

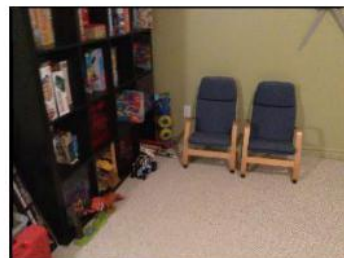
Deep Learning for 3D Point Cloud Analysis



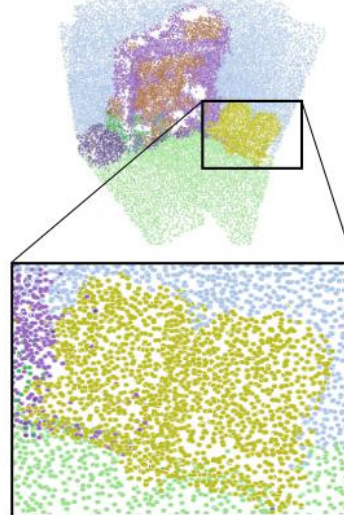
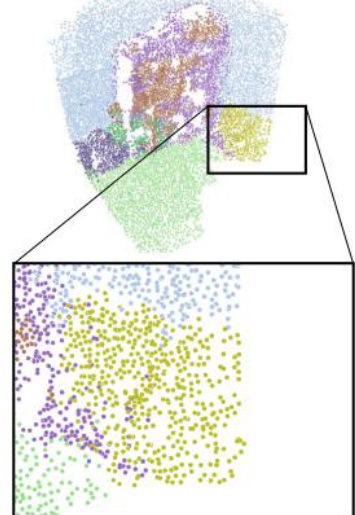
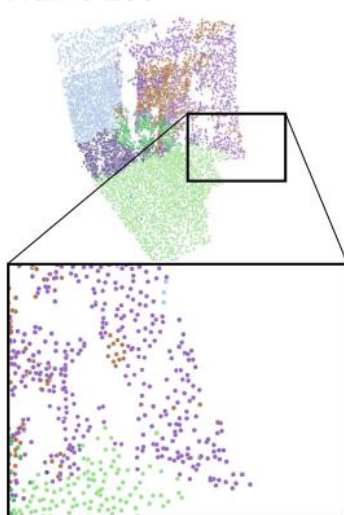
Frame 180



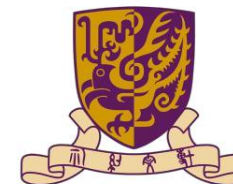
Frame 260



Frame 350



Xu Yan
2020.6.26



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Paper List

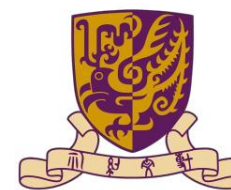


- Multi-Path Region Mining For Weakly Supervised 3D Semantic Segmentation on Point Clouds (CVPR2020)
- SESS: Self-Ensembling Semi-Supervised 3D Object Detection (CVPR2020 **oral**)
- Fusion-Aware Point Convolution for Online Semantic 3D Scene Segmentation (CVPR2020 **best paper final list**)

Paper List



- Multi-Path Region Mining For Weakly Supervised 3D Semantic Segmentation on Point Clouds (CVPR2020)
- SESS: Self-Ensembling Semi-Supervised 3D Object Detection (CVPR2020 oral)
- Fusion-Aware Point Convolution for Online Semantic 3D Scene Segmentation (CVPR2020 best paper final list)



Multi-Path Region Mining For Weakly Supervised 3D Semantic Segmentation on Point Clouds

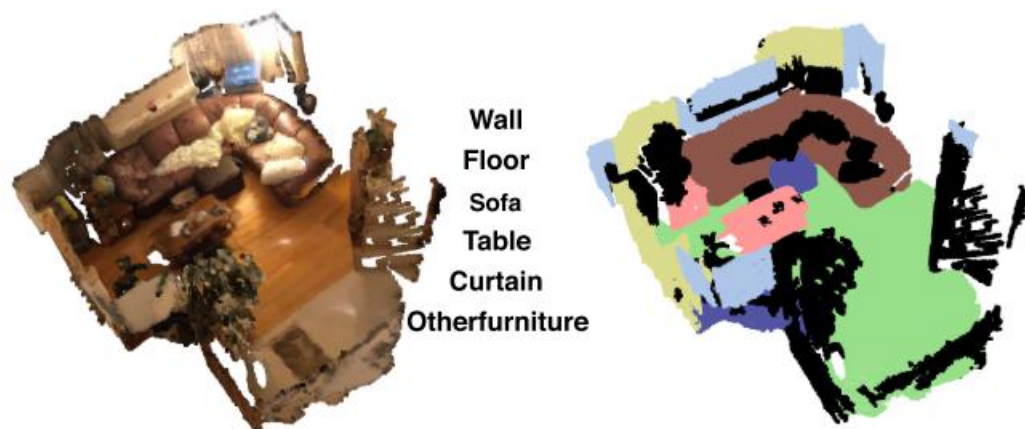
Jiacheng Wei¹ Guosheng Lin^{1*} Kim-Hui Yap¹ Tzu-Yi Hung² Lihua Xie¹

¹Nanyang Technological University, Singapore ²Delta Research Center, Singapore

{jiacheng002, gslin, ekhyap, elhxie}@ntu.edu.sg, tzuyi.hung@deltaww.com

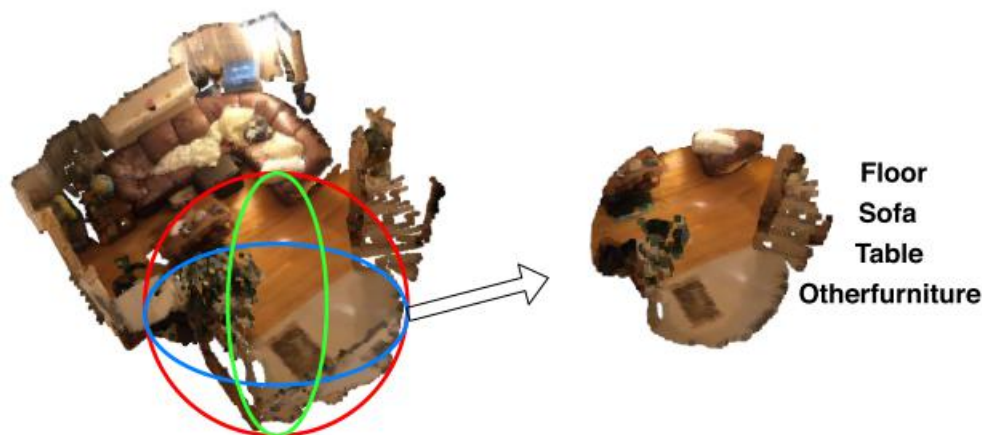
For *ScanNet* dataset

- Only **20 people** participated in the collection of 1513 3D scans.
- More than **500 workers** participated in the semantic annotation process.
- To ensure the annotation accuracy, each scene was annotated by 2 to 3 participants.
- The median and mean time for annotation per scan is **16.8 min** and **22.3 min**.
- The estimated annotation time for scene-level labels in a scene is **around 15 sec**, while the annotation time for subclouds from a scene is lower than **3 min**. The average number of subclouds is **18.4 (save 6-70 times the time)**.



(A) Scene-level Label

(B) Point-level Label



(C) Subcloud-level Label

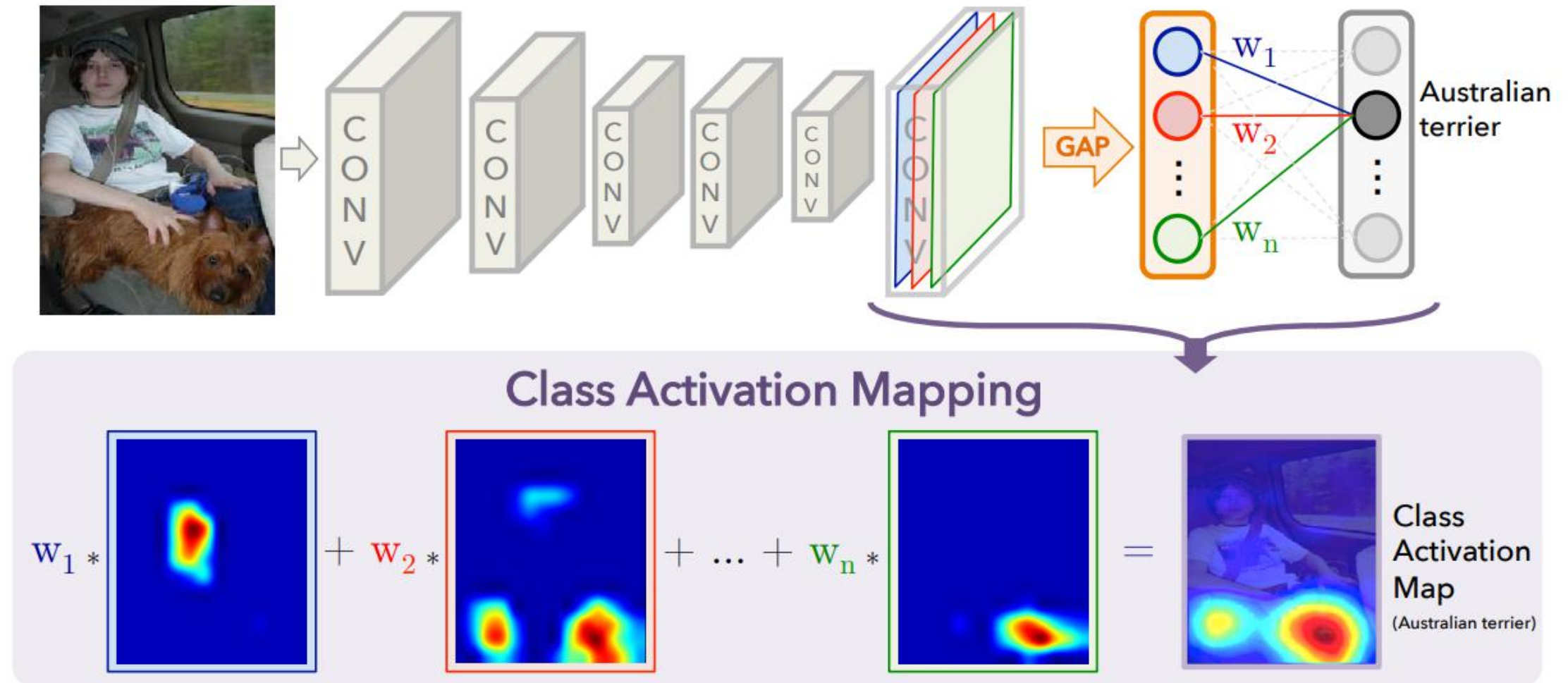
MPRM

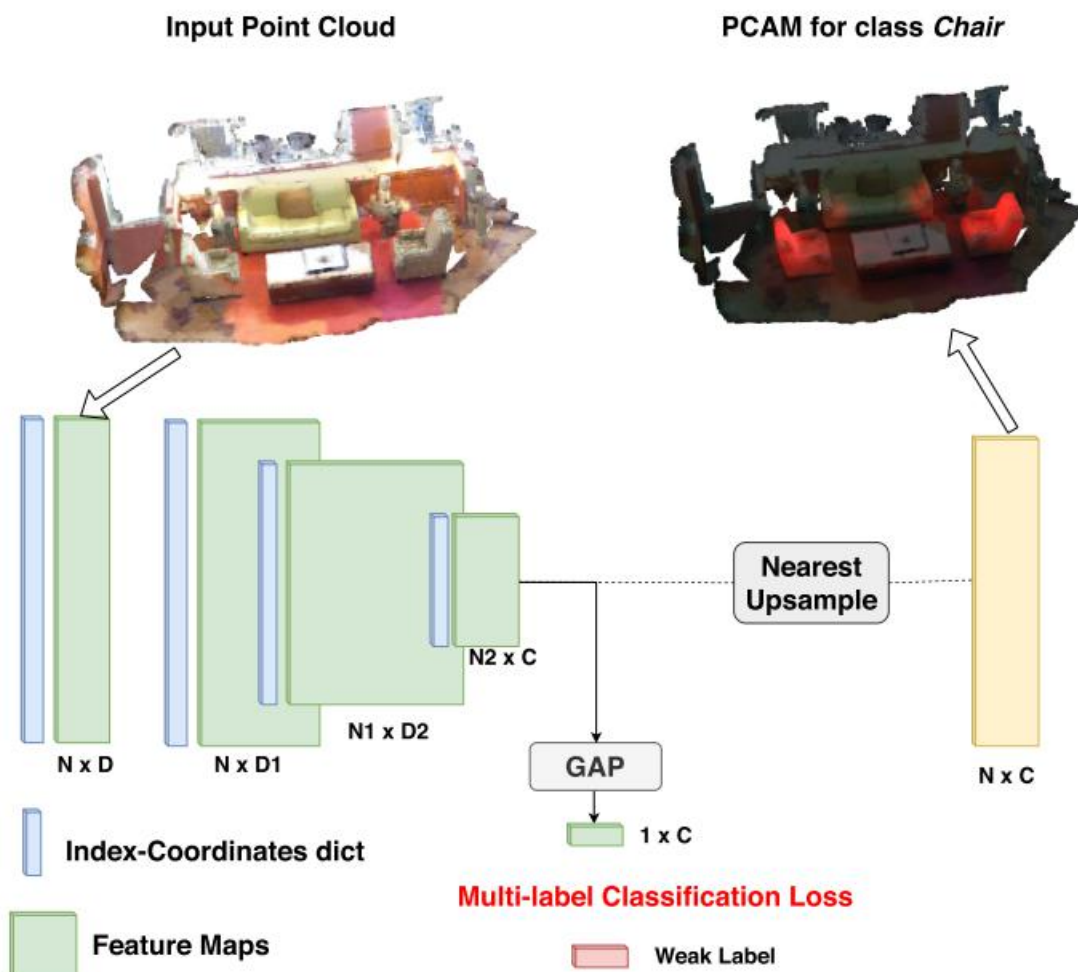


香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Class activation map (**CAM**) in 2D image:





- Point class activation map (PCAM):

We denote $f_{cam}(p)$ as the PCAM feature vector of point p before the global average pooling layer. For class c , the PCAM $M_c(p)$ for point p can be expressed as:

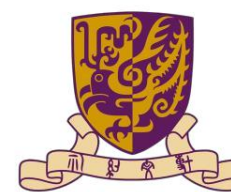
$$M_c(p) = \mathbf{w}_c^T \cdot f_{cam}(p) \cdot \mathbf{y}_c, \quad (2)$$

w_c is the classification weights for class c and $y_c \in \{0, 1\}$ indicates the one-hot subcloud ground truth for class c .

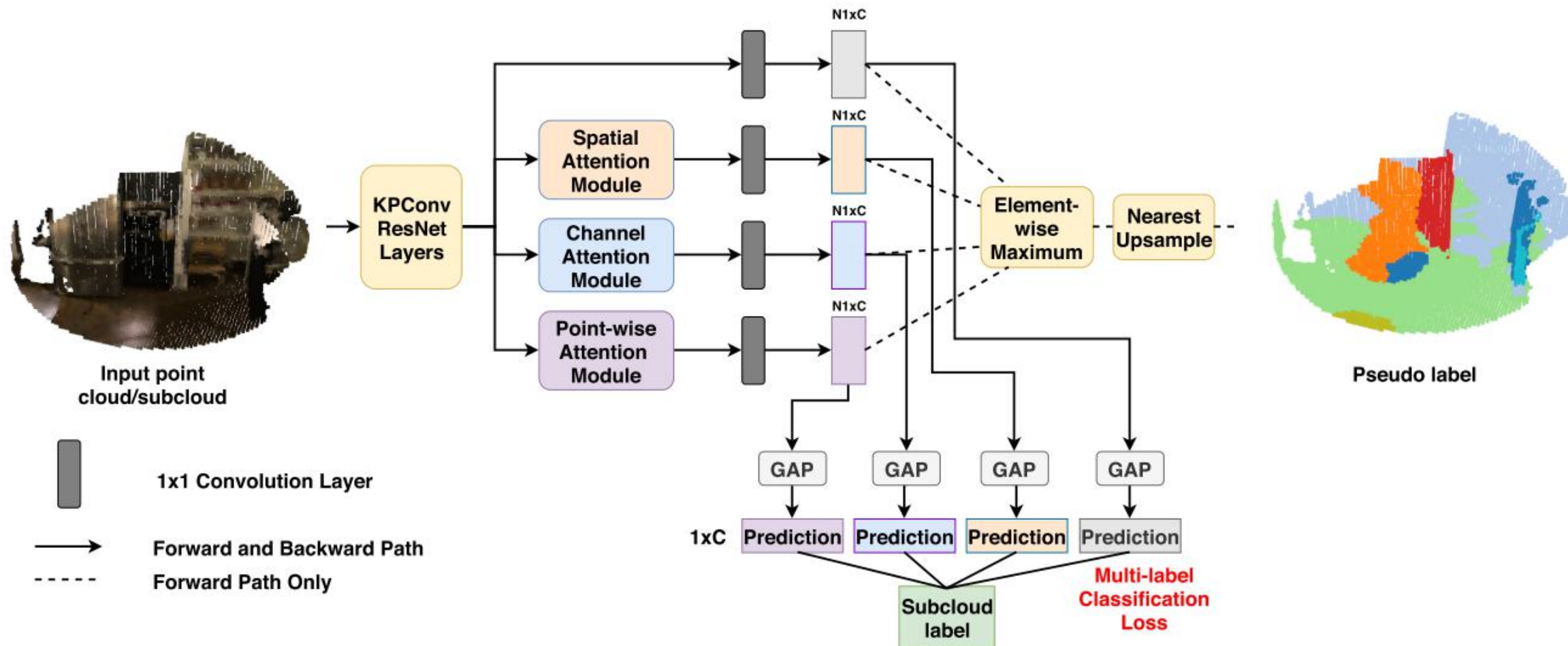
- Point-level pseudo masks:

$$\operatorname{argmax}(M(p))$$

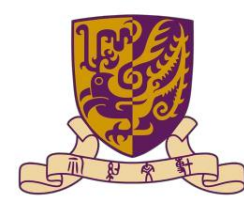
MPRM



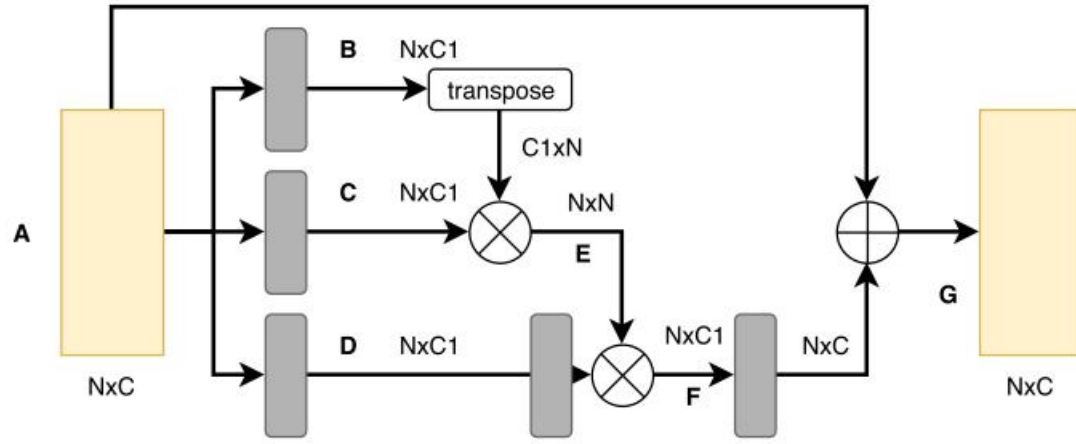
香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



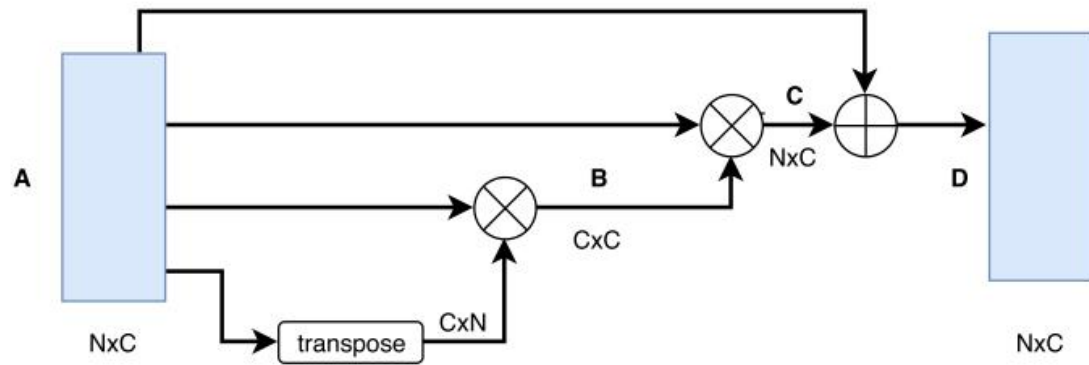
MPRM



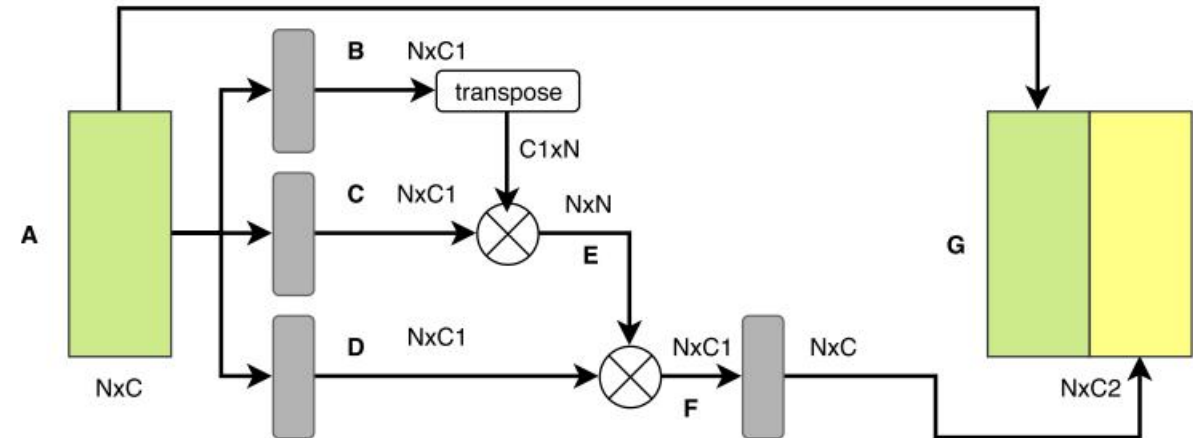
香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



A.Spatial Attention Module



B.Channel Attention Module



C.Point-wise Attention Module

Setting	Supervision	wall	floor	cabinet	bed	chair	sofa	table	door	window	B.S.	picture	cnt	desk	curtain	fridge	S.C.	toilet	sink	bathtub	other	mIoU
PCAM(Baseline)	Scene	54.9	48.3	14.1	34.7	32.9	45.3	26.1	0.6	3.3	46.5	0.6	6.0	7.4	26.9	0.0	6.1	22.3	8.2	52.0	6.1	22.1
MPRM(Ours)	Scene	47.3	41.1	10.4	43.2	25.2	43.1	21.5	9.8	12.3	45.0	9.0	13.9	21.1	40.9	1.8	29.4	14.3	9.2	39.9	10.0	24.4
PCAM(Baseline)	Subcloud	59.0	53.8	24.7	64.9	45.7	60.7	42.8	31.5	37.0	55.9	31.0	12.0	39.1	68.7	16.8	49.8	55.2	27.4	59.0	27.7	43.1
MPRM(Ours)	Subcloud	56.1	54.8	32.0	69.6	49.5	67.7	46.6	41.3	44.2	71.5	28.3	21.3	49.2	71.8	38.1	42.8	43.6	20.3	49.0	33.8	46.6

dCRF post-processing:

MPRM(Ours)	Subcloud	58.0	57.3	33.2	71.8	50.4	69.8	47.9	42.1	44.9	73.8	28.0	21.5	49.5	72.0	38.8	44.1	42.4	20.0	48.7	34.4	47.4
------------	----------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Table 2. The class-specific segmentation results (mIoU) of pseudo labels on training set generated with different settings and different supervision levels. We only show the dCRF post-processed result for MPRM with subcloud-level supervision since we use this pseudo label to train our final segmentation model. (Here B.S. stands for bookshelf; S.C. stands for shower curtain; cnt stands for counter.)

Fusion	PCAM	SA	CA	PSA	Training	Validation
-	✓				44.3	39.3
-		✓			44.8	39.4
-			✓		44.3	39.3
-				✓	44.7	39.5
Max	✓	✓			46.0	40.3
Max	✓		✓		45.9	40.0
Max	✓			✓	45.6	40.4
Max	✓	✓	✓	✓	46.6	41.0
Sum	✓	✓	✓	✓	45.9	39.7

Table 3. The mIoU of pseudo labels with different paths and their combinations on training and validation set.

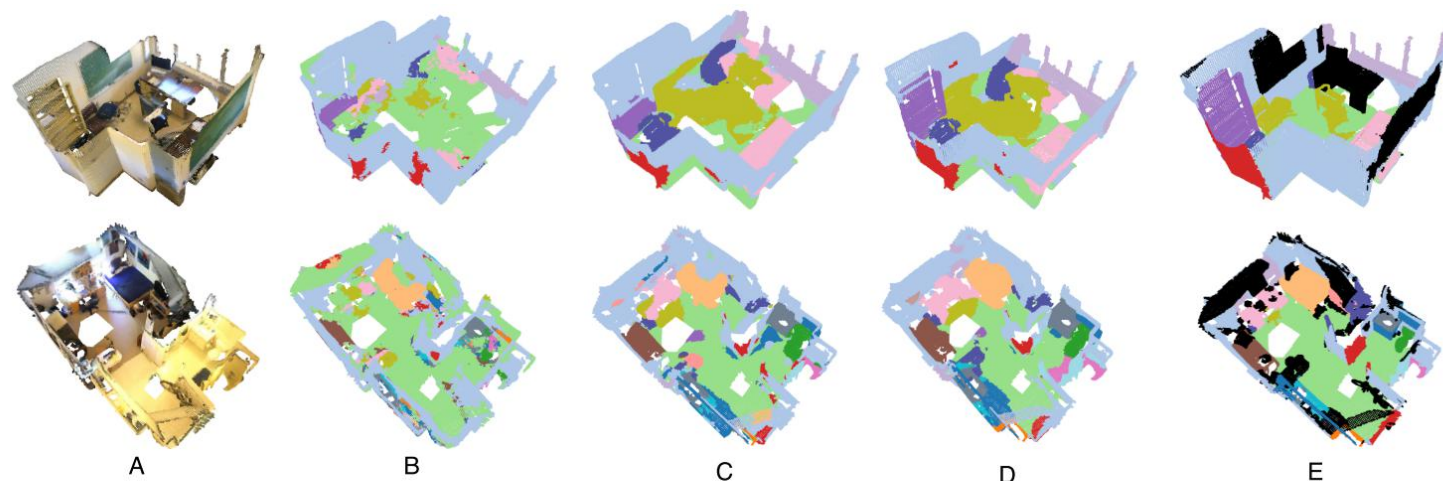
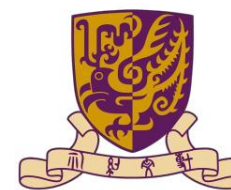


Figure 5. Visualizations of pseudo labels. (A)Input point clouds, (B)PCAMs trained with scene-level labels, (C)PCAMs trained with subcloud-level labels, (D)Multi-path region mining trained with subcloud-level labels, (E)Ground truth.

Paper List



- Multi-Path Region Mining For Weakly Supervised 3D Semantic Segmentation on Point Clouds (CVPR2020)
- **SESS: Self-Ensembling Semi-Supervised 3D Object Detection (CVPR2020 oral)**
- Fusion-Aware Point Convolution for Online Semantic 3D Scene Segmentation (CVPR2020 best paper final list)



SESS: Self-Ensembling Semi-Supervised 3D Object Detection

Na Zhao Tat-Seng Chua Gim Hee Lee

Department of Computer Science, National University of Singapore

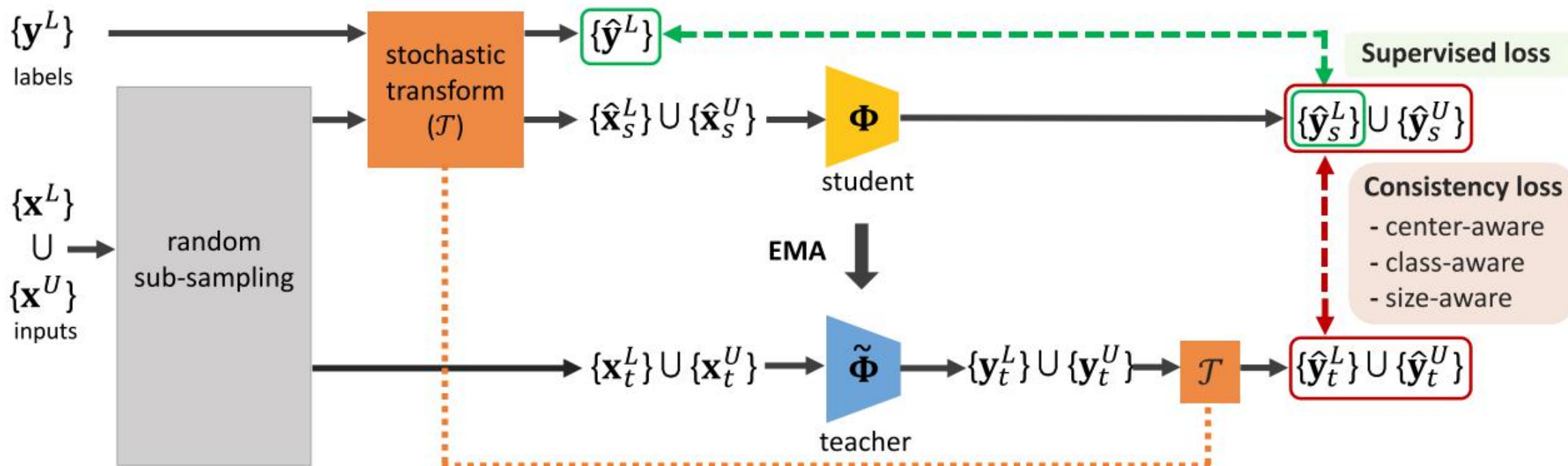
{nazhao, chuats, gimhee.lee}@comp.nus.edu.sg

Semi-supervised learning:

- **Self-training:** Pseudo-Label : The Simple and Efficient Semi-Supervised Learning Method for Deep Neural Networks (ICML 2013)

$$L = \sum_{m=1}^n \sum_{i=1}^C L(y_i^m, f_i^m) + \alpha(t) \sum_{m=1}^{n'} \sum_{i=1}^C L(y_i'^m, f_i'^m)$$

- **Temporal ensembling:** Temporal Ensembling for Semi-supervised Learning (ICLR 2017)
- **Self ensembling:** Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results (NIPS 2017)



- **Perturbation Scheme:**

Random Sub-sampling & Stochastic Transform (flipping, rotation and scaling)

- **Consistency Loss:**

Center-aware consistency loss:
$$\mathcal{L}_{center} = \frac{\sum_{\hat{c}_s} \|\hat{c}_s - \hat{c}_t^A\|_2 + \sum_{\hat{c}_t} \|\hat{c}_t - \hat{c}_s^A\|_2}{|\hat{C}_s| + |\hat{C}_t|}, \quad (5)$$

Class-aware consistency loss:
$$\mathcal{L}_{class} = \frac{1}{|\hat{P}_t|} \sum D_{KL}(\hat{p}_s^A \parallel \hat{p}_t). \quad (6)$$

Size-aware consistency loss:
$$\mathcal{L}_{size} = \frac{1}{|\hat{D}_t|} \sum (\hat{d}_s^A - \hat{d}_t)^2. \quad (7)$$

- **Implementation Details**

(1) VoteNet as backbone

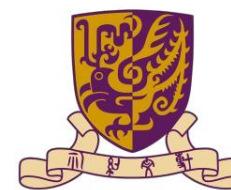
(2) Randomly sub-sampled points is 4,000

(3) Ramp up the coefficient of consistency cost from 0 to its maximum value of 10 during the first 30 epochs.

(4) **Pre-train** VoteNet with all the available labeled samples. Then initialize the student and teacher networks with the pre-trained weights, and train the student network on both the labeled and unlabeled data by minimizing the supervised loss as well as consistency loss.

Table 1: Comparison with VoteNet on SUN RGB-D val set and ScanNetV2 val set with varying ratios of labeled data. mAP@0.25 are reported as mean \pm standard deviation, based on 3 runs with random sampling. And the improvement (Improv.) is computed based on the mean performances over 3 runs. Note that our SESS is initialized by the VoteNet weights pre-trained on the corresponding labeled data.

Dataset	Model	10%	20%	30%	40%	50%	70%	100%
SUNRGB-D	VoteNet [11]	34.43 \pm 1.07	41.13 \pm 0.36	47.70 \pm 0.17	50.77 \pm 0.19	52.5 \pm 0.19	56.13 \pm 0.18	57.7
	SESS	42.87\pm1.01	47.87\pm0.48	53.17\pm0.63	54.73\pm0.26	56.37\pm0.22	58.97\pm0.17	61.1
	Improv.(%)	24.51 \uparrow	16.39 \uparrow	11.47 \uparrow	7.80 \uparrow	7.37 \uparrow	5.06 \uparrow	5.89 \uparrow
ScanNetV2	VoteNet [11]	30.97 \pm 0.79	41.60 \pm 0.46	45.57 \pm 0.38	49.2 \pm 0.33	52.57 \pm 0.07	54.97 \pm 0.07	58.6
	SESS	39.67\pm0.91	47.93\pm0.39	52.20\pm0.09	54.93\pm0.27	57.77\pm0.41	59.20\pm0.08	62.1
	Improv.(%)	28.09 \uparrow	15.22 \uparrow	14.55 \uparrow	11.64 \uparrow	9.89 \uparrow	7.70 \uparrow	5.97 \uparrow



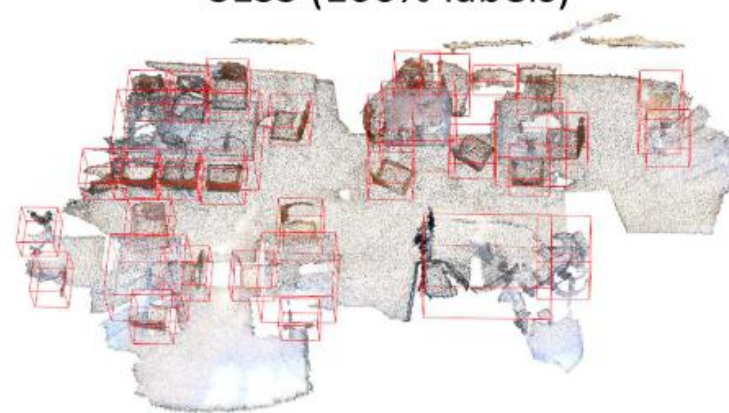
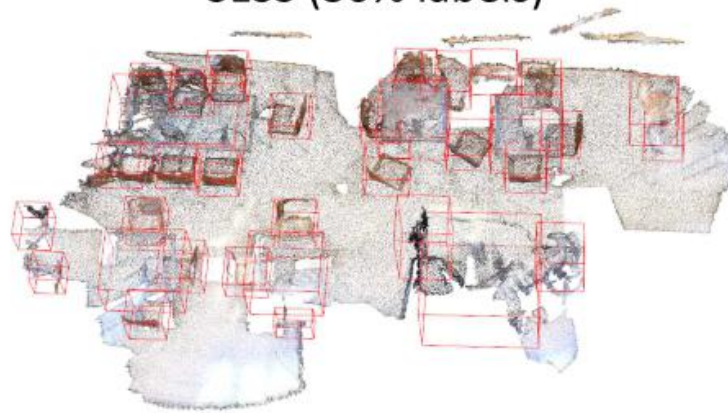
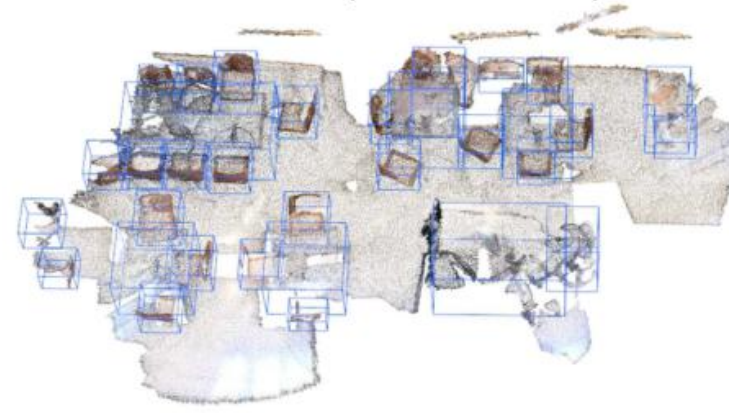
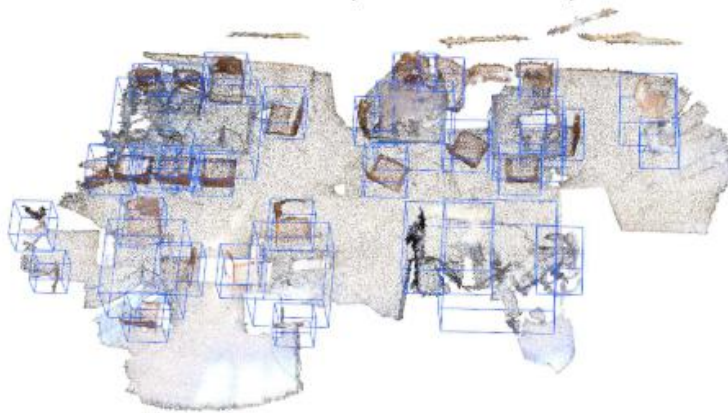
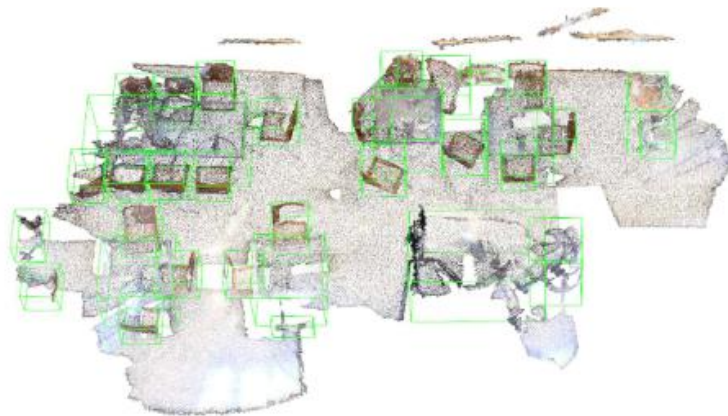
Ground truth

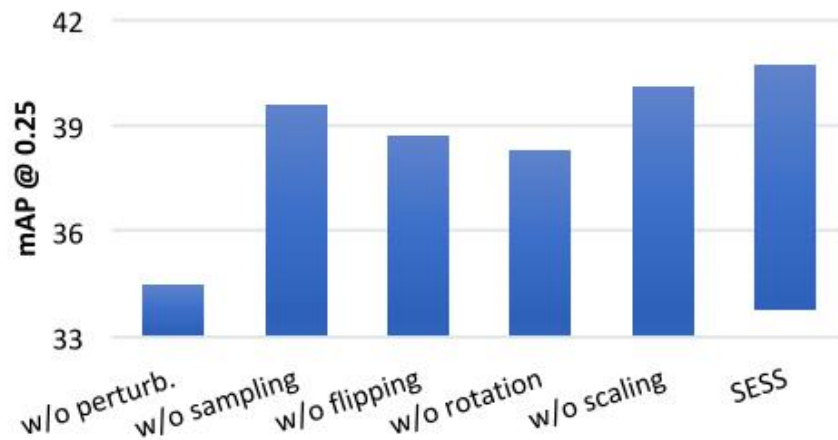
VoteNet (30% labels)

VoteNet (100% labels)

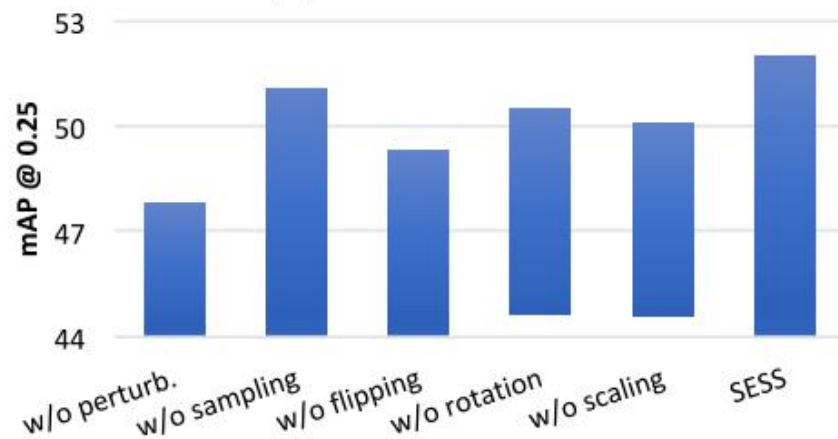
SESS (30% labels)

SESS (100% labels)





(a) SUN RGB-D



(b) ScanNetV2

Table 4: Ablation study on consistency losses.

center	class	size	SUN RGB-D	ScanNetV2
✓	✗	✗	38.2	50.0
✗	✓	✗	39.2	50.2
✗	✗	✓	38.1	49.2
✗	✓	✓	40.3	50.7
✓	✗	✓	38.9	50.5
✓	✓	✗	40.0	51.5
✓	✓	✓	40.7	52.0

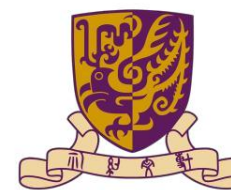
The training of ablation experiments is conducted on **SUN RGB-D with 10% labeled data** and **ScanNetV2 with of 30% labeled data**.

Figure 3: Effects of different perturbations.

Paper List



- Multi-Path Region Mining For Weakly Supervised 3D Semantic Segmentation on Point Clouds (CVPR2020)
- SESS: Self-Ensembling Semi-Supervised 3D Object Detection (CVPR2020 oral)
- Fusion-Aware Point Convolution for Online Semantic 3D Scene Segmentation (CVPR2020 best paper final list)

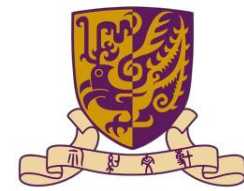
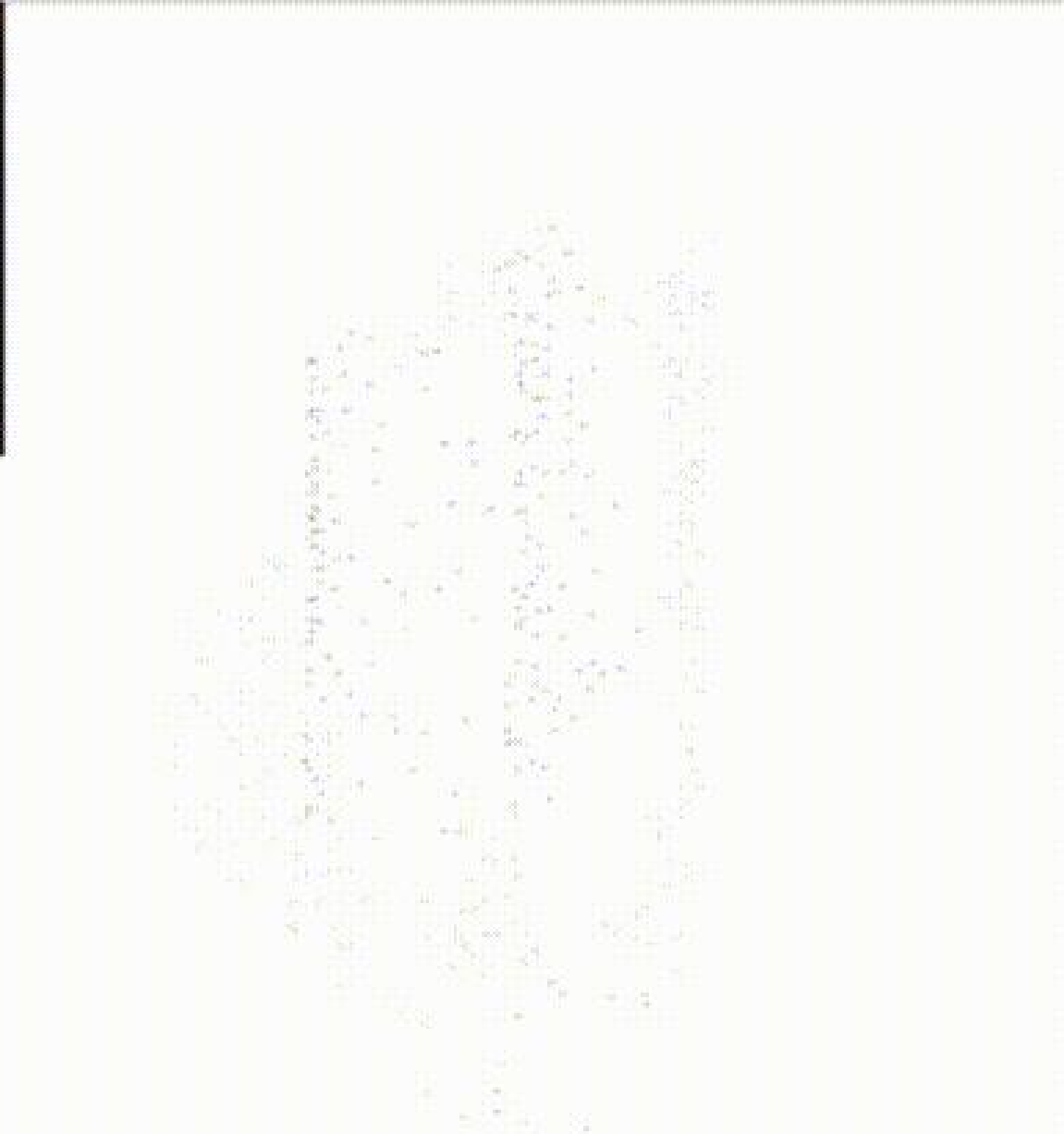


Fusion-Aware Point Convolution for Online Semantic 3D Scene Segmentation

Jiazhao Zhang^{1,*} Chenyang Zhu^{1,*} Lintao Zheng¹ Kai Xu^{1,2†}

¹National University of Defense Technology ²SpeedBot Robotics Ltd.

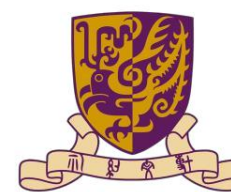
FusionAwareConv



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



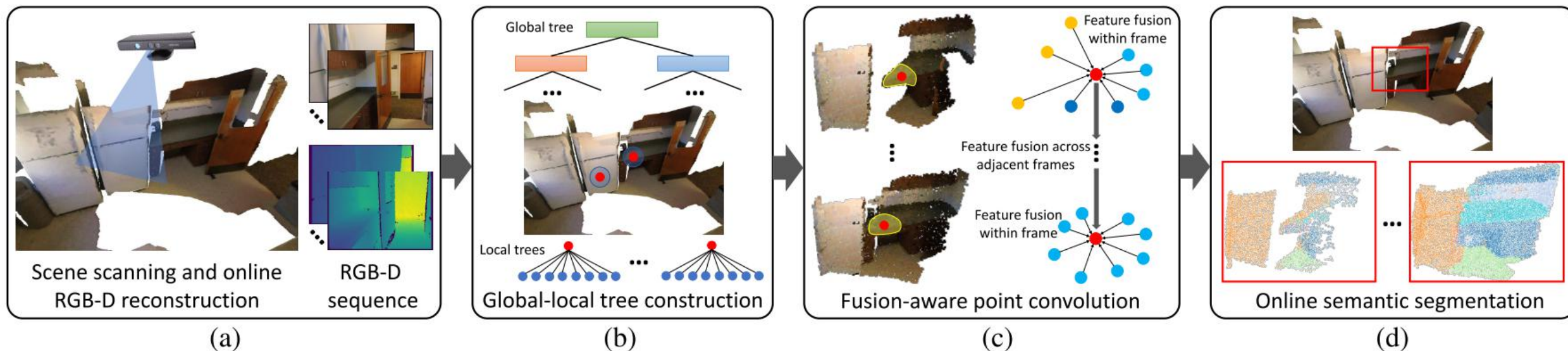
FusionAwareConv



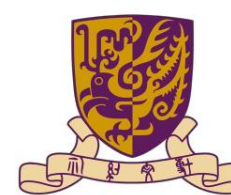
香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Pipeline



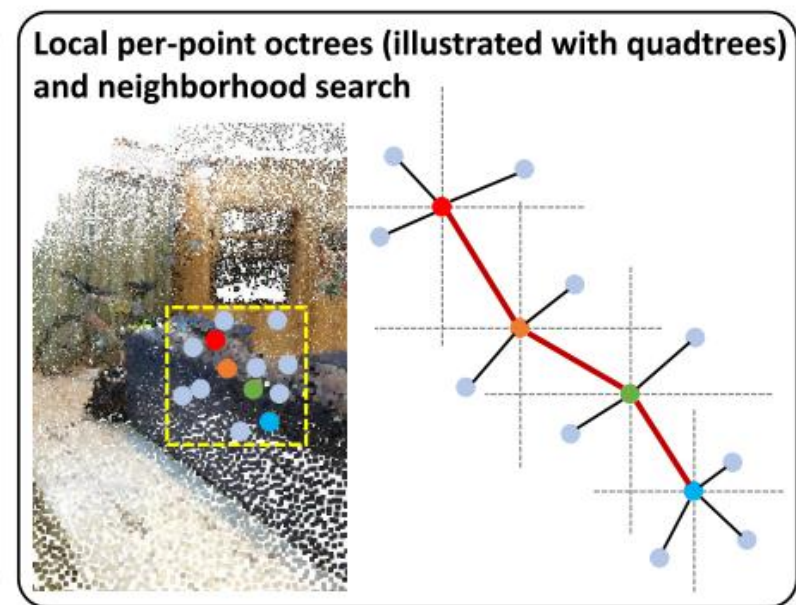
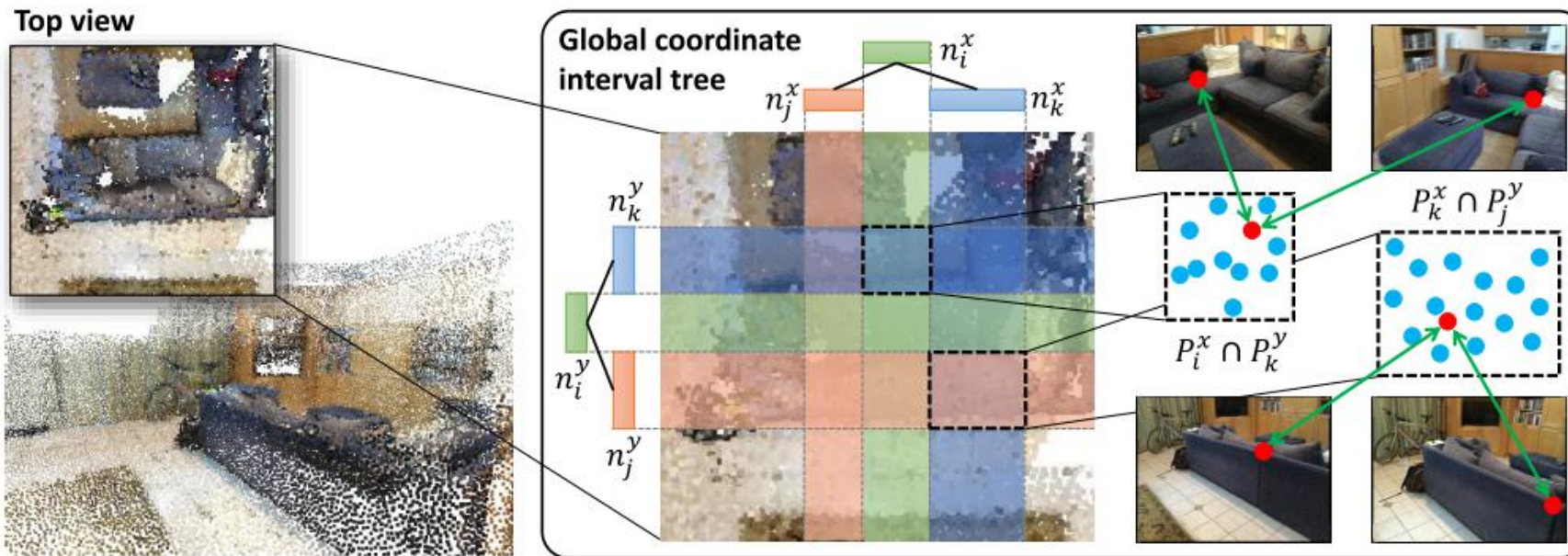
FusionAwareConv



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Global-local Tree



- Global Tree: **Coordinate interval trees** $x_{\max}(n_l) < x_{\min}(n_p)$, $x_{\max}(n_p) < x_{\min}(n_r)$,
- Local Tree: **Octree**

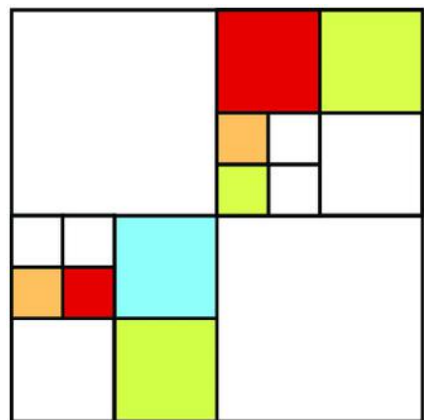
FusionAwareConv



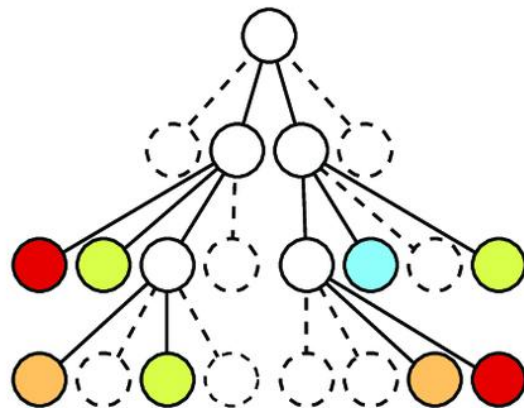
香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



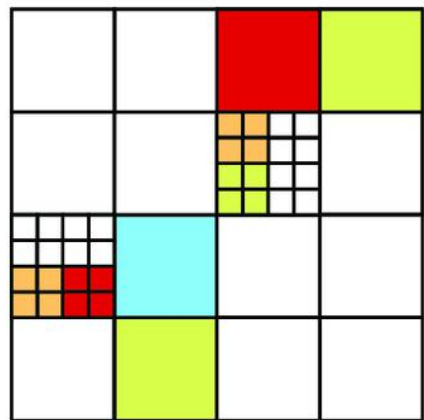
Octree



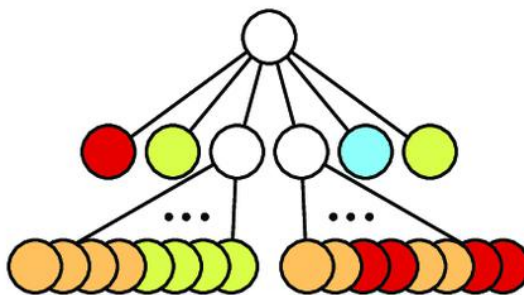
(a)



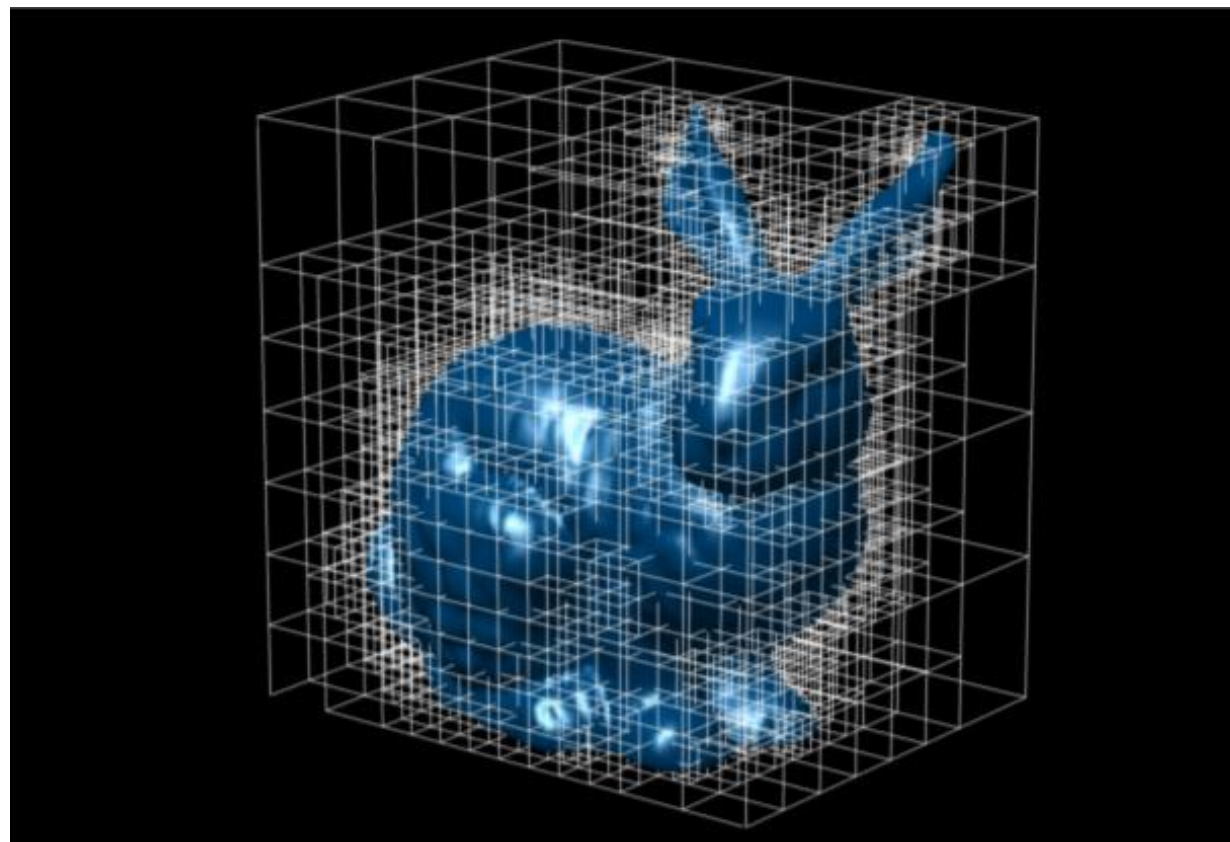
(b)



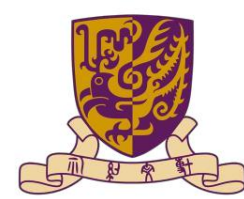
(c)



(d)



FusionAwareConv



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Fusion-aware Point Convolution

- 3D feature fusion: **PointConv**:

$$\text{PC}_p(W, F) = \sum_{\Delta p \in \Omega} W(\Delta p) F(p + \Delta p),$$

- 2D-3D feature fusion: **Pre-trained FuseNet**

$$F^{2D}(c^k) = \text{FuseNet}(f_k, c^k),$$

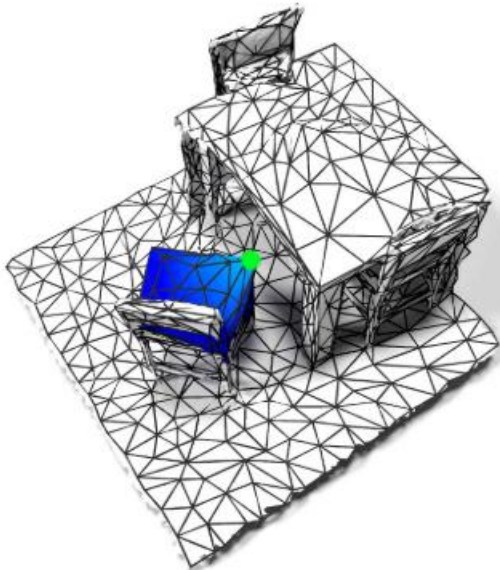
$$F^{2D3D}(p) = \text{maxpooling}\{F^{2D}(c^k) | c^k \in I(p)\}$$

Fusion-aware Point Convolution

- Octree-induced surface-aware 3D convolution. (Geodesic Neighbors)

$$\text{FPC}_p(W, F^{2\text{D}3\text{D}}) = \sum_{\Delta p \in \Omega^n(p)} W(\Delta p) F^{2\text{D}3\text{D}}(p + \Delta p).$$

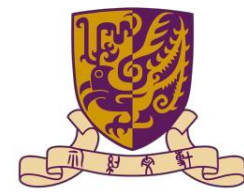
Geodesic Neighborhood



Euclidean Neighborhood



FusionAwareConv



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen

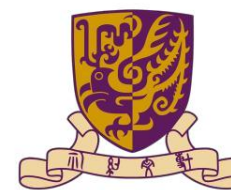


Frame-to-frame feature fusion.

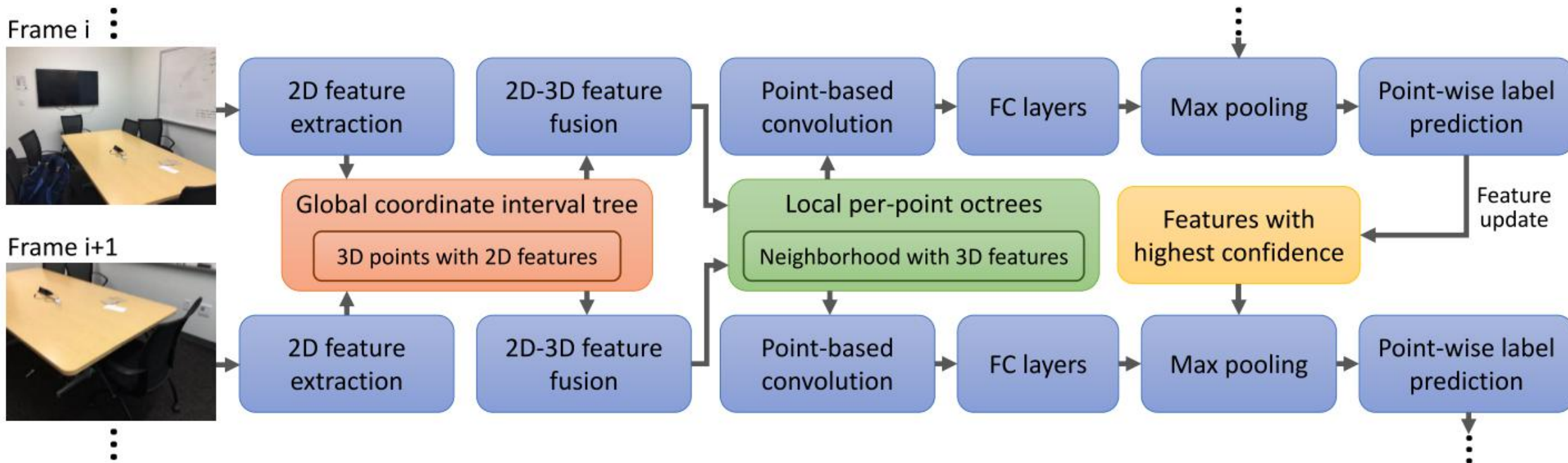
- If prediction has low segmentation uncertainty in the last frame, the current form of feature fusion should be useful in the future prediction.
- Record every uncertainty $U(p, i)$ when processing the frame sequence.

$$F^{\text{fused}}(p) = \text{maxpool}\{\text{FPC}_p^{\text{current}}(W, F^{2\text{D}3\text{D}}), \arg \min_{\text{FPC}_p^i} U(p, i)(W, F^{2\text{D}3\text{D}})\}$$

FusionAwareConv



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



FusionAwareConv



Results

- Online methods

Table 1: Accuracy comparison between our method and two state-of-the-art online scene segmentation methods.

Dataset	SemanticFusion [20]	ProgressiveFusion [23]	Ours
ScanNet	0.518	0.566	0.764
SceneNN	0.628	0.666	0.675

ScanNet: use 1200 sequences for training and the rest 312 for testing

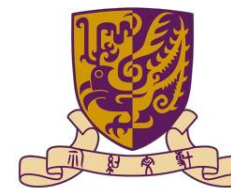
SceneNN: 15 sequences for evaluation

- Offline methods

Table 2: IOU comparison between our method and state-of-the-art offline scene segmentation methods. Our method has the highest mean IOU, outperforming the state-of-the-art methods for nine semantic categories.

model	mean	wall	floor	cabinet	bed	chair	sofa	table	door	window	bookshelf	picture	counter	desk	curtain	fridge	bathshade	toilet	sink	bathtub	others
SparseConvNet	0.685	0.828	0.950	0.620	0.805	0.894	0.825	0.707	0.633	0.588	0.788	0.252	0.601	0.592	0.681	0.428	0.607	0.928	0.596	0.881	0.504
MinkowskiNet	0.715	0.841	0.949	0.641	0.806	0.900	0.845	0.745	0.648	0.608	0.792	0.289	0.637	0.65	0.742	0.509	0.690	0.916	0.689	0.832	0.570
PointConv	0.580	0.741	0.948	0.474	0.672	0.813	0.633	0.651	0.346	0.446	0.713	0.067	0.568	0.525	0.551	0.370	0.520	0.840	0.590	0.750	0.387
Ours	0.720	0.862	0.924	0.615	0.848	0.716	0.804	0.637	0.680	0.698	0.724	0.513	0.617	0.588	0.764	0.734	0.696	0.870	0.681	0.885	0.556

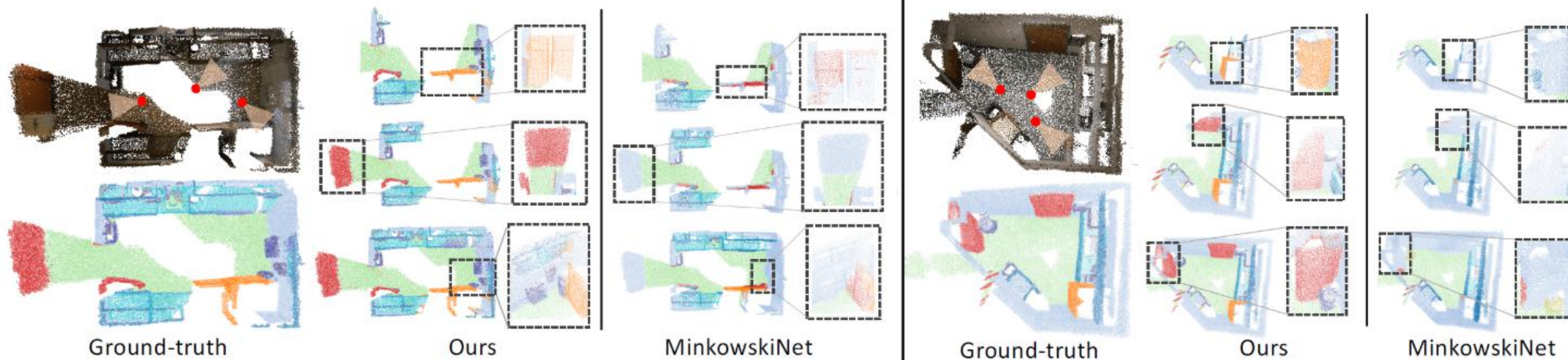
FusionAwareConv



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



- Visualization



- Ablation

2D3D feature	frame-to-frame feature	mean IOU
×	×	0.711
✓	×	0.718
×	✓	0.713
✓	✓	0.720



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Thanks for watching!

Xu Yan
2020.6.26