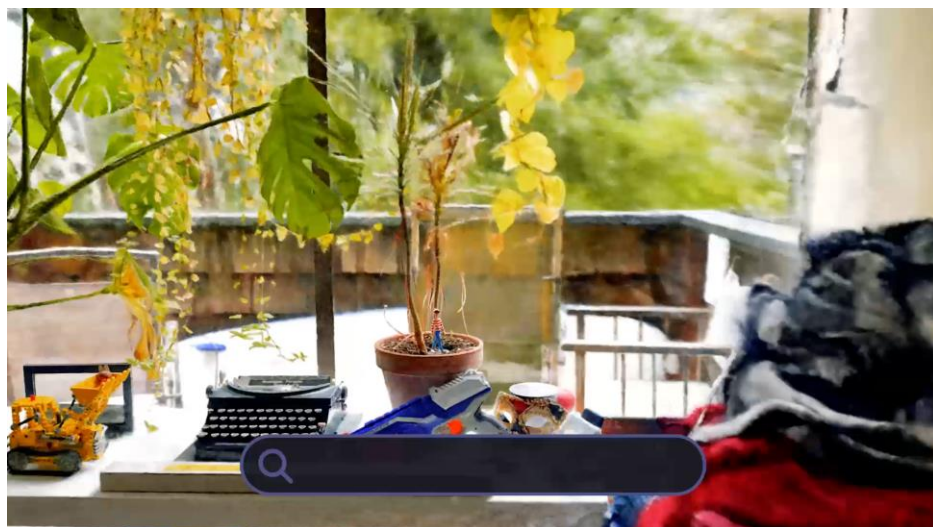
## a) Masked Pretraining Voxel-Grid Neural Radiance Fields



## b) Downstream 3D Tasks



# Neural Fields beyond showcasing high-rendering quality



Language-Embedded Radiance Fields  
(LeRF, Kerr et al)

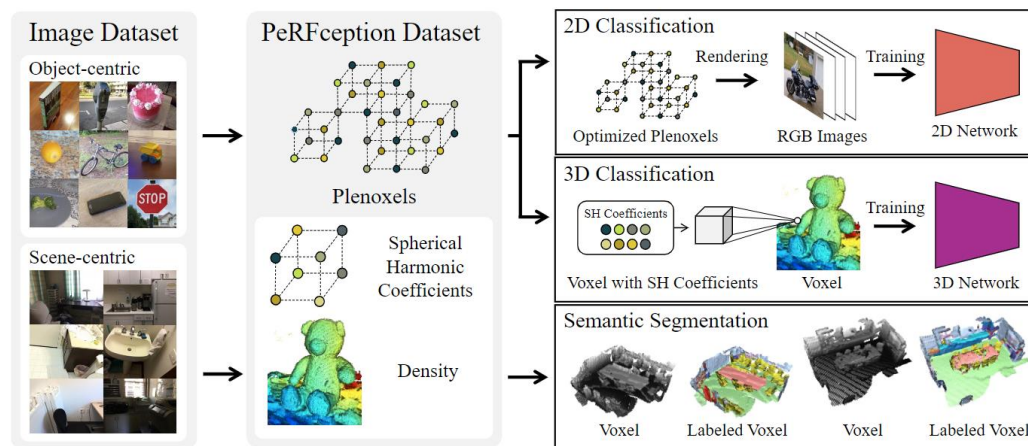


Inferring Accurate Geometry  
(NeRFMeshing, Rakotosaona et al)



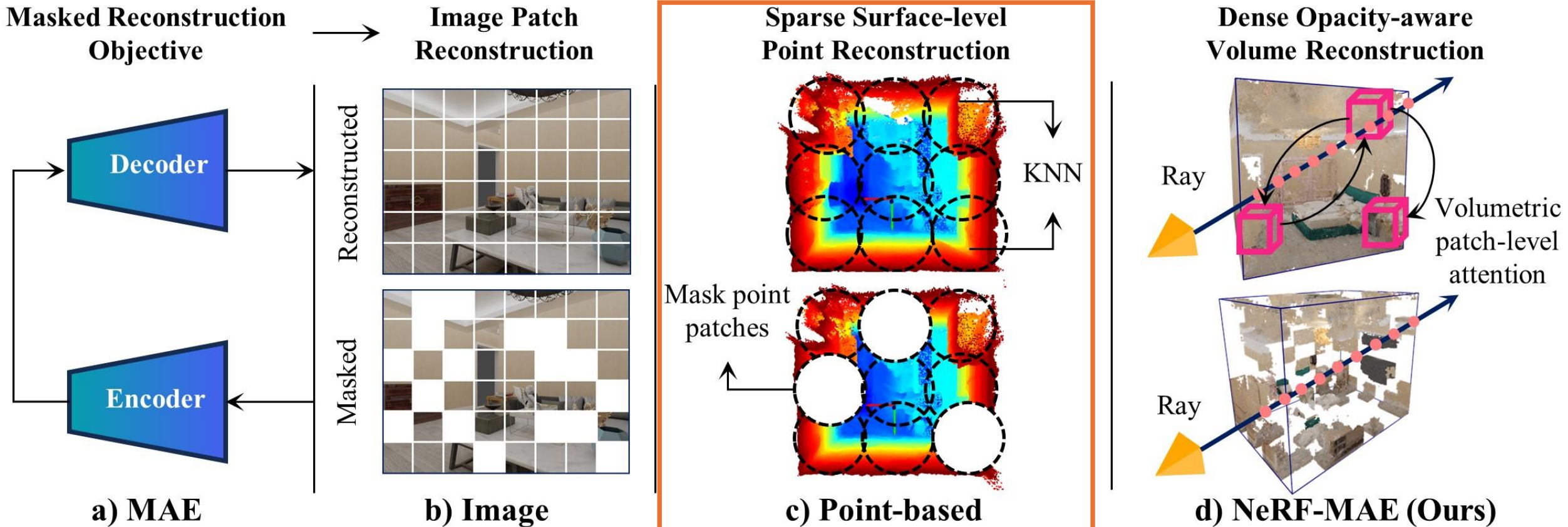
1. Scan Scene

Open-world Manipulation  
(F3RM, Shen et al)



Efficient Data Storage  
(PerFception, Jeong et al)

Existing **3D MAE Architectures** operate on pointclouds with **uneven information density** only model **surface-level information** and are **highly irregular data structures**



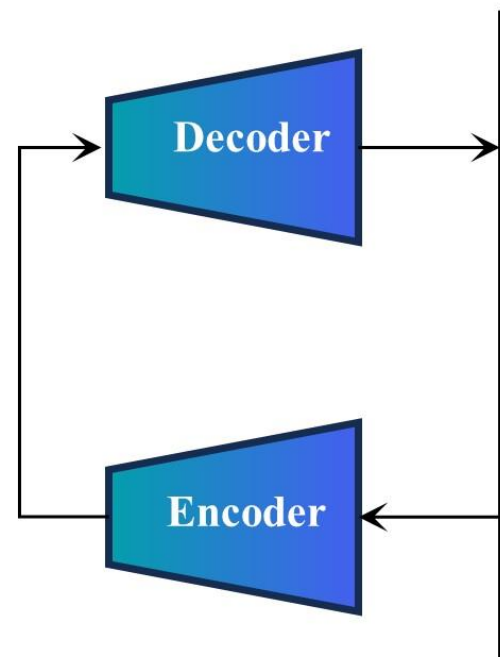
✗ Model Surface Level Information

✗ Irregular data Structures

✗ Uneven information density

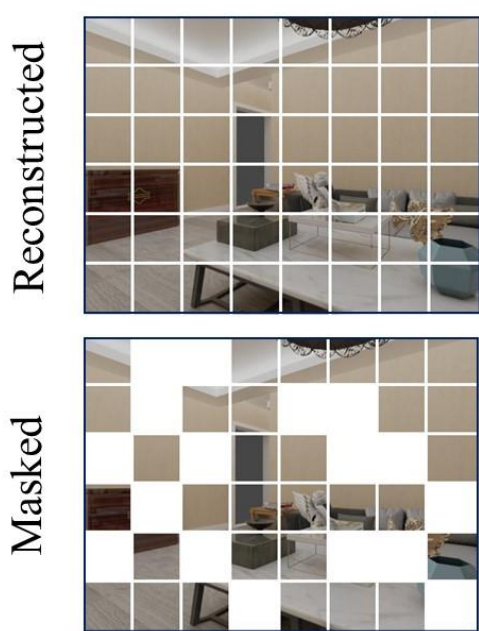
Our approach uses **NeRF's dense grid** as input to the Transformer. This makes our approach a direction extension of **image MAE to 3D**

Masked Reconstruction Objective



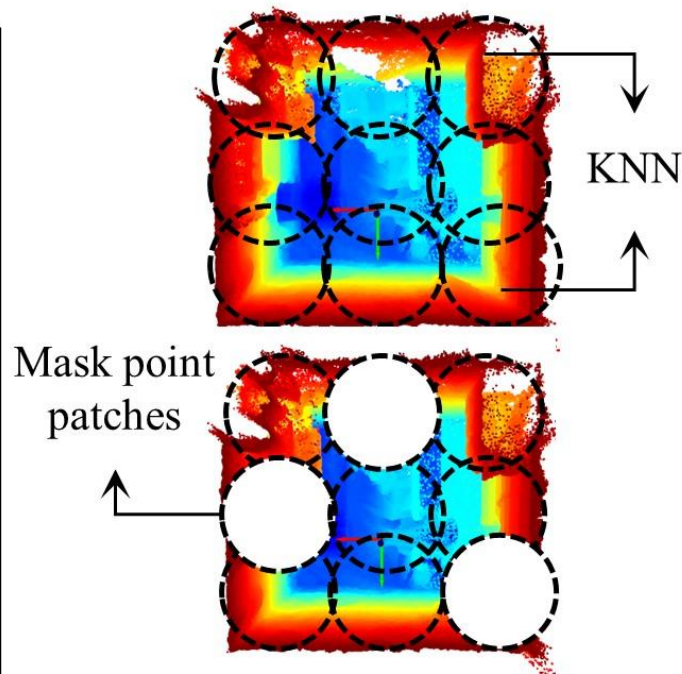
a) MAE

Image Patch Reconstruction



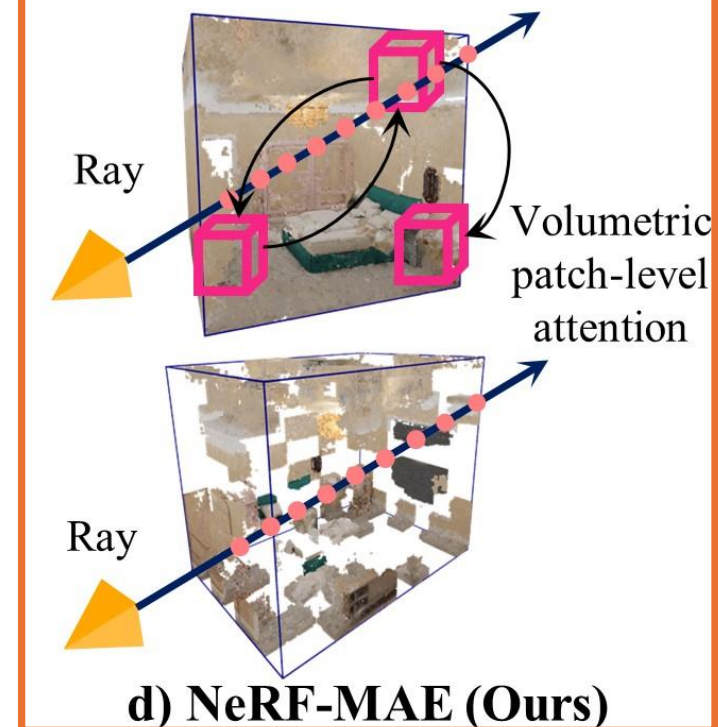
b) Image

Sparse Surface-level Point Reconstruction



c) Point-based

Dense Opacity-aware Volume Reconstruction



d) NeRF-MAE (Ours)

✓ High Information Density

✓ Regular unbiased Sampling

✓ Spatial data redundancy

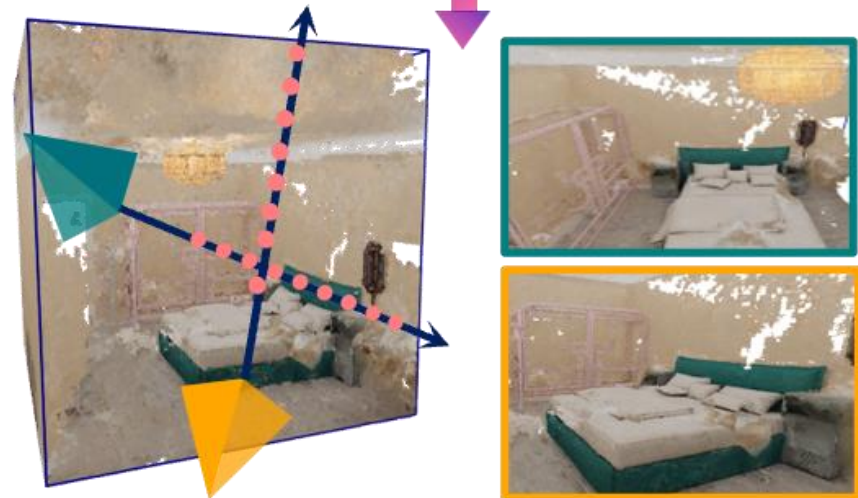
# Data Preprocessing flow for large-scale NeRF pretraining



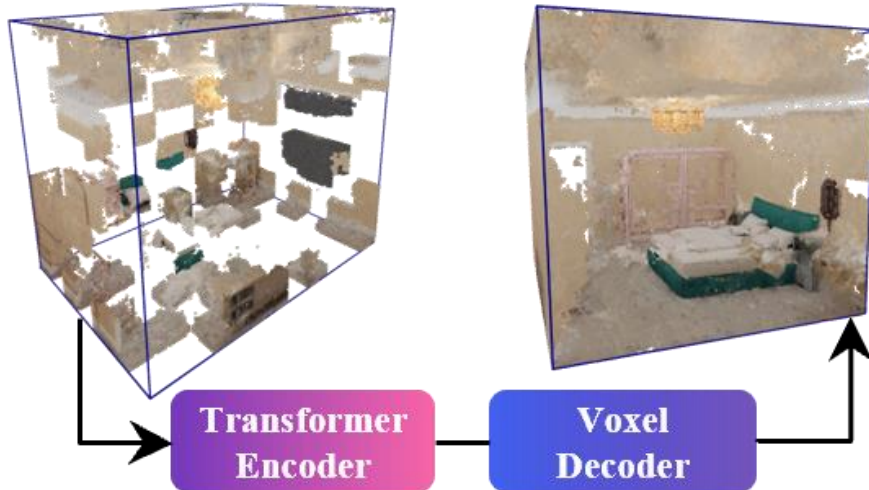
**a) Multi-view data**



**b) Trained NeRF**



**c) Extracted Radiance and Density Grid**

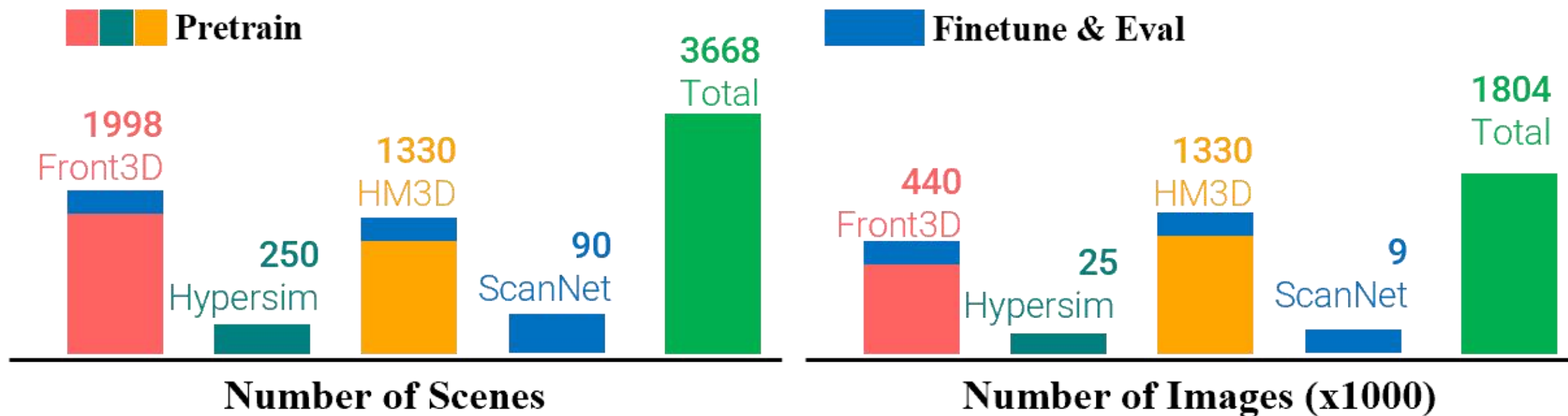


**d) NeRF-MAE Pretraining**

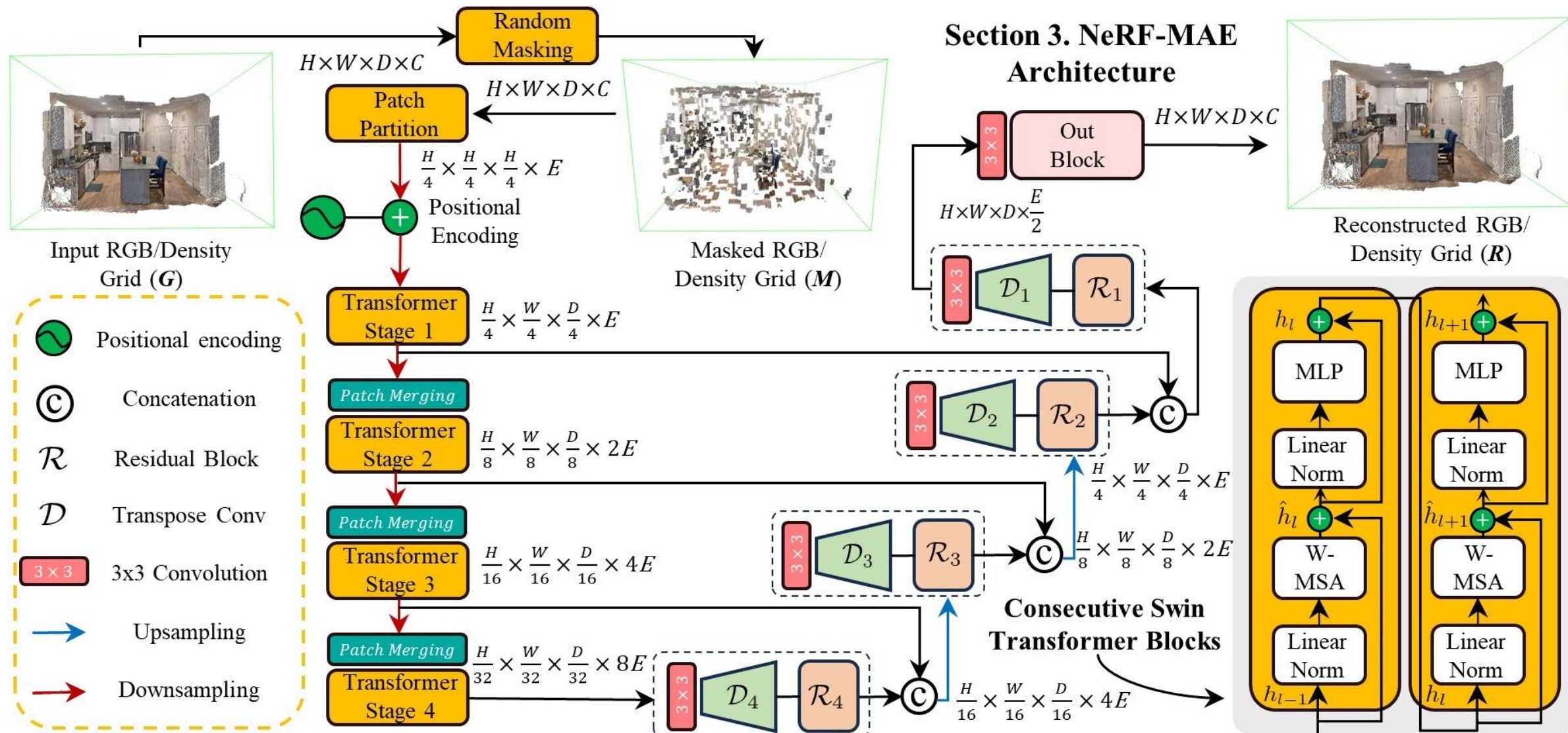
## a) Multi-view Dataset Setup



## b) NeRF-MAE Data Mix & Statistics

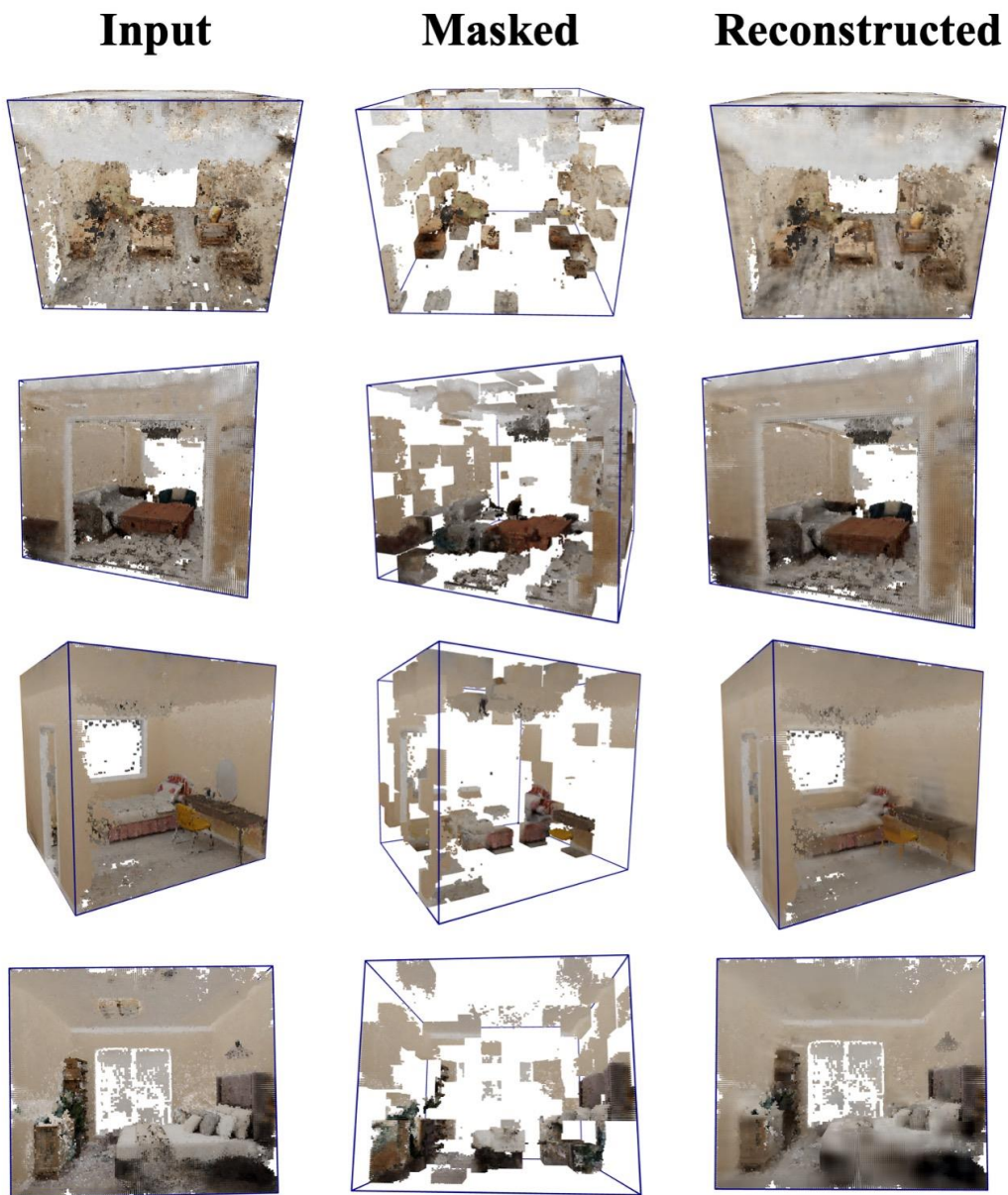


We utilize a **UNet architecture** employing SwinTransformer as encoder, lightweight voxel decoder to enforce **mask reconstruction objective in 3D**

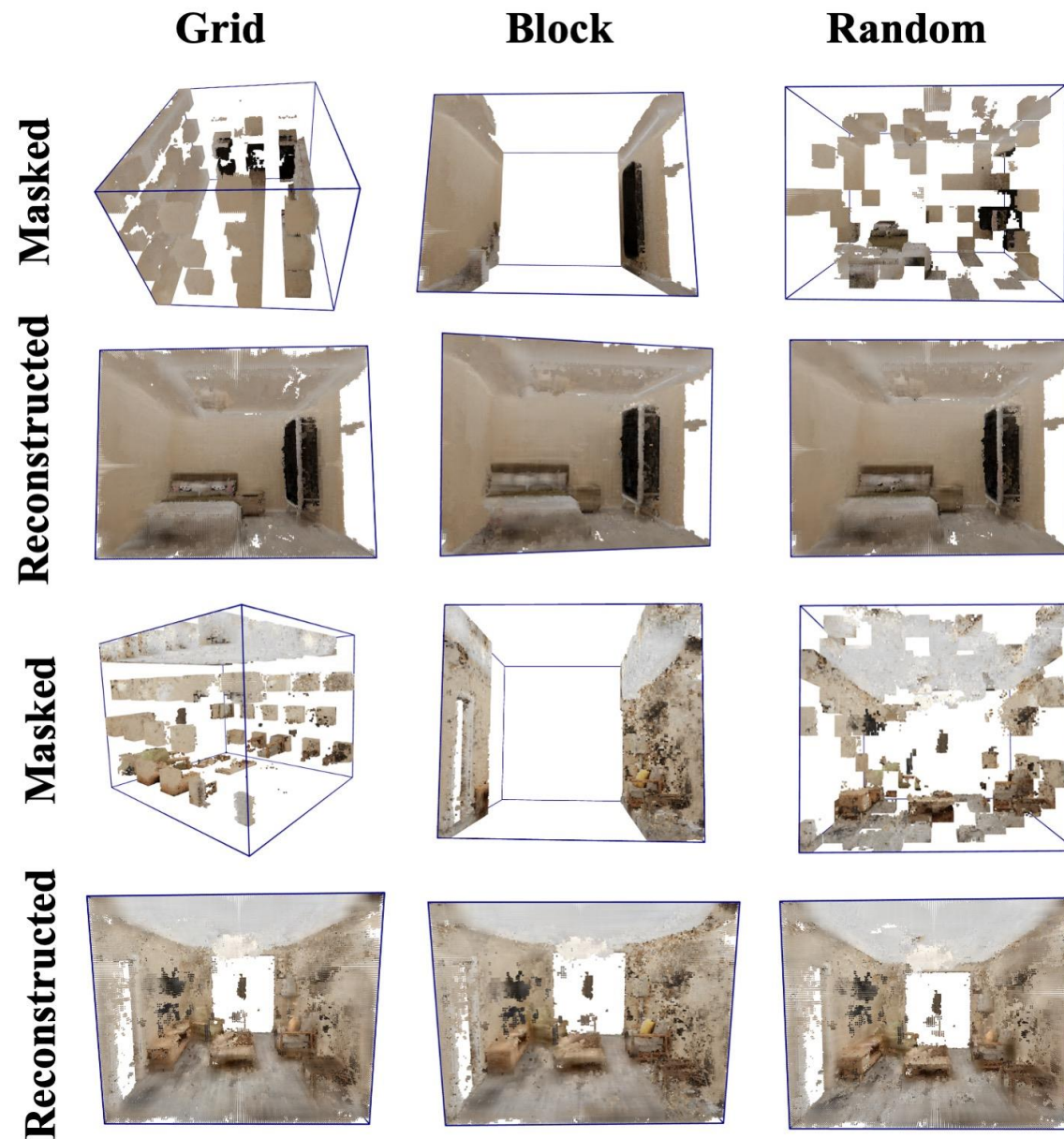




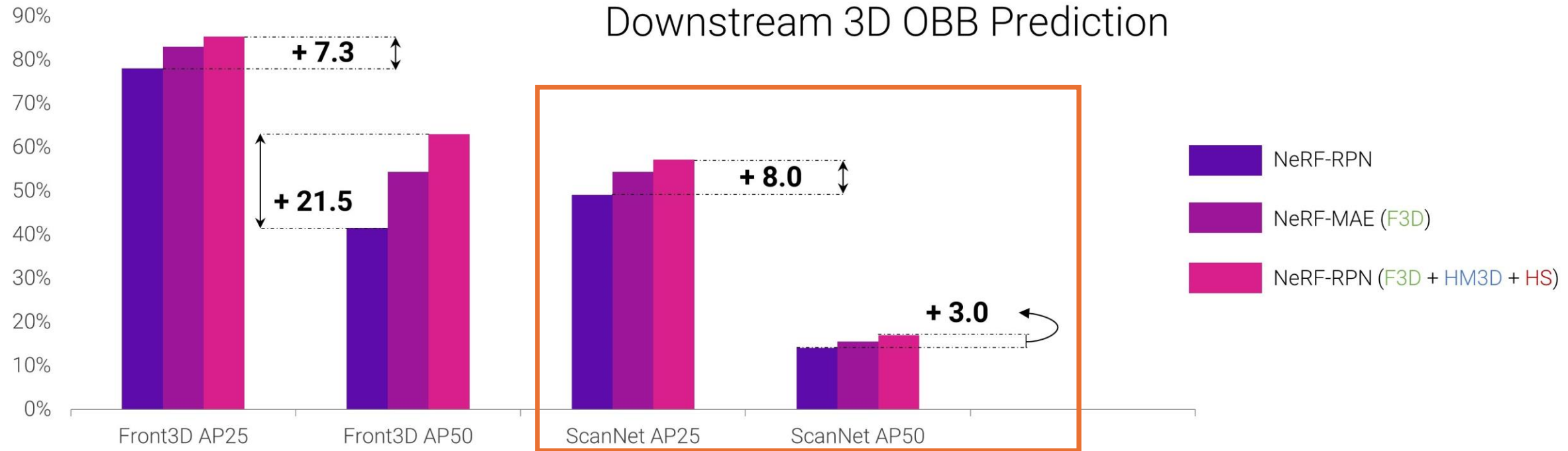
### a) *Qualitative Masked Reconstructions*



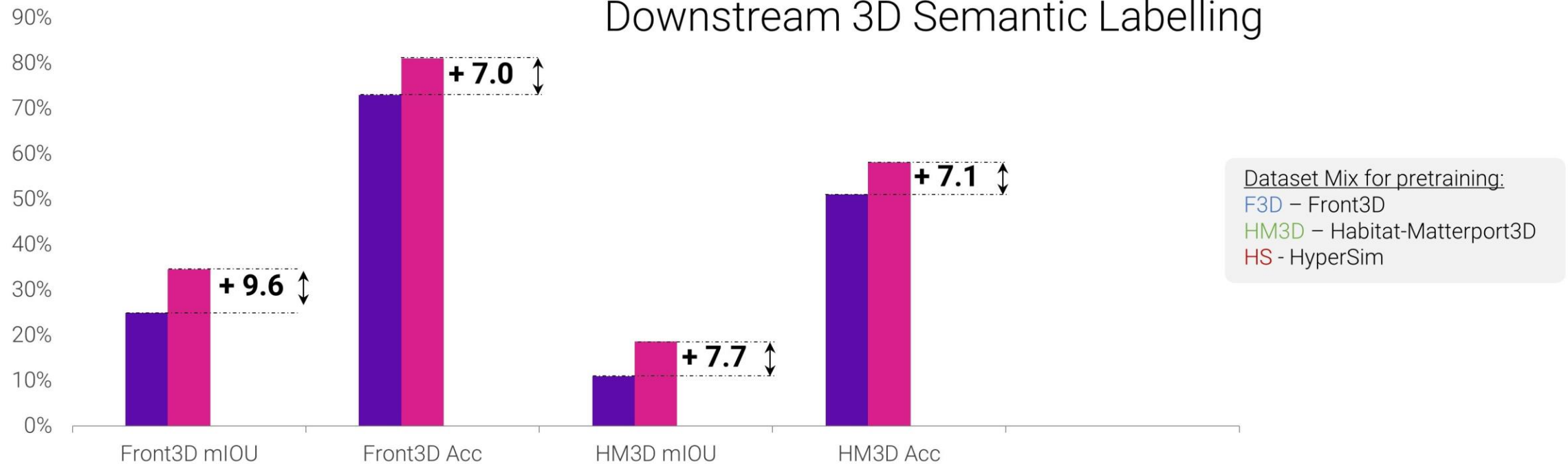
### b) *Masking Strategy Ablation*



## Downstream 3D OBB Prediction

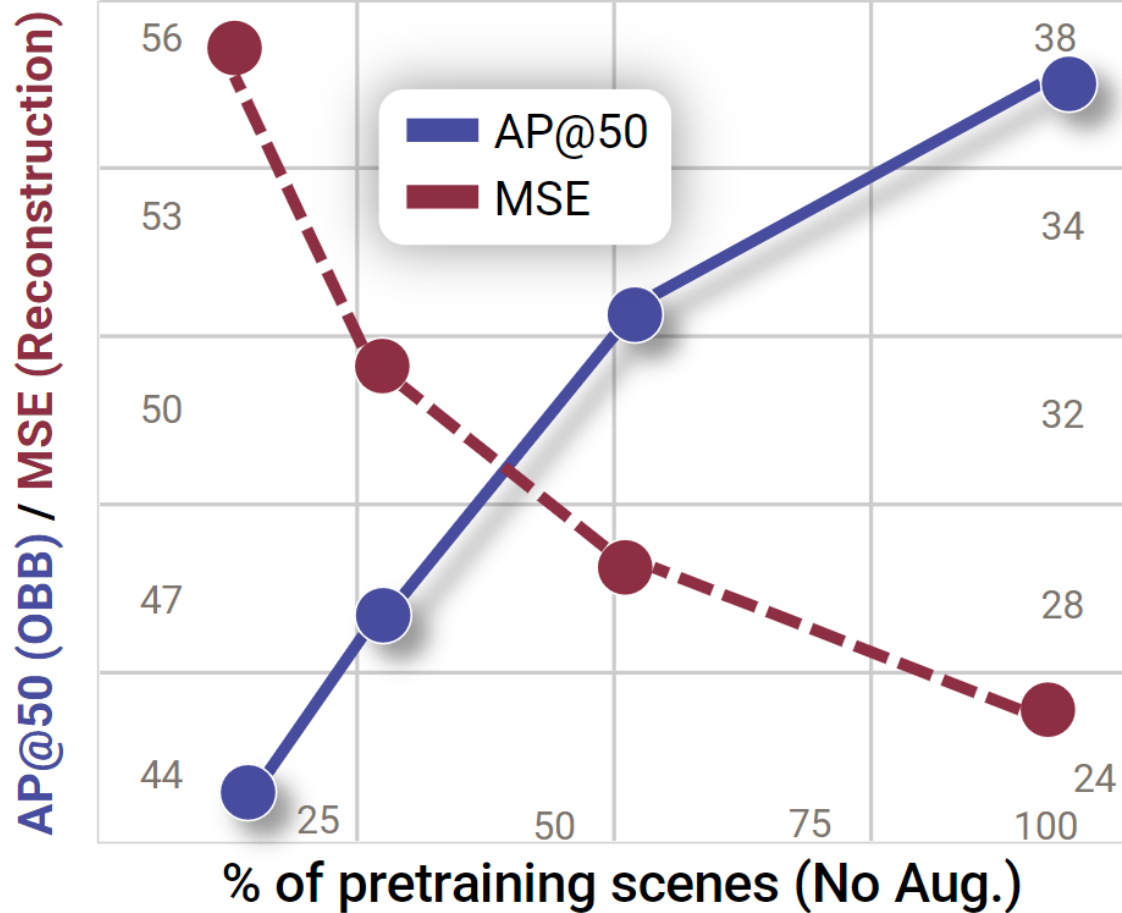


## Downstream 3D Semantic Labelling

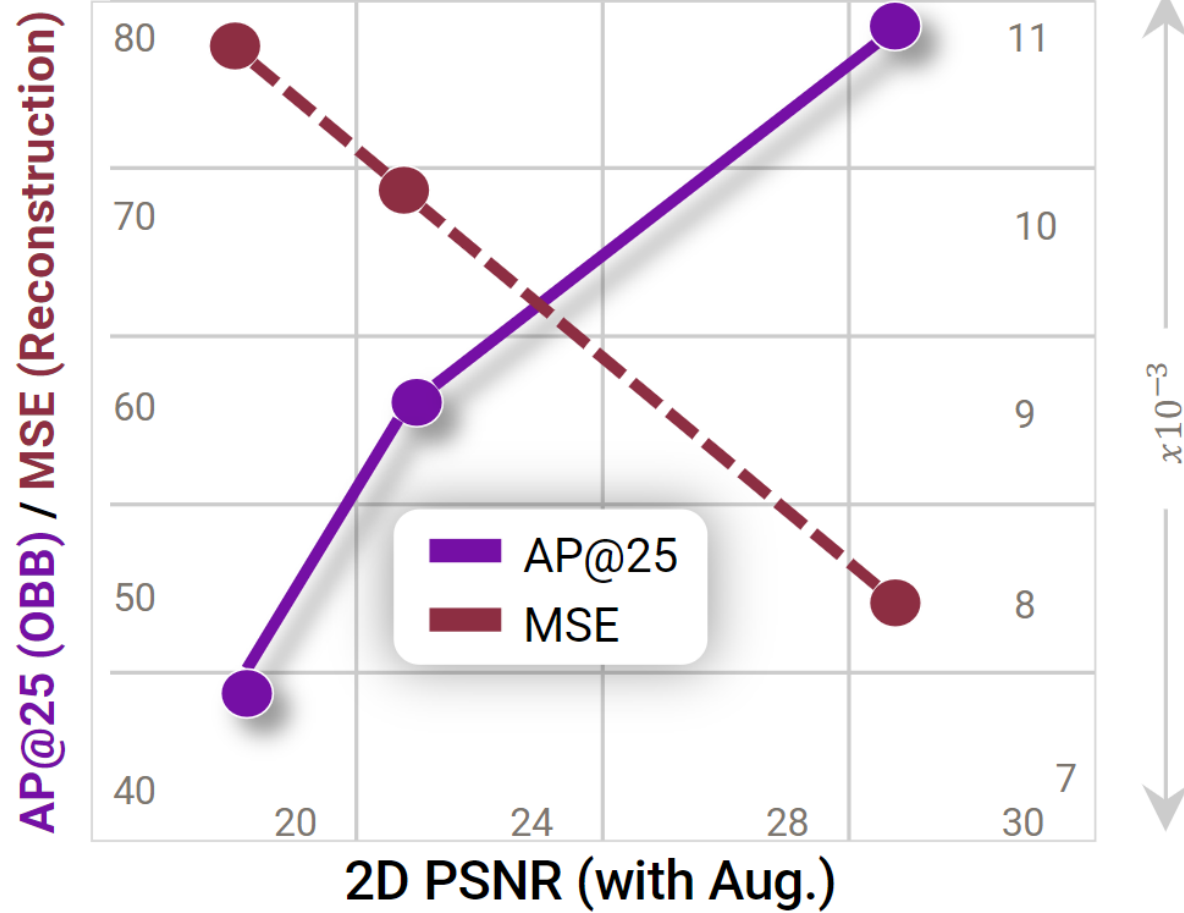


# Quantitative Results

## Scaling Performance of NeRF-MAE

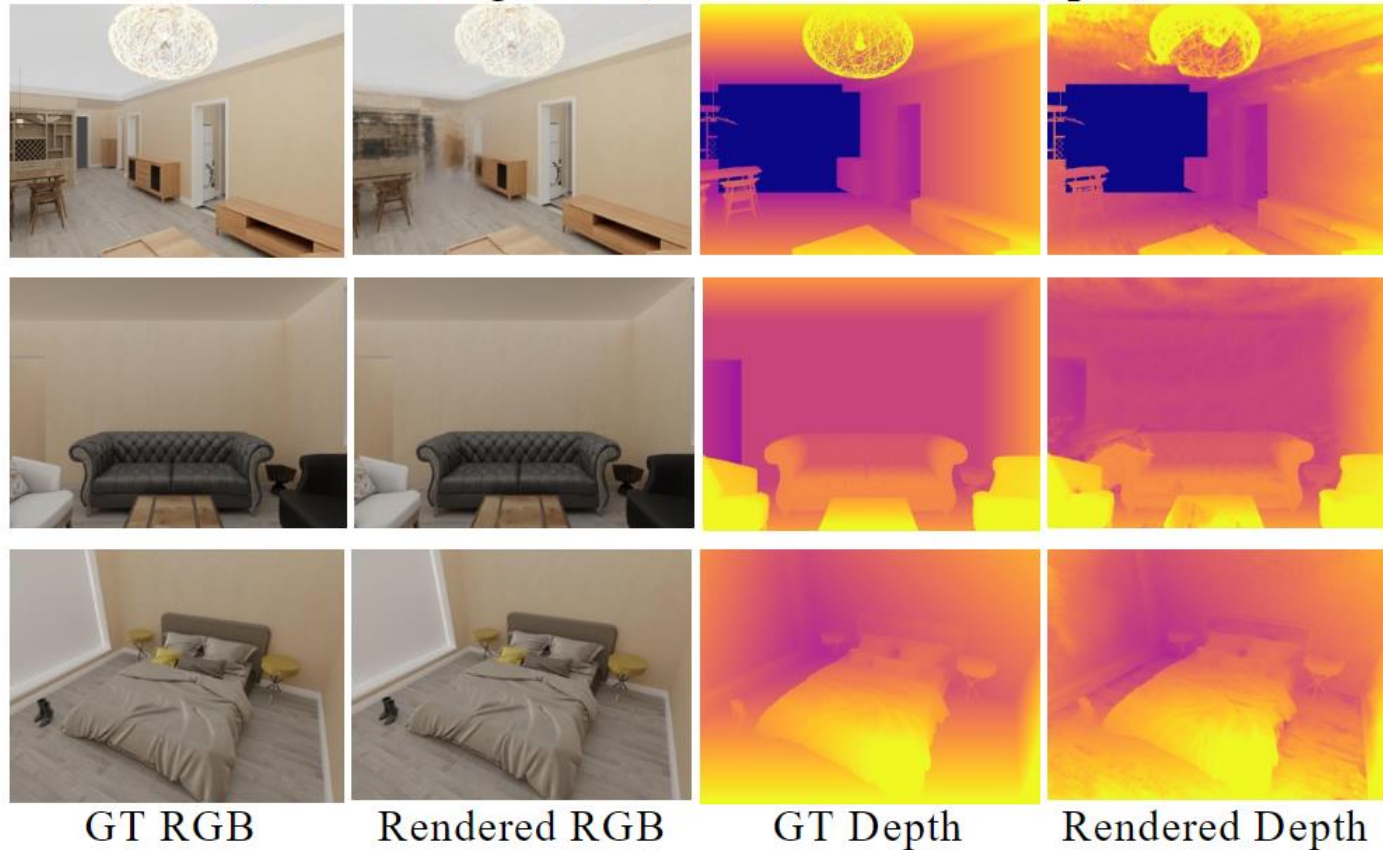


## NeRF Quality on Pretraining

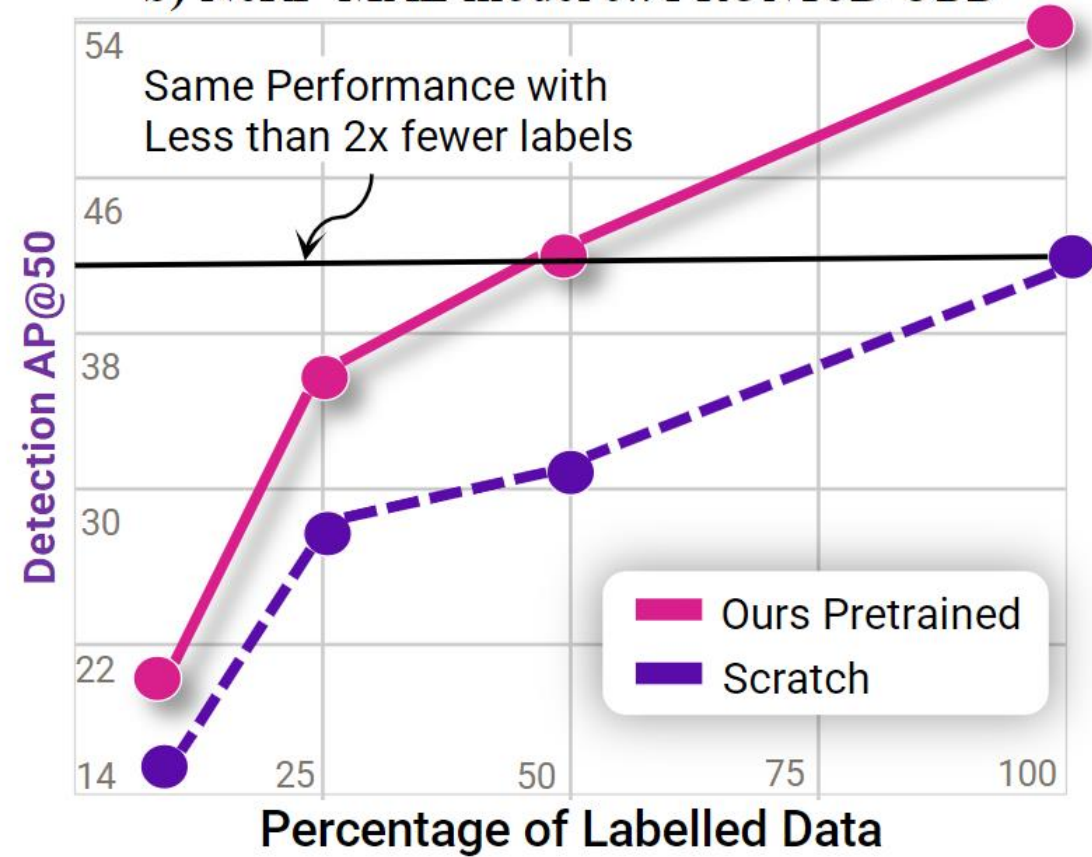


# Results Analysis

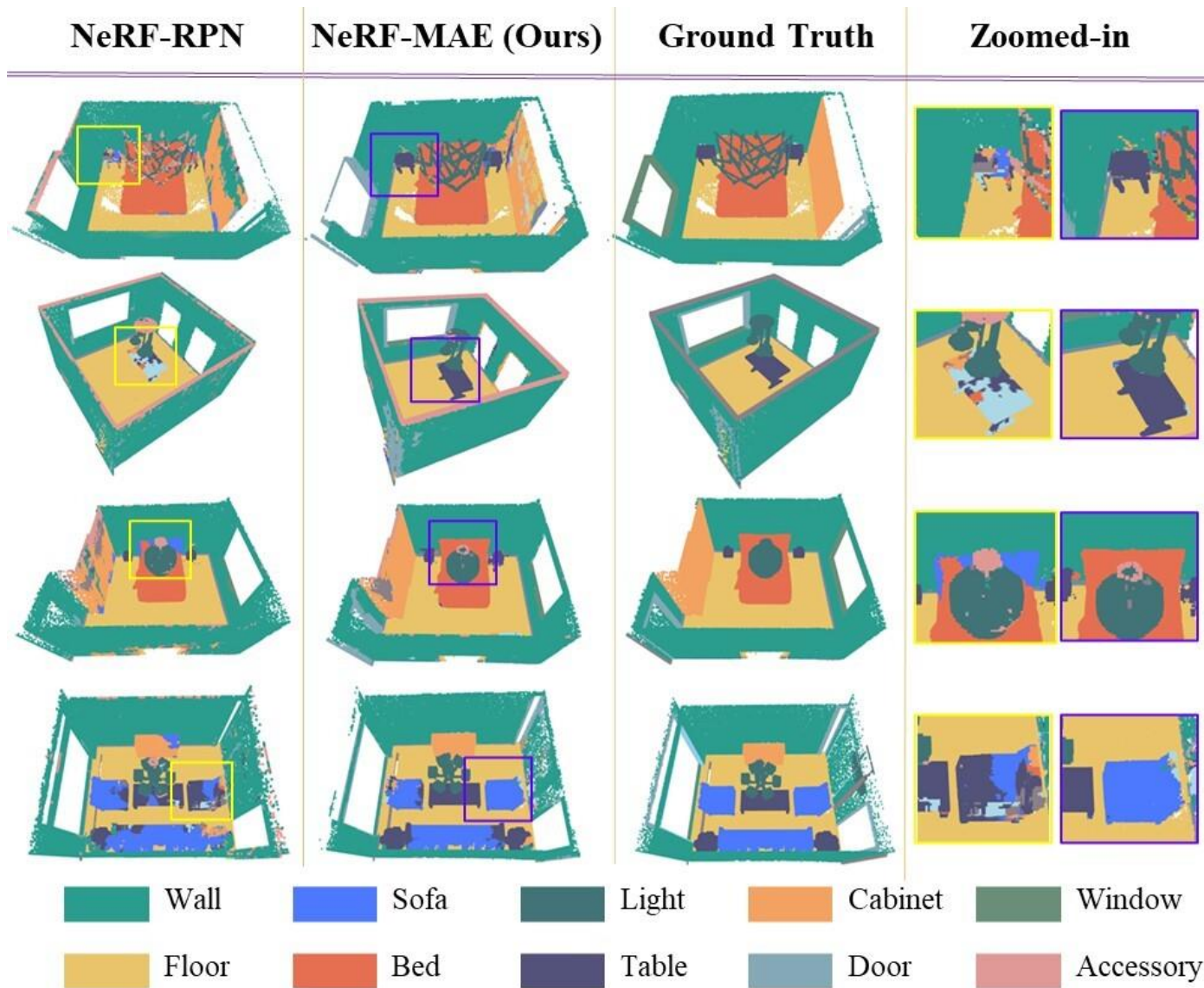
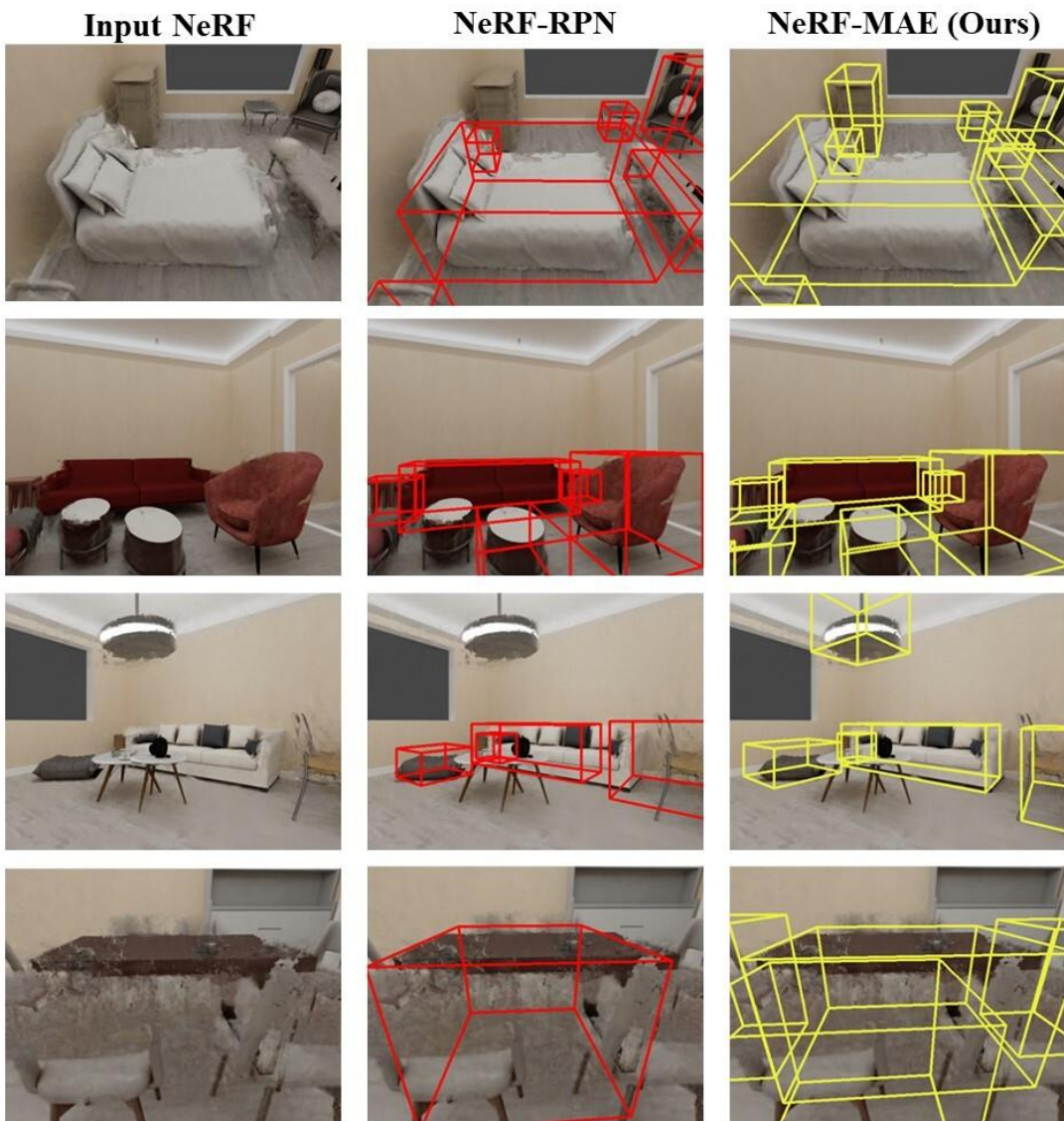
a) Pretraining Data (Rendered vs GT. Depth)



b) NeRF-MAE model *on FRONT3D OBB*



# Qualitative Results





# NeRF-MAE:

## Masked AutoEncoders for Self-Supervised 3D Representation Learning for Neural Radiance Fields

European Conference on Computer Vision, ECCV 2024

also appeared at CVPR Neural Rendering Intelligence Workshop, CVPR 2024

