



Neural Information  
Processing Systems  
(NeurIPS) 2019

# Learning Object Bounding Boxes for 3D Instance Segmentation on Point Clouds



B. Yang,



J. Wang,



R. Clark,



Q. Hu,



S. Wang,



A. Markham,

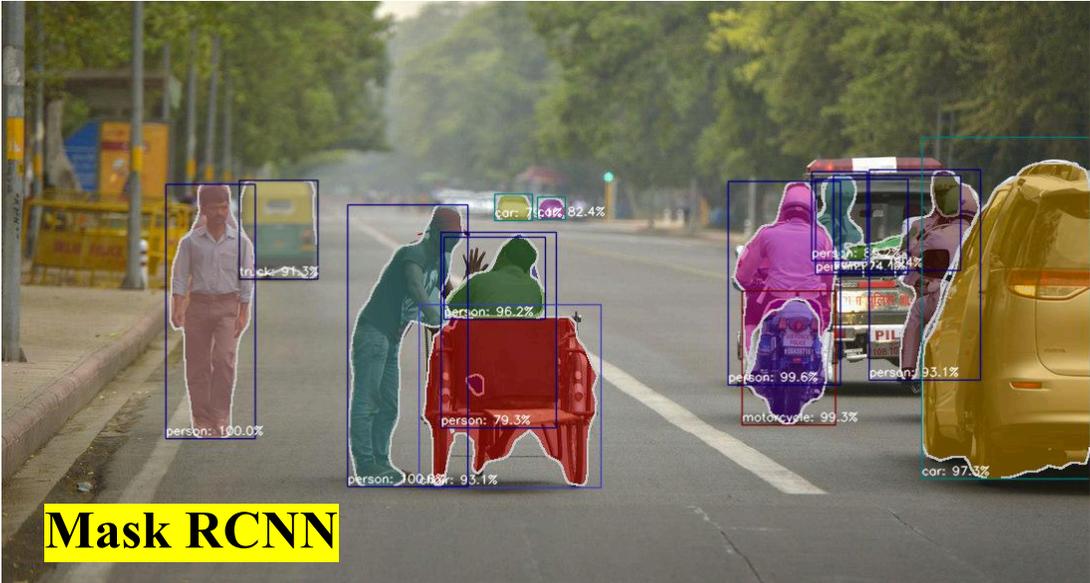


N. Trigoni

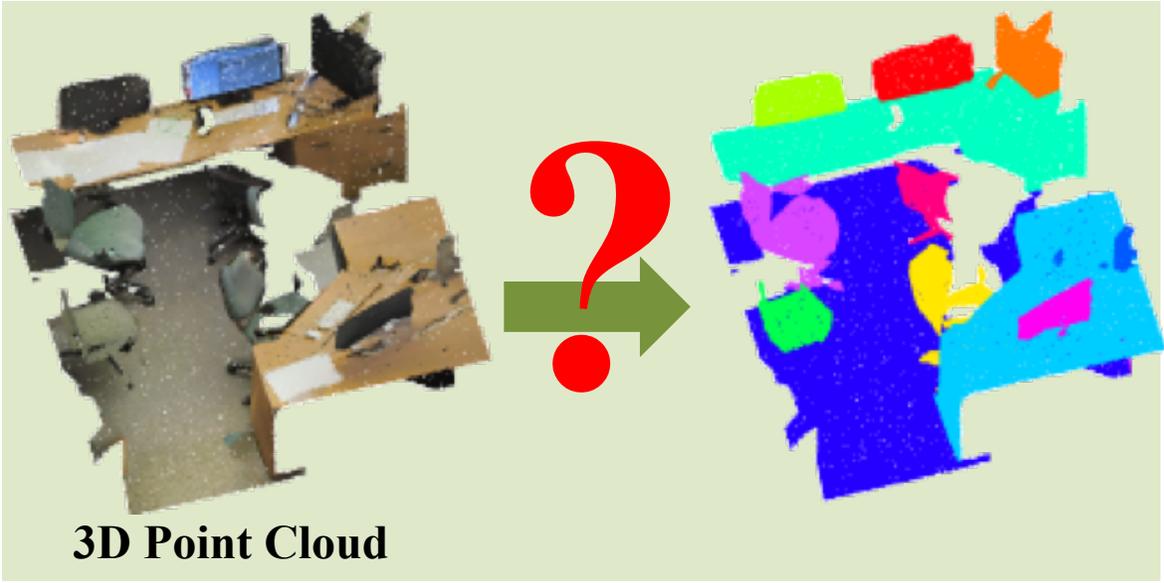


# Background:

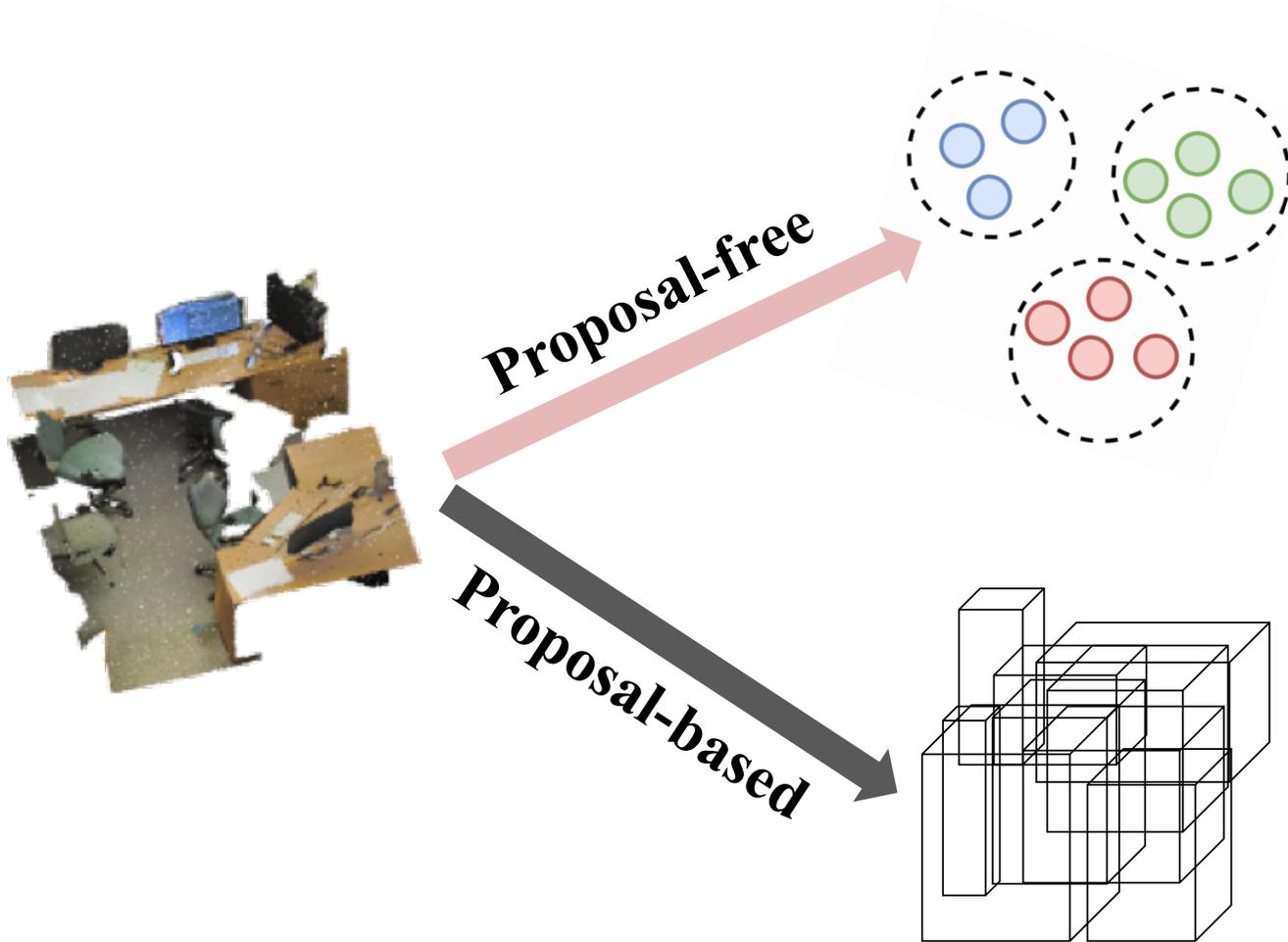
## 2D Instance Segmentation



## 3D Instance Segmentation



# Background:



## Limitations

- **Low objectness**
- **Heavy post-processing (grouping)**



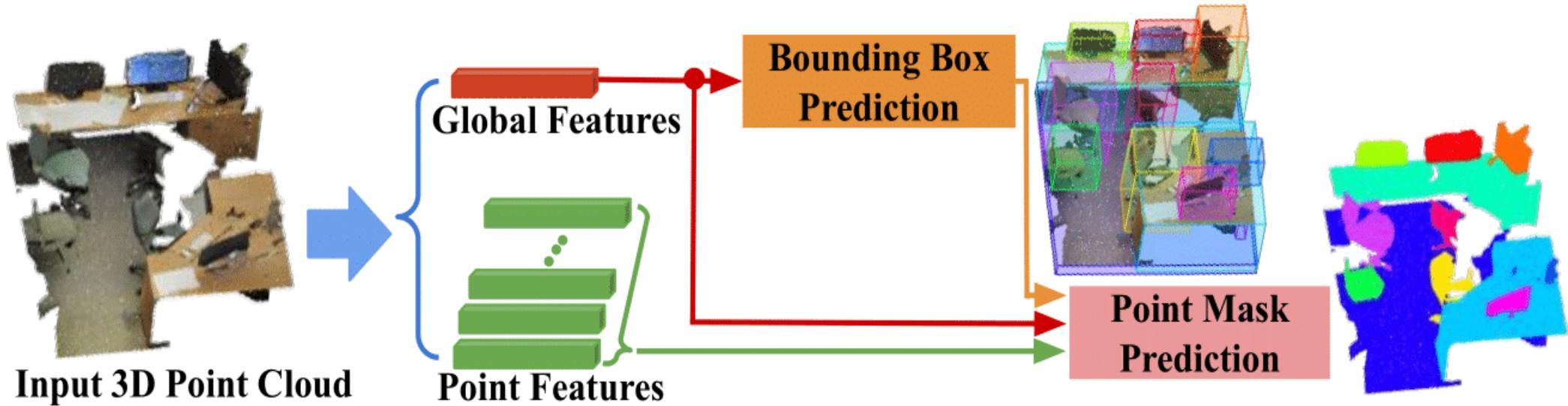
*SGPN (CVPR'18); ASIS (CVPR'19);  
JSIS3D (CVPR'19); 3D-BEVIS (GCPR'19);  
MTML (ICCV'19); MASC (arXiv'19)*

- **Two-stage training**
- **Heavy post-processing (NMS)**



*3D-SIS (CVPR'19); GSPN (CVPR'19)*

# Our Method (3D-BoNet):



## Highlights of our pipeline

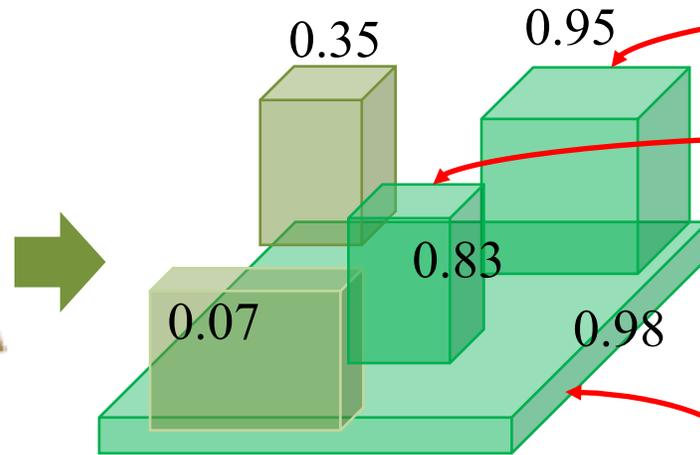
- Each object is **uniquely detected and segmented**.
- The learnt 3D bounding boxes guarantee **high objectness**.
- It's **end-to-end trainable, no post-processing**, and efficient.



# Our Method (3D-BoNet):

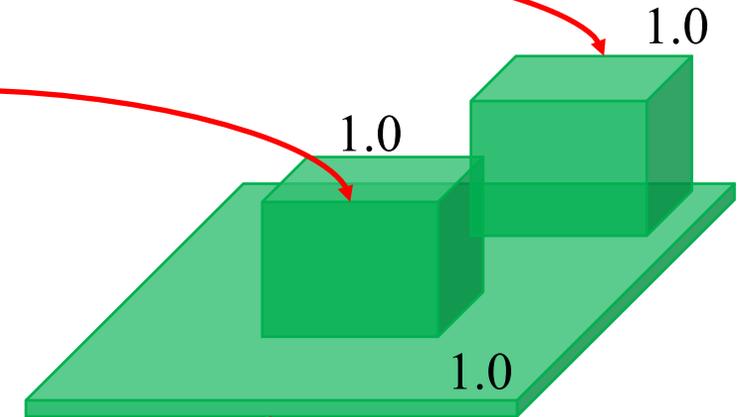


Input  
Point Cloud



- Predicted bounding boxes
- Predicted bbox scores

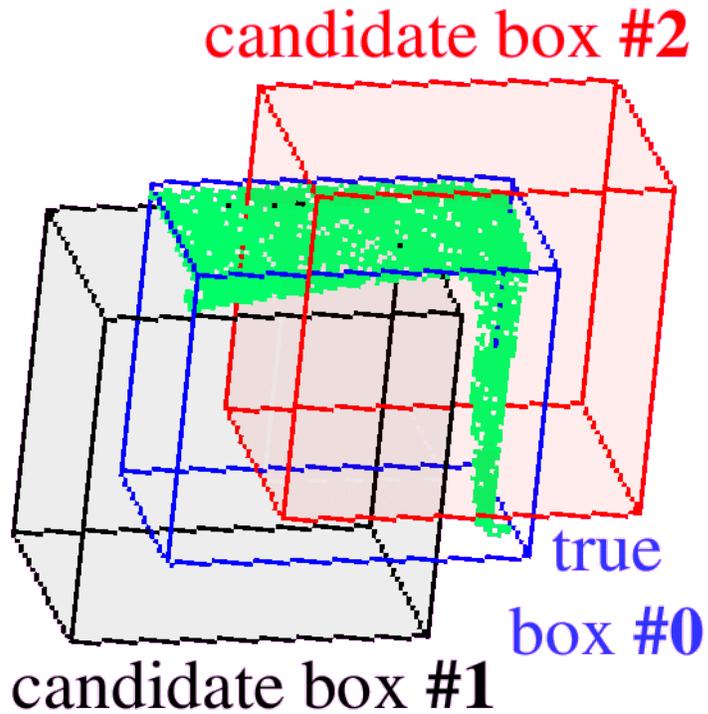
## Optimal Association



- ✓ GT bounding boxes
- ✓ GT bbox scores

$$\mathbf{A} = \arg \min_{\mathbf{A}} \sum_{i=1}^H \sum_{j=1}^T C_{i,j} A_{i,j} \quad \text{subject to} \quad \sum_{i=1}^H A_{i,j} = 1, \sum_{j=1}^T A_{i,j} \leq 1, j \in \{1..T\}, i \in \{1..H\}$$

# Our Method (3D-BoNet):



## Multiple criteria to match a pred bbox with a GT bbox

---

$$C_{i,j}^{ed} = \frac{1}{6} \sum (B_i - \bar{B}_j)^2$$

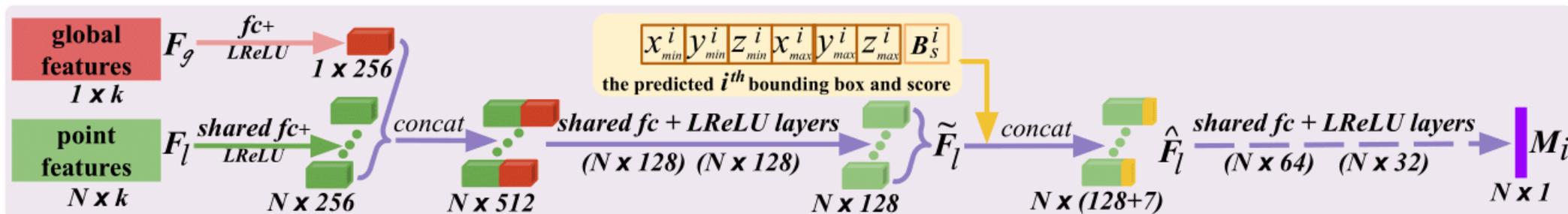
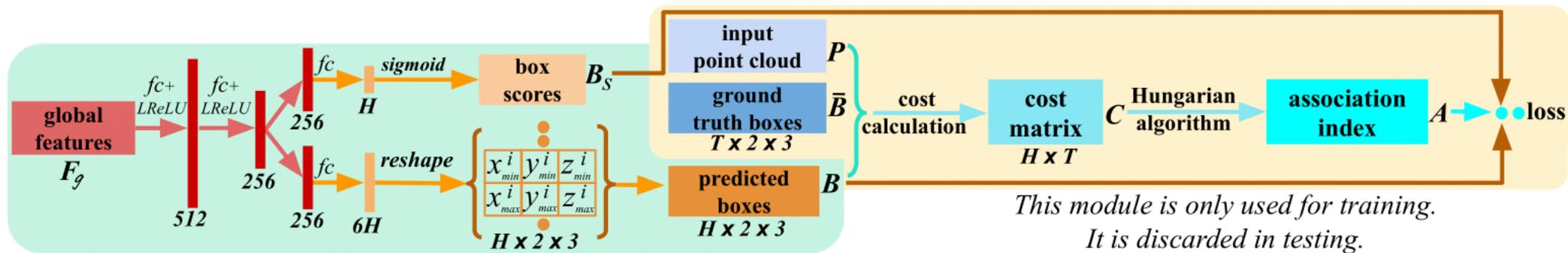
$$C_{i,j}^{sIoU} = \frac{-\sum_{n=1}^N (q_i^n * \bar{q}_j^n)}{\sum_{n=1}^N q_i^n + \sum_{n=1}^N \bar{q}_j^n - \sum_{n=1}^N (q_i^n * \bar{q}_j^n)}$$

$$C_{i,j}^{ces} = -\frac{1}{N} \sum_{n=1}^N [\bar{q}_j^n \log q_i^n + (1 - \bar{q}_j^n) \log(1 - q_i^n)]$$

---

$$C_{i,j} = C_{i,j}^{ed} + C_{i,j}^{sIoU} + C_{i,j}^{ces}$$

# Our Method (3D-BoNet):



## End-to-end training losses

$$l_{all} = l_{sem} + l_{bbox} + l_{bbs} + l_{pmask}$$

# Quantitative Results:

## ScanNet Benchmark

Benchmarks ▾ Documentation About Submit

Metric: AP 50% ▾

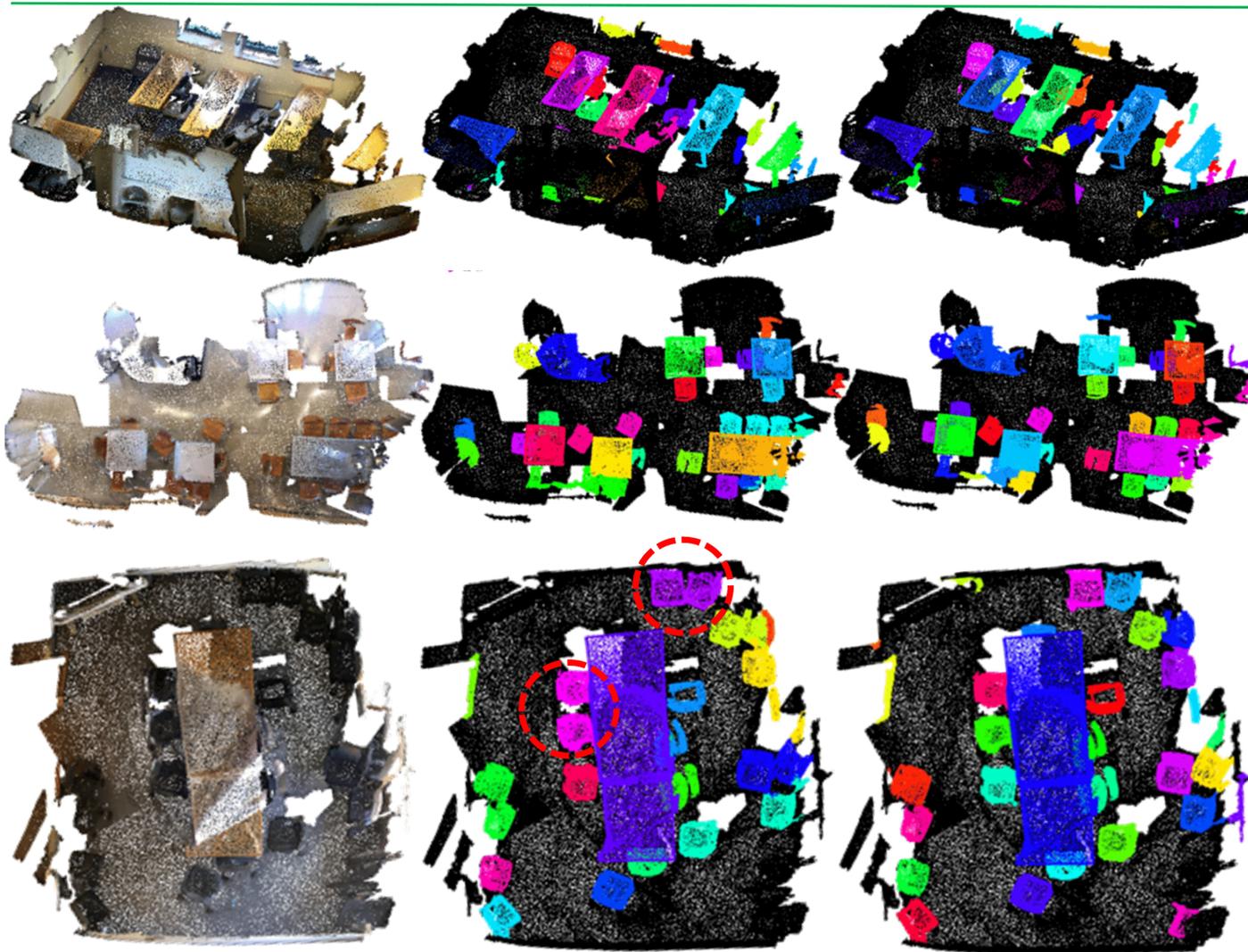
Method	Info	avg ap 50%	bathtub	bed	bookshelf	cabinet	chair	counter	curtain	desk	door	otherfurniture	picture	refrigerator	shower curtain	sink
<a href="#">3D-BoNet</a>		0.488 <sup>1</sup>	1.000 <sup>1</sup>	0.672 <sup>4</sup>	0.590 <sup>2</sup>	0.301 <sup>3</sup>	0.484 <sup>7</sup>	0.098 <sup>2</sup>	0.620 <sup>1</sup>	0.306 <sup>1</sup>	0.341 <sup>4</sup>	0.259 <sup>6</sup>	0.125 <sup>5</sup>	0.434 <sup>2</sup>	0.796 <sup>4</sup>	0.402 <sup>3</sup>
<a href="#">MTML</a>		0.481 <sup>2</sup>	1.000 <sup>1</sup>	0.666 <sup>5</sup>	0.377 <sup>4</sup>	0.272 <sup>4</sup>	0.709 <sup>1</sup>	0.001 <sup>11</sup>	0.579 <sup>3</sup>	0.254 <sup>3</sup>	0.361 <sup>3</sup>	0.318 <sup>4</sup>	0.095 <sup>7</sup>	0.432 <sup>3</sup>	1.000 <sup>1</sup>	0.184 <sup>6</sup>
<a href="#">PanopticFusion-inst</a>		0.478 <sup>3</sup>	0.667 <sup>5</sup>	0.712 <sup>3</sup>	0.595 <sup>1</sup>	0.259 <sup>6</sup>	0.550 <sup>6</sup>	0.000 <sup>12</sup>	0.613 <sup>2</sup>	0.175 <sup>5</sup>	0.250 <sup>7</sup>	0.434 <sup>1</sup>	0.437 <sup>1</sup>	0.411 <sup>5</sup>	0.857 <sup>2</sup>	0.485 <sup>1</sup>
Gaku Narita, Takashi Seno, Tomoya Ishikawa, Yohsuke Kaji: PanopticFusion: Online Volumetric Semantic Mapping at the Level of Stuff and Things. arXiv																
<a href="#">ResNet-backbone</a>		0.459 <sup>4</sup>	1.000 <sup>1</sup>	0.737 <sup>1</sup>	0.159 <sup>10</sup>	0.259 <sup>5</sup>	0.587 <sup>4</sup>	0.138 <sup>1</sup>	0.475 <sup>5</sup>	0.217 <sup>4</sup>	0.416 <sup>1</sup>	0.408 <sup>3</sup>	0.128 <sup>4</sup>	0.315 <sup>6</sup>	0.714 <sup>5</sup>	0.411 <sup>2</sup>
<a href="#">MASC</a>	<span>P</span>	0.447 <sup>5</sup>	0.528 <sup>8</sup>	0.555 <sup>7</sup>	0.381 <sup>3</sup>	0.382 <sup>1</sup>	0.633 <sup>2</sup>	0.002 <sup>9</sup>	0.509 <sup>4</sup>	0.260 <sup>2</sup>	0.361 <sup>2</sup>	0.432 <sup>2</sup>	0.327 <sup>2</sup>	0.451 <sup>1</sup>	0.571 <sup>6</sup>	0.367 <sup>4</sup>
Chen Liu, Yasutaka Furukawa: MASC: Multi-scale Affinity with Sparse Convolution for 3D Instance Segmentation.																
<a href="#">3D-SIS</a>	<span>P</span>	0.382 <sup>6</sup>	1.000 <sup>1</sup>	0.432 <sup>8</sup>	0.245 <sup>7</sup>	0.190 <sup>7</sup>	0.577 <sup>5</sup>	0.013 <sup>7</sup>	0.263 <sup>7</sup>	0.033 <sup>10</sup>	0.320 <sup>5</sup>	0.240 <sup>7</sup>	0.075 <sup>8</sup>	0.422 <sup>4</sup>	0.857 <sup>2</sup>	0.117 <sup>9</sup>
Ji Hou, Angela Dai, Matthias Niessner: 3D-SIS: 3D Semantic Instance Segmentation of RGB-D Scans. CVPR 2019																
<a href="#">UNet-backbone</a>		0.319 <sup>7</sup>	0.667 <sup>5</sup>	0.715 <sup>2</sup>	0.233 <sup>8</sup>	0.189 <sup>8</sup>	0.479 <sup>8</sup>	0.008 <sup>8</sup>	0.218 <sup>8</sup>	0.067 <sup>9</sup>	0.201 <sup>8</sup>	0.173 <sup>8</sup>	0.107 <sup>6</sup>	0.123 <sup>8</sup>	0.438 <sup>7</sup>	0.150 <sup>7</sup>
<a href="#">R-PointNet</a>		0.306 <sup>8</sup>	0.500 <sup>9</sup>	0.405 <sup>9</sup>	0.311 <sup>5</sup>	0.348 <sup>2</sup>	0.589 <sup>3</sup>	0.054 <sup>3</sup>	0.068 <sup>10</sup>	0.126 <sup>6</sup>	0.283 <sup>6</sup>	0.290 <sup>5</sup>	0.028 <sup>9</sup>	0.219 <sup>7</sup>	0.214 <sup>10</sup>	0.331 <sup>5</sup>
<a href="#">3D-BEVIS</a>		0.248 <sup>9</sup>	0.667 <sup>5</sup>	0.566 <sup>6</sup>	0.076 <sup>11</sup>	0.035 <sup>12</sup>	0.394 <sup>9</sup>	0.027 <sup>5</sup>	0.035 <sup>11</sup>	0.098 <sup>7</sup>	0.099 <sup>10</sup>	0.030 <sup>11</sup>	0.025 <sup>10</sup>	0.098 <sup>9</sup>	0.375 <sup>8</sup>	0.126 <sup>8</sup>
Cathrin Elich, Francis Engelmann, Jonas Schult, Theodora Kontogianni, Bastian Leibe: 3D-BEVIS: Birds-Eye-View Instance Segmentation.																
<a href="#">Seg-Cluster</a>	<span>P</span>	0.215 <sup>10</sup>	0.370 <sup>10</sup>	0.337 <sup>11</sup>	0.285 <sup>6</sup>	0.105 <sup>9</sup>	0.325 <sup>10</sup>	0.025 <sup>6</sup>	0.282 <sup>6</sup>	0.085 <sup>8</sup>	0.105 <sup>9</sup>	0.107 <sup>9</sup>	0.007 <sup>12</sup>	0.079 <sup>10</sup>	0.317 <sup>9</sup>	0.114 <sup>10</sup>

# Qualitative Results:

Input Point Clouds

Predicted Instance Labels

Ground Truth



# Intermediate Results:

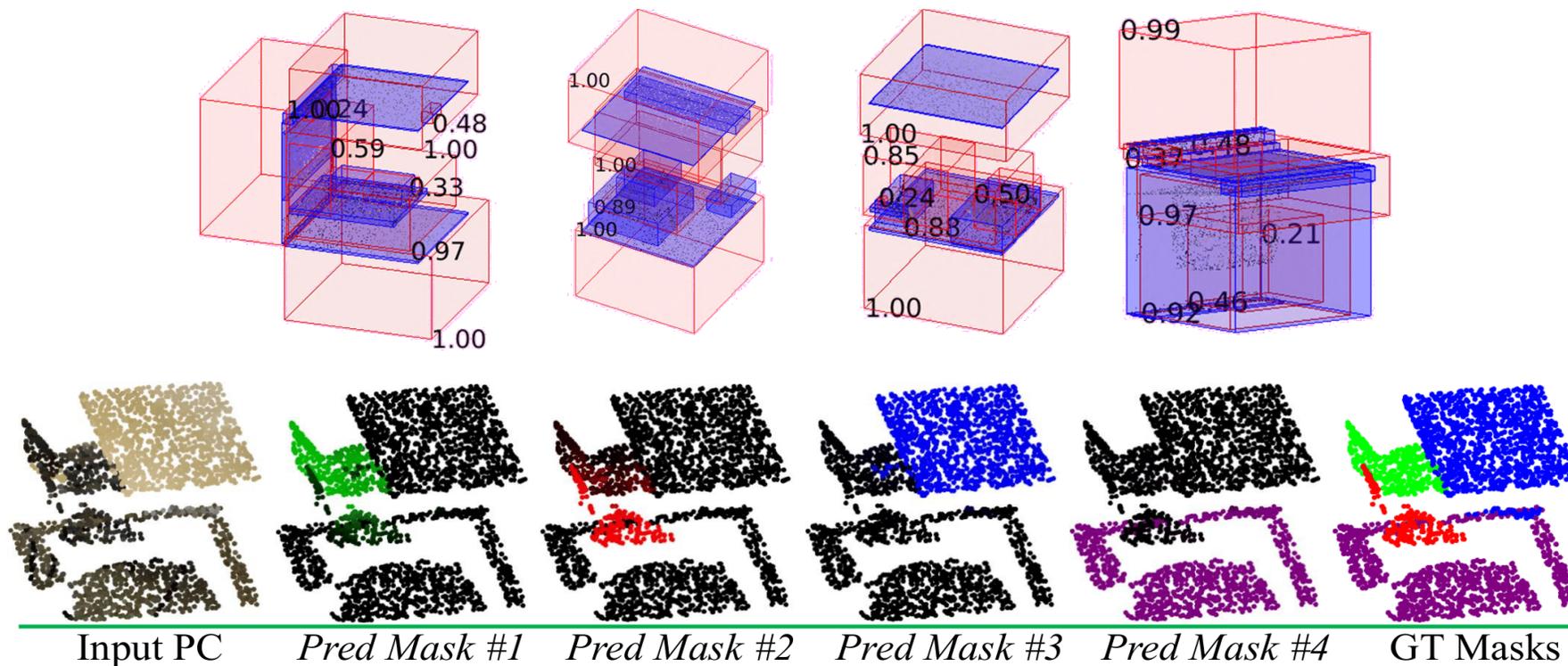


Table 4: Time consumption of different approaches on the validation split (312 scenes) of ScanNet(v2) (seconds).

	SGPN [50]	ASIS [51]	GSPN [58]	3D-SIS [15]	3D-BoNet(Ours)
	network(GPU): 650 group merging(CPU): 46562 block merging(CPU): 2221	network(GPU): 650 mean shift(CPU): 53886 block merging(CPU): 2221	network(GPU): 500 point sampling(GPU): 2995 neighbour search(CPU): 468	voxelization, projection, network, etc. (GPU+CPU): 38841	network(GPU): 650 <i>SCN (GPU parallel): 208</i> block merging(CPU): 2221
total	49433	56757	3963	38841	<b>2871</b>

**Thank You**