

## Introduction

### Goal

- Recover the 3D shape of an object from single or multiple viewpoints.
- Unknown shape priors (e.g., lighting condition) and camera intrinsic and extrinsic parameters.

### The Problems of 3D-R2N2

- RNN is not permutation invariant
- Long-short memory loss of RNN
- RNN-based methods are time-consuming

### How to Solve it?

- Propose a new module, called *Context-aware Fusion*.

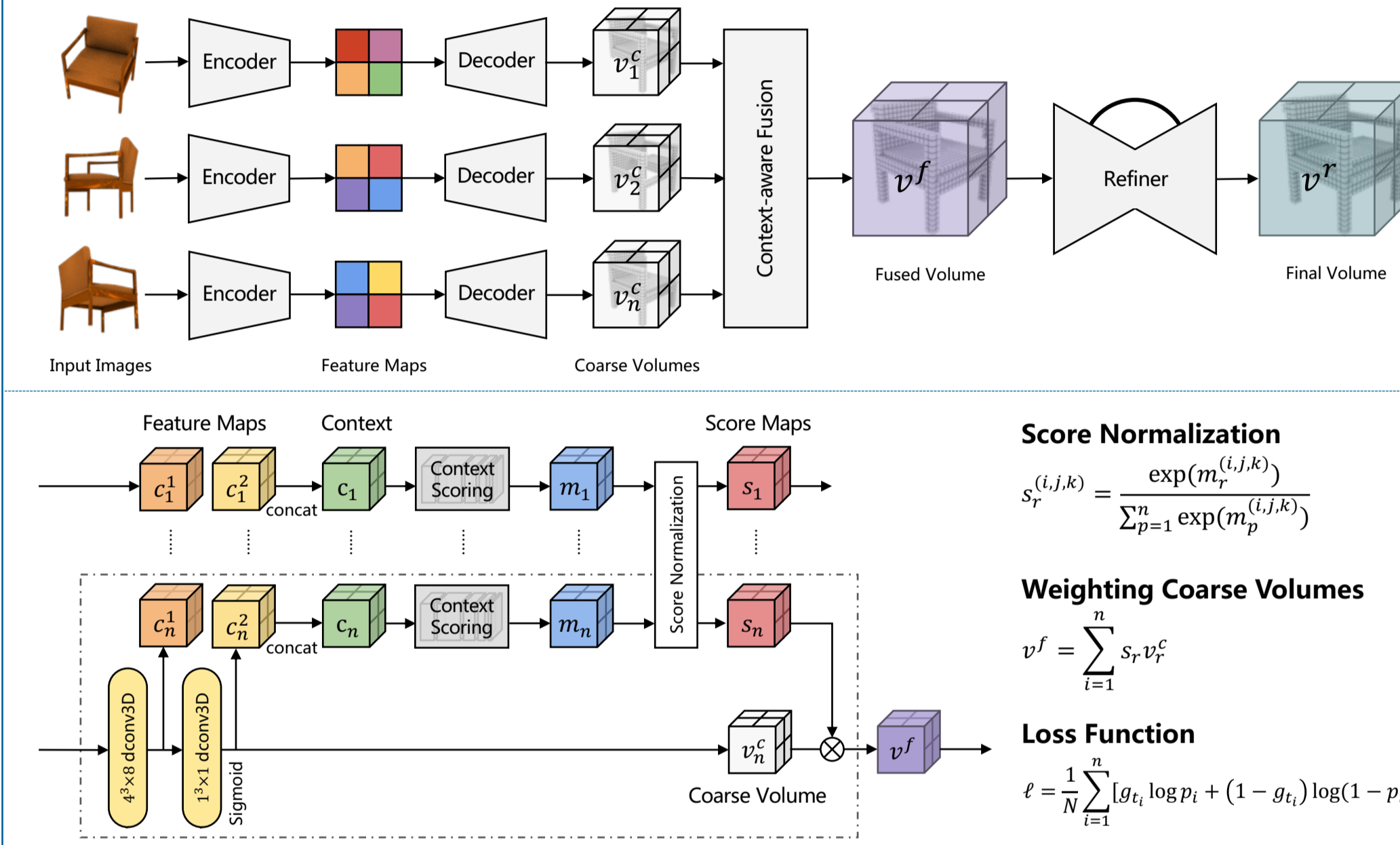
### The Advantages of Context-aware Fusion

- Selects high-quality reconstruction parts from different coarse 3D volumes in parallel.
- Better preserves multi-view spatial constraints when selecting reconstruction parts.

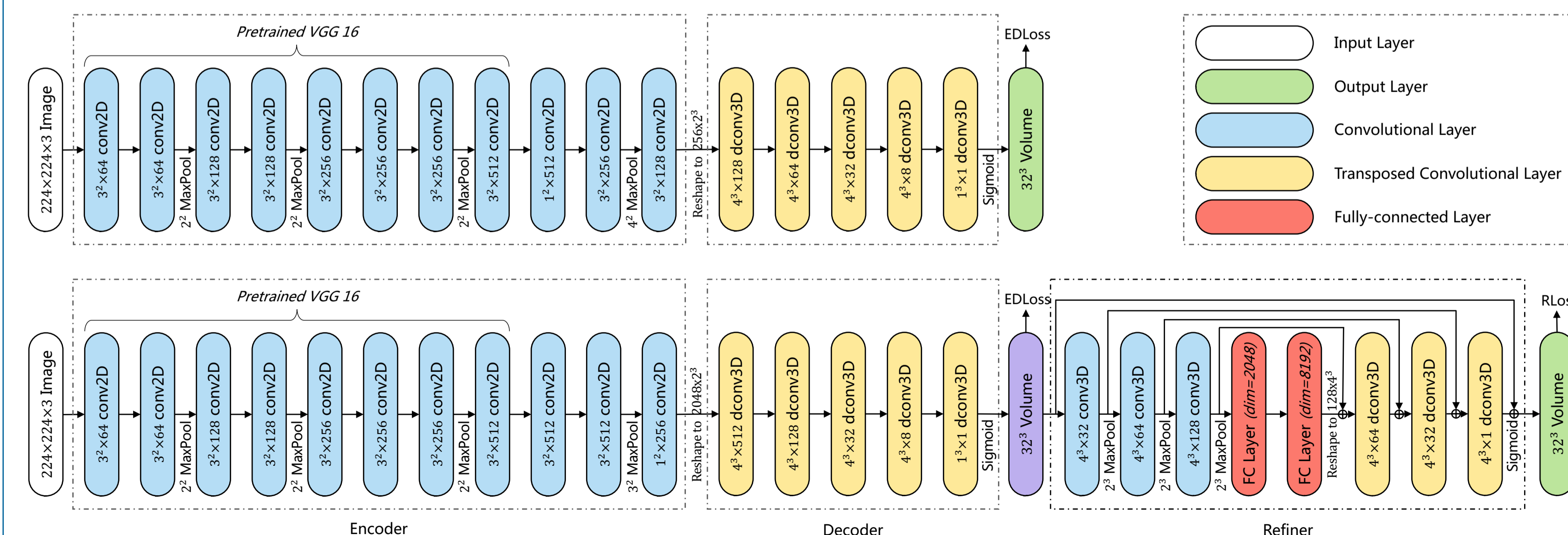
## Contributions

- A unified framework for both single-view and multi-view 3D object reconstruction, which consists of a well-designed encoder, decoder and refiner.
- A context-aware fusion module that adaptively selects high-quality reconstruction parts from different coarse 3D volumes to produce a fused reconstruction of the whole object.
- Outperform all state-of-the-arts methods on the ShapeNet and Pix3D datasets in terms of both accuracy and efficiency.
- Additional experiments show strong generalization abilities in reconstructing unseen 3D objects.

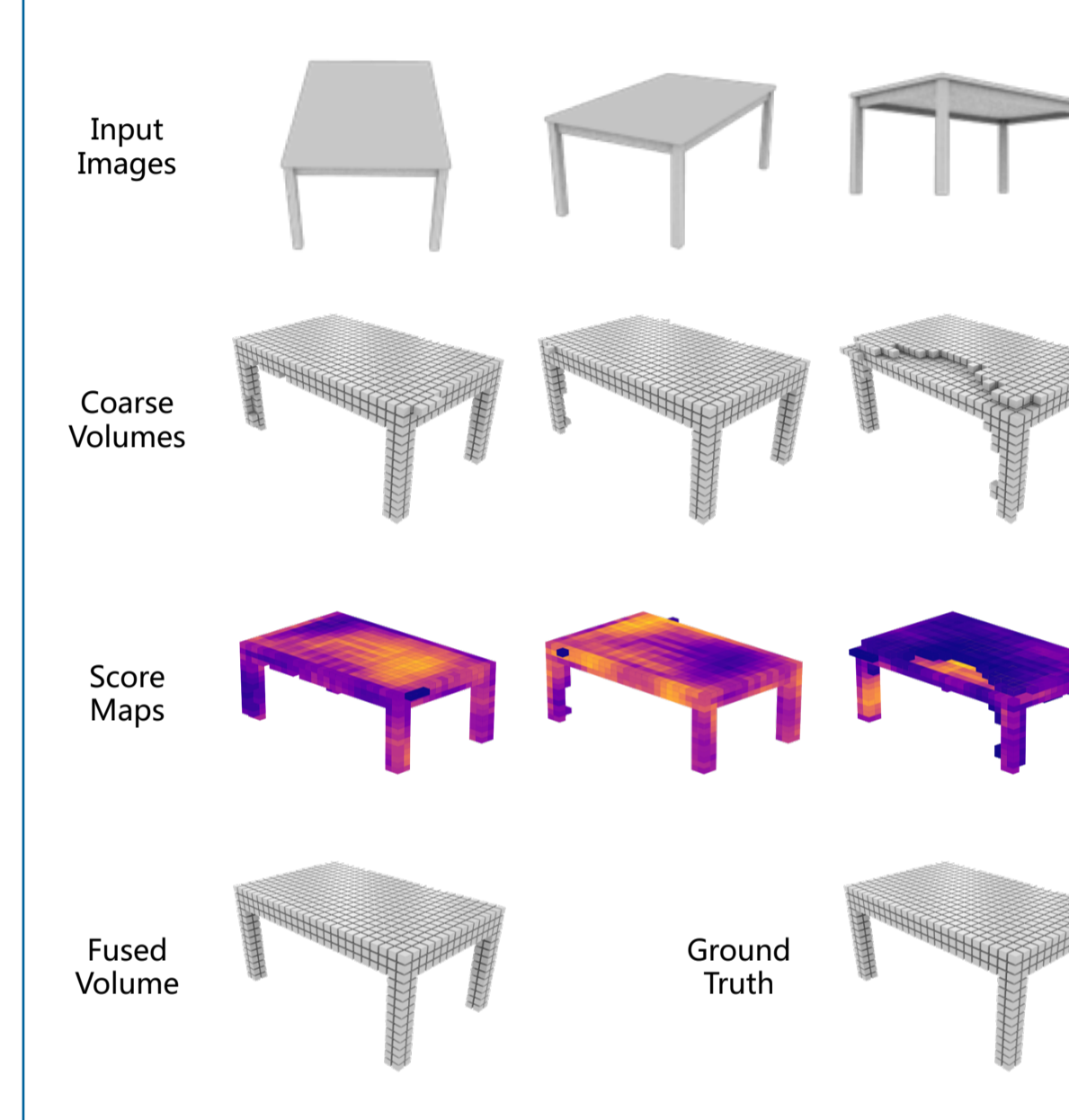
## Proposed Method



## Network Architecture



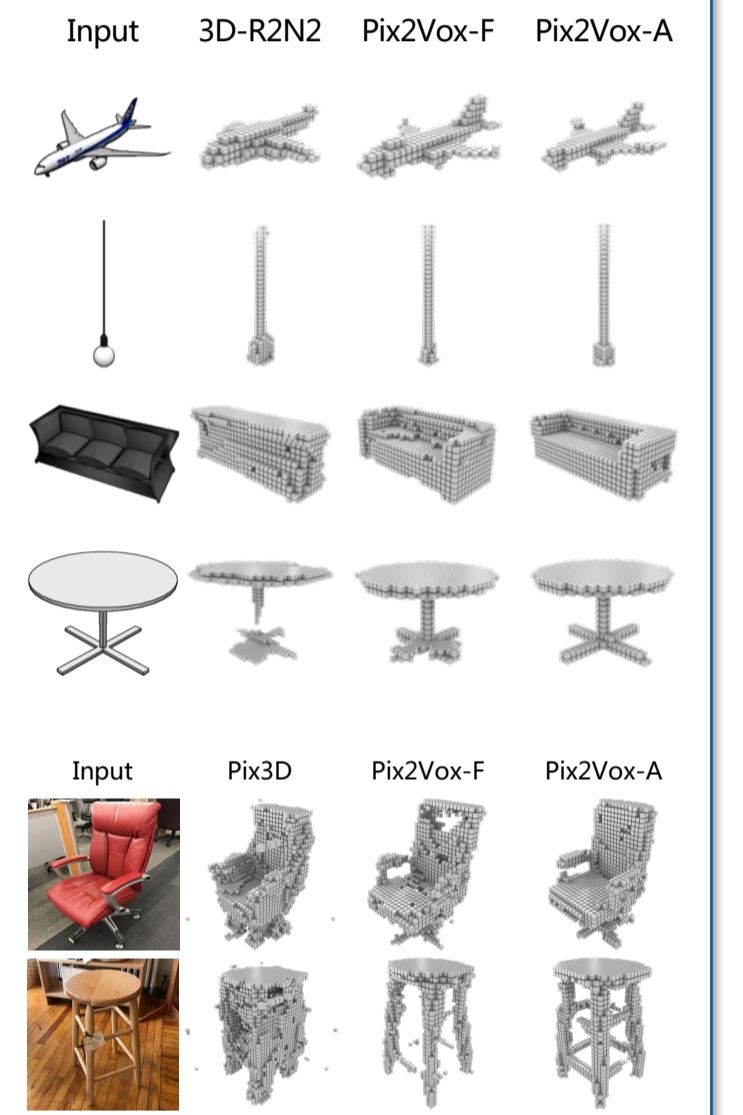
## Context-aware Fusion



## Experimental Results

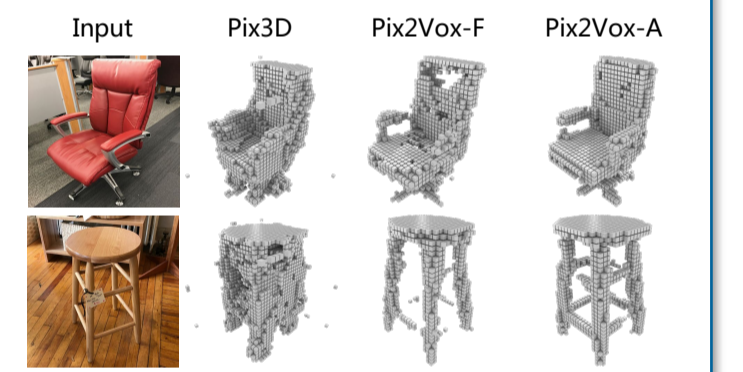
### Single-view Reconstruction on ShapeNet

Category	3D-R2N2	OGN	DRC	PSGN	Pix2Vox-F	Pix2Vox-A
airplane	0.513	0.587	0.571	0.601	0.600	0.684
bench	0.421	0.481	0.453	0.550	0.538	0.616
cabinet	0.716	0.729	0.635	0.771	0.765	<b>0.792</b>
car	0.798	0.828	0.755	0.831	0.837	<b>0.854</b>
chair	0.466	0.483	0.469	0.544	0.535	<b>0.567</b>
display	0.468	0.503	0.419	<b>0.552</b>	0.511	0.537
lamp	0.381	0.398	0.415	<b>0.462</b>	0.435	0.443
speaker	0.662	0.637	0.609	<b>0.737</b>	0.707	0.714
rifle	0.544	0.593	0.608	0.604	0.598	<b>0.615</b>
sofa	0.628	0.646	0.606	0.708	0.687	<b>0.709</b>
table	0.513	0.536	0.424	<b>0.606</b>	0.587	0.601
telephone	0.661	0.702	0.413	0.749	0.770	<b>0.776</b>
watercraft	0.513	<b>0.632</b>	0.556	0.611	0.582	0.594
IoU	0.560	0.596	0.545	0.640	0.634	<b>0.661</b>



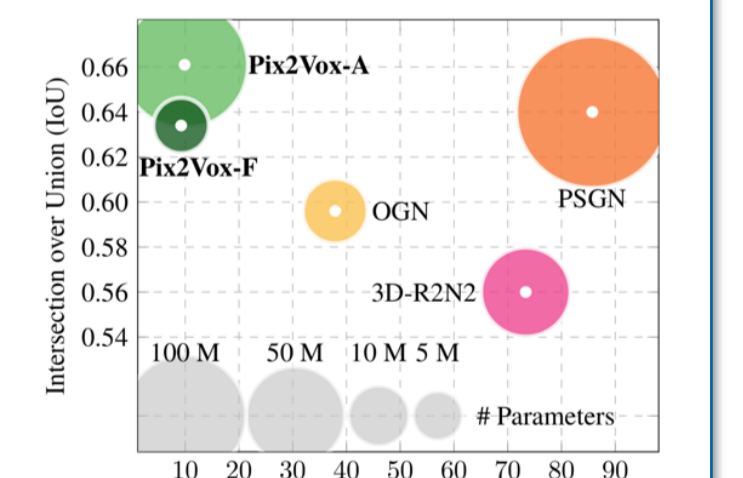
### Multi-view Reconstruction on ShapeNet

# views	1	2	3	4	5	8	12	16	20
3D-R2N2	0.560	0.603	0.617	0.625	0.634	0.635	0.636	0.636	0.636
Pix2Vox-F	0.634	0.660	0.668	0.673	0.676	0.680	0.682	0.684	0.684
Pix2Vox-A	<b>0.661</b>	<b>0.686</b>	<b>0.693</b>	<b>0.697</b>	<b>0.699</b>	<b>0.702</b>	<b>0.704</b>	<b>0.705</b>	<b>0.706</b>



### Memory usage and Running Time

Methods	3D-R2N2	OGN	PSGN	Pix2Vox-F	Pix2Vox-A
# Parameters (M)	35.97	12.46	210.46	7.41	114.24
Memory (MB)	1407	793	2733	673	2729
Training Time (hours)	169	192	50	12	25
Backward Time (ms)	312.50	312.25	97.90	12.93	72.01
Forward Time, 1 view (ms)	73.35	37.90	85.73	9.25	9.90
Forward Time, 2 views (ms)	108.11	N/a	N/a	12.05	13.69
Forward Time, 4 views (ms)	112.36	N/a	N/a	23.26	26.31



## References

- [3D-R2N2] 3D-R2N2: A Unified Approach for Single and Multi-view 3D Object Reconstruction. *ECCV 2016*.
- [OGN] Octree Generating Networks: Efficient Convolutional Architectures for High-resolution 3D Outputs. *JCCV 2017*.
- [DRC] Multi-view Supervision for Single-View Reconstruction via Differentiable Ray Consistency. *CVPR 2017*.
- [PSGN] A Point Set Generation Network for 3D Object Reconstruction from a Single Image. *CVPR 2017*.