

## ORCA - Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository:https://orca.cardiff.ac.uk/id/eprint/149167/

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Zhou, Feng, Lai, Yu-Kun, Rosin, Paul L., Zhang, Fengquan and Hu, Yong 2022. Scale-aware network with modality-awareness for RGB-D indoor semantic segmentation. Neurocomputing 492, pp. 464-473. 10.1016/j.neucom.2022.04.025

Publishers page: http://dx.doi.org/10.1016/j.neucom.2022.04.025

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See http://orca.cf.ac.uk/policies.html for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



# Scale-aware Network with Modality-awareness for RGB-D Indoor Semantic Segmentation

Feng Zhou<sup>a,\*</sup>, Yu-Kun Lai<sup>b</sup>, Paul L. Rosin<sup>b</sup>, Fengquan Zhang<sup>a</sup>, Yong Hu<sup>c,d</sup>

<sup>a</sup>North China University of Technology <sup>b</sup>Cardiff University <sup>c</sup>State Key Laboratory of Virtual Reality Technology and Systems, Beihang University <sup>d</sup>School of New Media Art and Design, Beihang University

## Abstract

This paper focuses on indoor semantic segmentation based on RGB-D images. Semantic segmentation is a pixellevel classification task that has made steady progress based on fully convolutional networks (FCNs). However, we find there is still room for improvements in the following three aspects. The first is related to multi-scale feature extraction. Recent state-of-the-art works forcibly concatenate multi-scale feature representations extracted by spatial pyramid pooling, dilated convolution or other architectures, regardless of the spatial extent for each pixel. The second is regarding RGB-D modal fusion. Most successful methods treat RGB and depth as two separate modalities and force them to be joined together regardless of their different contributions to the final prediction. The final aspect is about the modeling ability of extracted features. Due to the "local grid" defined by the receptive field, the learned feature representation lacks the ability to model spatial dependencies. In addition to these modules, we design a depth estimation module to encourage the RGB network to extract more effective features. To solve the above challenges, we propose four modules to address them: scale-aware module, modality-aware module, attention module and depth estimation module. Extensive experiments on the NYU-Depth v2 and SUN RGB-D datasets demonstrate that our method is effective against RGB-D indoor semantic segmentation.

Keywords: Semantic segmentation, Scale selection, Attention, RGB-D, Depth estimation

## 1 1. Introduction

The purpose of semantic segmentation is to assign 2 specific class labels to regions in the input images. This 3 is a fundamental task for scene understanding [1], video 4 5 analysis [1, 2], clothing retrieval [3], and such of those intelligent applications [4]. However, scene understand-6 ing is a daunting task, especially for indoor scenes, due 7 to the varying illuminations and cluttered backgrounds. 8 With the development of commercial depth cameras, 9 such as Kinect and Prime-Sense, we are able to cap-10 ture high-quality, synchronized RGB and depth images. 11 RGB data provides rich visual information such as color 12 and texture. In contrast to RGB data, the depth modal-13 ity data provides pure shape and geometry information, 14

\*Corresponding author

*Email addresses:* zhoufeng@ncut.edu.cn (Feng Zhou), laiy4@cardiff.ac.uk (Yu-Kun Lai), RosinPL@cardiff.ac.uk (Paul L. Rosin), fqzhang@ncut.edu.cn (Fengquan Zhang), huyong@buaa.edu.cn (Yong Hu)

Preprint submitted to Neurocomputing

which is invariant to lighting and reflectance. Combining these two complementary modalities provides us with an opportunity to dramatically improve the performance of semantic segmentation of indoor scenes.

Extensive studies have been conducted for the task 19 of indoor semantic segmentation. [5] proposes a patch-20 wise model, and [6] utilizes an R-CNN (Region-based 21 Convolutional Neural Network) scheme to learn an 22 RGB-D multi-modal feature representation to boost the 23 performance. Recently, [7] proposes an end-to-end 24 FCN (Fully Convolutional Network) for semantic seg-25 mentation and achieves significant improvement. How-26 ever, there are still many problems with indoor semantic 27 segmentation. Towards the problem of multi-scale ob-28 jects, many successful methods [8, 9, 10, 11, 12] adopt 29 pyramid layers to extract multi-scale feature represen-30 tations. Towards the problem of modeling long-range 31 contextual information, [9, 13, 14] utilize global pool-32 ing techniques to obtain global context feature, and [15] 33 subdivides images into super-pixels and uses LSTM 34



Figure 1: Limitations of the baseline on indoor scene semantic segmentation with RGB-D data. The depth image in this paper is encoded to three channel HHA (horizontal disparity, height above ground, and angle with gravity). The baseline consists of two-stream atrous spatial pyramid pooling networks trained on RGB and depth data respectively. These two networks are combined together by late fusion with equal-weight sum.

(Long Short-Term Memory) to aggregate and enlarge 35 contextual information by multi-scale context intertwin-36 ing. Towards RGB-D fusion, three levels of fusion are 37 often adopted. The first one is early fusion [5], which 38 simply concatenates the input of two complementary 39 modalities, RGB and depth, together as four-channel in-40 put. The second one is middle fusion [6], which lever-41 ages the two modalities, RGB and depth, as two inde-42 pendent inputs and extracts different modality feature 43 representations, and then concatenates them together to learn a final classifier. The third one is late fusion (also 45 called score map fusion) [7], which utilizes RGB and 46 depth as two separate inputs to learn two different mod-47 els, and obtain two different score maps. Then, the two 48 score maps are fused together by equal weights. 49

In this paper, the model proposed by [8] is extended 50 by using the late fusion strategy. This extended model 51 is used as our baseline for indoor semantic segmenta-52 tion in this study. Compared with [7], our extended 53 model achieves better performance. However, we have 54 found that there are three aspects that can be improved. 55 The first is that the pixel itself does not have enough 56 information for semantic prediction, it needs to learn 57 the appropriate scale information. As shown in Figure 58 1(a), since the appearance of the board object is very 59 similar to the back wall, multi-scale feature representa-60 tions extracted using multiple atrous convolutional lay-61 62 ers do not provide proper surrounding scale information for pixels on the board. Likewise, for the wire object, 63 since it is so thin the extracted multi-scale features do 64 not capture suitable information for it. The second is 65

about the fusion of two complementary modalities. As 66 shown in Figure 1(b), the appearance cues are beneficial 67 for classifying the object as a chair, whereas the depth 68 cue would confuse the identification of chair parts (most 69 sofa objects in the dataset are near a wall). The third one 70 is that the extracted features lack the ability to model 71 long-range dependencies. As shown in Figure 1(c), part 72 of the refrigerator is misclassified as a door (due to be-73 ing confused by the rectangle shape). 74

This paper aims to discuss the problem of indoor 75 semantic segmentation based on the two complemen-76 tary modalities of RGB and depth. In particular, we 77 propose a scale-aware module, a modality-aware mod-78 ule, and an attention module, which address the above 79 three-aspect problems. The scale-aware module learns 80 a proper scale feature representation for each object in 81 the input. It learns a weighted mask for each extracted 82 multi-scale feature, and then multiplies these masks by 83 the multi-scale features to generate a scale-aware fea-84 ture representation to address the first problem. Towards 85 the second problem, the modality-aware module is pro-86 posed to combine the two complementary modalities of 87 RGB and depth using different weights instead of equal 88 weights. Towards the third problem, an attention mod-89 ule is introduced to complement the scale-aware mod-90 ule, which can capture long-range dependencies in the 91 generated scale-aware feature representation. Besides 92 the above three modules, since we have the ground truth 93 depth value of the input RGB image, we can thus de-94 sign an encoder-decoder depth estimation module on the 95 RGB network to encourage the RGB backbone network 96 to extract better and more precise features. The contri-97 butions of this paper can be summarized in the follow-98 ing aspects. 99

- An efficient scale-aware module with modalityawareness, an attention module, and a depth estimation network is proposed for semantic segmentation.
- Within the network, a scale-aware module is used to select the appropriate scale feature for each pixel, which enables a proper scale feature representation to be learned for each object in the input.
- In order to improve the segmentation performance, a modality-aware module is proposed, which adaptively combines the RGB module and depth module to obtain useful features.
- To further improve the segmentation performance, 112 the attention module and depth estimation module 113



Figure 2: The overall architecture of our SAMD model for RGB-D indoor semantic segmentation. It is a two-stream convolutional neural network, one for RGB and the other one for depth (HHA). SAMD consists of four parts: 1) the encoder feature extractor part. It is a standard two-stream convolutional neural network, which leverages atrous spatial pyramid pooling to learn multi-scale feature representations; 2) the scale-aware module, which is used to learn features maps of an appropriate scale; 3) the modality-aware module, which is proposed to effectively combine RGB and depth networks based on the contributions of the two modalities; 4) the attention and depth estimation module, which is used to extract more plausible features. Best viewed in color.

are proposed to extract better feature representations. The former is to obtain long-range dependent features, and the latter is to force the RGB
module to extract more plausible feature representations.

The rest of the paper is organized as follows. Section 2 briefly covers related work, highlighting current work. Then we give the details of the proposed approach in Section 3. Experimental results and analysis are provided in Section 4. Finally, the conclusions are drawn in Section 5.

#### 125 2. Related Work

The proposed method relates to a lot of work on 126 scale-aware selection, attention method, modal combi-127 nation and depth estimation. CNN-based semantic seg-128 mentation has achieved great advances in recent years 129 [16, 8, 13, 9, 17, 18, 19, 20]. Most of the existing work 130 has employed fully convolutional networks (FCNs) [7]. 131 However, objects in indoor scenes cover a huge range of 132 scales due to both their range of actual sizes in the real 133 world, as well as by their differences in distance to the 134 camera. The methods above only forcibly stack the ex-135 tracted multi-scale features together. This is not enough 136 137 for real-world cluttered indoor scene understanding.

Selecting the appropriate scale feature for each pixelis particularly important. Many successful works have

investigated this problem. [21] proposes a channel at-140 tention scheme to boost the performance of semantic 141 segmentation. [22] exploits a scale-space to select a 142 properly scaled feature. However, as far as we know, 143 there is little work on RGB-D feature scale selection. In 144 this paper, we propose a scale-aware module that com-145 bines RGB and depth modal features to build a scale-146 aware module to improve the performance of RGB-D 147 semantic segmentation. Although the scale-aware mod-148 ule can generate features that fit the scale for each neu-149 ron, the feature cannot reflect the contributions of each 150 modality. 151

The synchronized RGB and depth pair images pro-152 vide useful multi-modal information for the task of 153 computer vision. Most successful methods simply com-154 bine the extracted multi-modal feature representation 155 using early fusion [5], middle fusion [6], or late fusion 156 [7]. However, in the final prediction layer, RGB and 157 depth contribute unequally in most cases. An exam-158 ple is shown in Figure 1 (b) where the chair object is 159 misclassified by concatenating the two complementary 160 modalities with the same weight. 161

Recently, the attention mechanism has been proposed to model and capture long-range dependencies, and it has become an integral part of many successful works [23, 24, 25, 26]. [27] proposes a self-attention mechanism to capture long-range dependencies of inputs and achieves the state-of-the-art performance in machine translation. The attention mechanism has not only been used in the Natural Language Processing (NLP) field,

but has also been utilized in the computer vision field. 170 [28] utilizes a self-attention scheme to obtain better per-171 formance on the image generation task. [29] adopts 172 an attention mechanism in object recognition to boost 173 performance. [30] proposes a MAT (Motion-Attentive 174 Transition) module comprised of a soft attention unit 175 and an attention transition unit to learn more specific 176 and useful feature representations. 177

The combination of semantic segmentation and depth 178 estimation was studied in many previous works, with 179 the goal of improving both semantic segmentation and 180 depth estimation. [31] proposes three ways to improve 181 semantic segmentation performance with depth estima-182 tion, and [32] adopts knowledge from a semantic seg-183 mentation network to teach the depth estimation task. 184 In our paper, we introduce depth estimation as an auxil-185 iary task to help improve semantic segmentation. 186

This paper adopts a scale-aware module, a modalityaware module, a self-attention and a depth estimation module to address the above problems. As shown in the experiments, the proposed four modules can achieve performance gains on many publicly RGB-D semantic segmentation datasets.

## 193 **3. Our Approach**

In the following section, we mainly focus on the 194 learning details of the proposed SAMD method. SAMD 195 is composed of four modules: the scale-aware module, 196 the modality-aware module, attention, and depth esti-197 mation (as shown in Figure 2). The scale-aware mod-198 ule is used to generate a scale-aware feature representa-199 tion which predicts the scale information for each pixel 200 from the learned multi-scale feature representation. The 201 modality-aware module is to learn an effective fusion 202 way for the two modal networks. To further improve 203 the performance, we propose the attention and depth es-204 timation modules. The attention module is used to cap-205 ture the global feature dependencies in the spatial do-206 main for the input feature. The depth estimation module 207 is used to push the RGB network to extract more precise 208 and useful features. 209

We adopt atrous spatial pyramid pooling (ASPP) as our feature encoder to extract multi-scale features. To be specific, let  $\mathcal{L} = \{(\mathcal{R}_1, \mathcal{D}_1, Y_1), ..., (\mathcal{R}_n, \mathcal{D}_n, Y_n)\}$ be the *n* pairwise RGB-D training data, where  $\mathcal{R} =$  $\{r_i\}_{i=1}^{H \times W}$  is the RGB modality training image whose size is  $H \times W$ , and  $\mathcal{D} = \{d_i\}_{i=1}^{H \times W}$  is the corresponding depth training image, whose size is  $H \times W$ , and  $Y = \{y_i\}_{i=1}^{H \times W}$  is the label image, in which  $r_i$  and  $d_i$  are corresponding pixels in the pairwise image, la-218 bel  $y_i \in \{0, 1, ..., C\}$  gives the per-pixel label, C de-219 notes the number of the categories. In our approach, 220 given an  $H \times W$  pair RGB-D image, through the en-221 coder part, we obtain features  $f_e^r$  and  $f_e^d$  whose sizes are 222  $\frac{H}{8} \times \frac{W}{8}$  (ignoring the channel size), where  $f_e^r$  is from RGB modality, and  $f_e^d$  is from depth modality. These 223 224 two features serve two purposes. The first is to gener-225 ate subsequent multi-scale feature representations and 226 the second is used in our scale-aware module for scale 227 selection. 228

#### 3.1. Scale-aware Module

The output of the feature encoder part is a multi-scale 230 feature of the forced concatenation, but the learned feature still does not hold the correct scale feature representation. To this end, we employ a scale-aware module to 233 enable our model to learn a feature map of proper scale 234 for all neurons in the input. 235

Specifically, let the multi-scale feature set generated from  $f_e^r$  and  $f_e^d$  be  $\{f_{a_i}^r, f_{a_i}^d\}$ , where  $f_{a_i}^r$  denotes the multi-scale feature extracted from the RGB modality feature encoder part by dilated convolution (a.k.a. atrous convolution) [8] (kernel size  $a_i$ ),  $f_{a_i}^d$  denotes the feature from the depth modality. In the experiments, we adopt four dilated kernel sizes (6, 12, 18, 24) for each modality, and  $a_i$  stands for the 4 different dilated kernels. We concatenate  $f_e^r$  and  $f_e^d$  to generate  $f_e^f$ , and feed it into a  $1 \times 1$  convolutional layer conv(.) and output  $f_c^f$  whose size is  $8 \times \frac{H}{8} \times \frac{W}{8}$ . Then we use a softmax operation to normalize  $f_c^f$  to obtain  $f_m^f$ . For the RGB modality, we split the four channels of  $f_e^r$ )  $f_{m_j}^f$  in  $f_m^f$  and then calculate the scale-aware features  $f_{sa}^r$  and  $f_{sa}^d$  as follows:

$$f_{sa}^r = \sum_{j=1}^4 f_{a_i}^r \odot f_{m_j}^f \tag{1}$$

to the depth modality,

$$f_{sa}^{d} = \sum_{j=5}^{8} f_{a_{i}}^{d} \odot f_{m_{j}}^{f}$$
(2)

where the operator  $\odot$  represents the Hadamard product. 237

#### 3.2. Modality-aware Module

The modality-aware module is proposed to combine feature representations of RGB and depth modalities for semantic segmentation. The structure of the module is

236

238

229



Figure 3: Illustration of the scale-aware confidence map. The two images in the first column are RGB and depth images. The two images in the second column are scale-aware confidence maps upon the two modalities. The remaining 4 images in the first row are each scale channel confidence map of RGB modality, in the second row are each channel confidence map of depth modality. From left to right, they are  $a_6$ ,  $a_{12}$ ,  $a_{18}$ ,  $a_{24}$ . For the sake of simplicity, we show the confidence maps (average value) by using the "COLORMAP\_JET" color map (where blue is low value and red is high value) upon RGB. Best viewed in color.



Figure 4: Illustration of the scale-aware module. From the illustration, we can find that with the scale-aware module, our model can effectively focus on larger and smaller objects, as shown by the red dashed boxes. From left to right and top to bottom, they are RGB and ground truth; prediction results with and without the scale-aware module; confidence maps with and without scale-aware module (including the colorbar where the value increases from blue to red); error maps with and without scale-aware module. Best viewed in color.

similar to the scale-aware module and it is composed of 242 four layers. The first one is a concatenation layer which 243 is used to combine the  $f_{sa}^r$  feature and  $f_{sa}^d$  feature. The 244 second one is a  $1 \times 1$  convolutional layer which is used 245 to produce an  $M^{2 \times h \times w}$  modal mask. The last two are 246 a softmax layer and a matrix multiplication layer. The 247 former is used to generate a normalized modal mask and 248 the latter is used for element-wise multiplication. For 249 brevity and clarity, the layers are not illustrated in Fig-250 ure 2. 251

To be more specific regarding the structure of the modality-aware module, after the concatenation layer, we obtain RGB-D fusion feature representation  $f_f \in \mathbb{R}^{(2c \times h \times w)}$ , and then feed it to the  $1 \times 1$  convolutional layer to produce the mask M. Then M is fed into the softmax layer to produce a normalized modal mask  $M' \in \mathbb{R}^{2 \times h \times w}$ . Let  $M^{rgb} \in \mathbb{R}^{1 \times h \times w}$  and  $M^{depth} \in \mathbb{R}^{1 \times h \times w}$  denote the modal masks on RGB and depth respectively. When the modal masks are generated, we calculate the predictions based on RGB and

depth using the Hadamard product as follows:

$$P^{rgb} = Conv(f^r_{at}) \odot M^{rgb}$$

$$P^{depth} = Conv(f^d_{at}) \odot M^{depth}$$
(3)

where Conv(.) denotes a  $1 \times 1$  convolutional layer, and  $Conv(f_{at}^r) \in \mathbb{R}^{C \times h \times w}, Conv(f_{at}^d) \in \mathbb{R}^{C \times h \times w}$ . The elements in  $P(i, j)^{rgb}$  and  $P(i, j)^{depth}$  imply how confidently we can rely on RGB and depth respectively to predict the pixel (i, j) in the input. 252

Finally, we generate the final prediction result as follows:

$$P^f = P^{rgb} + P^{depth} \tag{4}$$

257

## 3.3. Attention and Depth-estimation Modules

To further improve the segmentation performance, <sup>258</sup> we propose attention and depth estimation modules to <sup>259</sup> obtain long-range dependencies and more plausible feature representations. In order to enlarge the context relationship of the above-obtained  $f_{sa}^r$  and  $f_{sa}^d$  features, <sup>262</sup> inspired by [33], we introduce a self-attention module to improve the obtained feature modeling ability. <sup>264</sup>

For the sake of simplicity, let the size of the image feature obtained from the RGB modality scale-aware layer be  $f_{sa}^r \in \mathbb{R}^{c \times h \times w}$ , where c denotes the feature channel,  $h = \frac{H}{8}, w = \frac{W}{8}$ . Take the RGB modality as an example for explanation. We copy the  $f_{sa}^r$  and reshape it into three feature spaces,  $\Theta(f) \in \mathbb{R}^{c \times N}$ ,  $\theta(f) \in \mathbb{R}^{c \times N}$ , and  $\vartheta(f) \in \mathbb{R}^{c \times N}$ , respectively, where  $N = h \times w$ . Then we calculate self-attention using the above two reshaped features, as follows:

$$at = softmax(\Theta(f)^T \cdot \theta(f))$$
(5)

each item at(i, j) in the module at is the dot-product similarity, which indicates the effect of the model at

the *i*th position to the *j*th position. To make it more implementation-friendly, we normalize the attention module before the softmax operation. Then we obtain the scale-aware attention feature representation as follows:

$$f_{at}^r = f_{sa}^r + \beta(at \cdot \vartheta(f)).$$
(6)

For the depth modality,  $f_{at}^d$  is similarly defined as  $f_{at}^r$ , where  $\beta$  is a learnable parameter, and is initialized to 0 during training inspired by [28]. The scheme mentioned above makes our model rely on non-attention features in the initial stages of training. For the depth modality, we utilize the same operation as the depth image feature representation. For the depth estimation module, we adopt the struc-

For the depth estimation module, we adopt the structure of Monodepth2 [34], which is a successful depth estimation model. For simplicity of training, we adopt Depth Loss and Gradient Loss.

$$L_{depth} = \frac{1}{n} \sum_{i=1}^{n} \sqrt{\log^2(d_i) - \log^2(d_i^{Gt})}$$
(7)

$$L_{grad} = \frac{1}{n} \sum_{i=1}^{n} \left\| \nabla(d) - \nabla(d^{Gt}) \right\|_{1}$$
(8)

where n is the number of pixels in the input image,  $d_i$ 272 and  $d_i^{Gt}$  denote the predicted depth value and the corre-273 sponding ground truth depth value, respectively. In the 274 experiments, the main purpose of the task is to obtain 275 a per-pixel semantic segmentation label, and the depth 276 estimation module is to encourage the RGB network 277 to extract a more effective feature representation. We 278 use the pre-trained Monodepth2 model to initialize our 279 depth estimation module, and then use a small learning 280 rate (1e-4) to fine-tune it in the final experiments. 281

#### 282 4. Experiments

In this section, we perform extensive experiments on two publicly available datasets, NYU-Depth v2 and SUN RGB-D to evaluate our method. All of our implementations are made using the popular PyTorch framework.

## 288 4.1. Datasets

 NYU-Depth V2 is one of the most popular RGB-D indoor scene datasets, consisting of 1449 finely
 labeled RGB and depth image pairs. The entire
 dataset is divided into two parts, of which 795 are
 for training and 654 are for testing. SUN RGB-D is a large-scale RGB-D dataset recently used for indoor scene understanding. It contains 10335 pairs of RGB and depth images captured by four kinds of commercial depth sensors. Of these finely labeled image pairs, 5285 pairs are used for training and the remaining 5050 pairs are used for testing.

301

319

320

321

## 4.2. Metrics

Following recent methods [10, 35], performance in 302 our experiments is quantitatively measured by pixel ac-303 curacy (Acc), mean intersection over union (mIoU), 304 mean pixel accuracy of different categories (mAcc) and 305 frequency weighted IoU (f.w. IoU), which are widely 306 used in indoor semantic segmentation. To be concrete, 307 let  $n_{ij}$  be the number of pixels which are misclassified 308 as class j when the ground truth is category i.  $t_i$  is 309 the number of pixels which belong to the *i*th category, 310 where  $t_i = \sum_j n_{ij}$ , and the total number of pixels in the dataset is t. The above four metrics are defined as 311 312 follows: 313

- pixel accuracy =  $\sum_{i} \frac{n_{ii}}{t}$  314
- mean intersection over union =  $\frac{1}{C}\sum_{i} \frac{n_{ii}}{t_i + \sum_{i} n_{ji} n_{ii}}$  316
- mean pixel accuracy =  $\frac{1}{C} \sum_{i} \frac{n_{ii}}{t_i}$  317
- frequency weighted IoU =  $\frac{1}{t}\sum_{i} \frac{t_i n_{ii}}{t_i + \sum_j n_{ji} n_{ii}}$  318

## 4.3. Training Protocol

In the following, we will provide details of the experimental implementation.

Learning rate policy The training procedure con-322 sists of two stages. In the first stage, we adopt the Adam 323 optimizer to train two independent networks of RGB 324 and depth modalities respectively for semantic segmen-325 tation, excluding the scale-aware and modality-aware 326 modules. For each modality network, we adopt "poly" 327 learning rate policy, where the current learning rate is 328 calculated by multiplying the initial learning rate with 329  $(1 - \frac{iter}{max.iter})^{power}$ , power = 0.9, the initial learning rate is set to 0.01. We use ResNet50 and ResNet101 330 331 as our backbone network and combine atrous spatial 332 pyramid pooling as our feature encoder to extract multi-333 scale features. Each of the backbones is initialized by 334 the model pre-trained on ImageNet, and the other lay-335 ers are initialized by random weights. In the second 336 stage, we add the scale-aware module and the modality-337 aware module and then fine-tune our RGB-D model on 338 the synchronized RGB and depth training data. Each 339



Figure 5: The visualization of results on the NYU-Depth v2 dataset. The comparison results of (d) and (f) demonstrate that our SAMD module is effective for indoor semantic segmentation. For the detailed analysis, please refer to Section 4.4. Best viewed in color.

modality network is initialized by the trained models
obtained from the first stage. During the training, we
discard the classification layer in the already-trained
network in every single modality, and then combine
them together via the added scale-aware and modalityaware modules. In the second stage, we set the initial
learning rate to 0.001.

Data preprocessing and data augmentation In 347 the experiment, the depth modality image is encoded to 348 three-channel HHA (horizontal disparity, height above 349 ground, and angle with gravity) image as the approach 350 [6]. In both training stages, our two separate modal-351 ity networks and our RGB-D model are trained on the 352 cropped images of size  $417 \times 417$ . To avoid over-fitting, 353 common data augmentations such as random brightness 354 jittering, random left-right flipping, and random scaling 355 in the range of [0.5, 2.0] to the input training samples 356 are used. 357

**Loss** In the experiments, the overall loss is as follows:

$$L = L_{seg} + \lambda_1 \cdot L_{aux} + \lambda_2 \cdot L_{dep} \tag{9}$$

where  $\lambda_1$  and  $\lambda_2$  are the balancing weights for the semantic segmentation and depth estimation. To enhance the feature representation extracted from the backbone, we adopt an auxiliary loss after the 4th blocks (as used in [13]) to supervise the training process. In the experiments,  $\lambda_1$  is set to 0.5, and  $\lambda_2$  is set to 0.1.  $L_{dep}$  is composed by  $L_{Depth}$  and  $L_{grad}$ .



Figure 6: Illustration of the self-attention module. We observe that, with this module, the extracted feature representation is better (as shown in the red dashed bounding boxes). From left to the right and top to down, they are RGB and ground truth images; the RGB feature and HHA feature without attention module; RGB and depth feature with attention module; prediction results with and without attention module.

#### 4.4. Ablation Studies and Discussion

In order to demonstrate that our SAMD model 366 does not depend on any particular feature encoder ar-367 chitecture, we embed scale-aware module, attention 368 module and modality-aware module into two standard 369 fully convolutional backbone networks, ResNet50 and 370 ResNet101. We provide the quantitative results on the 371 NYU-Depth v2 dataset of these two backbone networks 372 in Table 3. Through the results, we find that using 373 our SAMD module significantly improves the perfor-374 mance of semantic segmentation throughout the differ-375 ent backbones. In the experiments, we use ResNet50 376 and ResNet101 as alternatives for our backbone net-377 work, and the default choice is ResNet101 if not explic-378 itly specified. 379

365

In order to show that our scale-aware module does not 380

Table 1: Category-wise IoU results on the NYU-Depth v2 dataset. The Baseline and SAMD rows show the results of our baseline and SAMD model respectively. The class of background is ignored during performance evaluation. The top two results are shown in red and blue respectively. Cheng<sup> $\ddagger$ </sup> pre-trains their model on SUN RGB-D dataset, and then fine-tunes it on NYU-Depth v2 dataset.

	wall	floor	cabinet	bed	chair	sofa	table	door	window	bookshelf	picture	counter	blinds	desk	shelves	curtain	dresser	pillow	mirror	floormat	clothes	ceiling
Long [7]	69.9	79.4	50.3	66.0	47.5	53.2	32.8	22.1	39.0	36.1	50.5	54.2	45.8	11.9	8.6	32.5	31.0	37.5	22.4	13.6	18.3	59.1
Gupta [6]	68.0	81.3	44.9	65.0	47.9	47.9	29.9	20.3	32.6	18.1	40.3	51.3	42.0	11.3	3.5	29.1	34.8	34.4	16.4	28.0	4.7	60.5
Deng [36]	65.6	79.2	51.9	66.7	41.0	55.7	36.5	20.3	33.2	32.6	44.6	53.6	49.1	10.8	9.1	47.6	27.6	42.5	30.2	32.7	12.6	56.7
He [37]	72.7	85.7	55.4	73.6	58.5	60.1	42.7	30.2	42.1	41.9	52.9	59.7	46.7	13.5	9.4	40.7	44.1	42.0	34.5	35.6	22.2	55.9
Cheng <sup>‡</sup> [10]	78.5	87.1	56.6	70.1	65.2	63.9	46.9	35.9	47.1	48.9	54.3	66.3	51.7	20.6	13.7	49.8	43.2	50.4	48.5	32.2	24.7	62.0
Wang [35]	-	_	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Daniel [38]	-	_	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Gu [ <mark>39</mark> ]	-	_	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Zhou [40]	80.1	88.3	61.7	72.8	63.9	65.4	48.0	46.5	48.3	44.4	61.4	69.9	59.5	27.2	16.8	59.3	50.6	50.9	51.3	38.6	25.1	79.5
Lin [41]	80.5	87.6	63.0	72.3	63.9	68.7	51.1	37.6	52.1	44.7	60.0	69.2	63.1	30.5	15.6	60.3	49.3	47.3	58.7	42.6	30.4	70.0
Cao [42]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
Baseline	77.4	86.5	59.4	76.6	63.1	65.4	39.8	38.1	51.2	37.5	59.5	67.5	60.8	17.0	13.6	46.2	51.6	44.5	55.1	29.9	16.5	73.4
SAMD	82.3	89.7	62.3	73.0	64.8	67.7	51.0	46.8	51.8	46.9	65.5	71.3	62.2	23.4	19.7	60.1	48.8	49.7	51.7	42.0	26.6	81.2
	books	fridge	tv	paper	towel	shower	box	board	person	nightstand	toilet	sink	lamp	bathhub	bag	ostuct	ofurn	oprops	Acc	mloU	mAcc	f.w. IoU
Long [7]	syooq 27.3	eline firidge	<u>≥</u> 41.9	baber 15.9	1000 tower	shower 14.1	ход 6.5	poard 12.9	berson 57.6	nightstand	toilet 61.3	yuis 44.8	dung 32.1	pathhub 39:5	seq 4.8	ostroct 15.2	unjo 7.7	obrops 30.0	8 V 65.4	Doju 34.0	оружи 29 ш 46.1	1001 .w. 19.5
Long [7] Gupta [6]	syooq 27.3 6.4	eggi 27.0 14.5	≥ 41.9 31.0	Laber 15.9 14.3	26.1 16.3	shower 14.1 4.2	<u>ко</u> 6.5 2.1	pogr 12.9 14.2	beison 57.6 0.2	uightstand 30.1 27.2	toilet 61.3 55.1	¥ 44.8 37.5	du <u>32.1</u> 34.8	qnquq paququ 39.2 38.2	бер 4.8 0.2	tontso 15.2 7.1	umjo 7.7 6.1	sdoudo 30.0 23.1	8 ¥ 65.4 60.3	Doju 34.0 28.6	2000 m 46.1	nor »: 49.5 47.0
Long [7] Gupta [6] Deng [36]	27.3 6.4 8.9	27.0 14.5 21.6	≥ 41.9 31.0 19.2	Ladiad 15.9 14.3 28.0	26.1 16.3 28.6	14.1 4.2 22.9	6.5 2.1 1.6	prod 12.9 14.2 1.0	uosiad 57.6 0.2 9.6	30.1 30.6	to 61.3 55.1 48.4	<u>ч</u> 44.8 37.5 41.8	32.1 34.8 28.1	qnutupped 39.2 38.2 27.6	50 4.8 0.2 0	tongo 15.2 7.1 9.8	umjo 7.7 6.1 7.6	sdoudo 30.0 23.1 24.5	8 4 65.4 60.3 63.8	Dolm 34.0 28.6 31.5	200 46.1 - -	noI 49.5 47.0 48.5
Long [7] Gupta [6] Deng [36] He [37]	<sup>sy</sup> ooq 27.3 6.4 8.9 29.8	27.0 14.5 21.6 41.7	≥ 41.9 31.0 19.2 52.5	15.9 14.3 28.0 21.1	26.1 16.3 28.6 34.4	14.1 4.2 22.9 15.5	6.5 2.1 1.6 7.8	prood 12.9 14.2 1.0 29.2	57.6 0.2 9.6 60.7	30.1 30.6 42.2	tojet 61.3 55.1 48.4 62.7	¥uis 44.8 37.5 41.8 47.4	32.1 34.8 28.1 38.6	qnyy 39.2 38.2 27.6 28.5	50 eq 4.8 0.2 0 7.3	tonto 15.2 7.1 9.8 18.8	7.7 6.1 7.6 15.1	30.0 23.1 24.5 31.4	8 4 65.4 60.3 63.8 70.1	Dofu 34.0 28.6 31.5 40.1	оуни 46.1 - 53.8	001 ·
Long [7] Gupta [6] Deng [36] He [37] Cheng <sup>‡</sup> [10]	27.3 6.4 8.9 29.8 34.2	27.0 14.5 21.6 41.7 45.3	≥ 41.9 31.0 19.2 52.5 53.4	15.9 14.3 28.0 21.1 27.7	26.1 16.3 28.6 34.4 42.6	land land land land land land land land	6.5 2.1 1.6 7.8 11.2	ppp 12.9 14.2 1.0 29.2 58.8	57.6 0.2 9.6 60.7 53.2	30.1 27.2 30.6 42.2 54.1	tə ioi 61.3 55.1 48.4 62.7 80.4	¥uis 44.8 37.5 41.8 47.4 59.2	32.1 34.8 28.1 38.6 45.5	qnuutpeq 39.2 38.2 27.6 28.5 52.6	4.8 0.2 0 7.3 15.9	15.2 7.1 9.8 18.8 12.7	7.7 6.1 7.6 15.1 16.4	30.0 23.1 24.5 31.4 29.3	8 4 65.4 60.3 63.8 70.1 71.9	Dolm 34.0 28.6 31.5 40.1 45.9	46.1 - 53.8 60.7	Pol
Long [7] Gupta [6] Deng [36] He [37] Cheng <sup>‡</sup> [10] Wang [35]	27.3 6.4 8.9 29.8 34.2 -	27.0 14.5 21.6 41.7 45.3 -	≥ 41.9 31.0 19.2 52.5 53.4 -	15.9 14.3 28.0 21.1 27.7 -	26.1 16.3 28.6 34.4 42.6 -	14.1 4.2 22.9 15.5 23.9 -	6.5 2.1 1.6 7.8 11.2	ppog 12.9 14.2 1.0 29.2 58.8 –	uosiad 57.6 0.2 9.6 60.7 53.2 –	30.1 27.2 30.6 42.2 54.1 -	tojet 61.3 55.1 48.4 62.7 80.4 -	¥US 44.8 37.5 41.8 47.4 59.2 –	32.1 34.8 28.1 38.6 45.5 -	qnutup 39.2 38.2 27.6 28.5 52.6 -	4.8 0.2 0 7.3 15.9	15.2 7.1 9.8 18.8 12.7 -	umjo 7.7 6.1 7.6 15.1 16.4 –	sdoudd 30.0 23.1 24.5 31.4 29.3 -	8 65.4 60.3 63.8 70.1 71.9 -	Dofm 34.0 28.6 31.5 40.1 45.9 43.9	25 WH 46.1 - 53.8 60.7 53.5	nol wj 49.5 47.0 48.5 55.7 59.3 -
Long [7] Gupta [6] Deng [36] He [37] Cheng <sup>‡</sup> [10] Wang [35] Daniel [38]	xyooq 27.3 6.4 8.9 29.8 34.2 - -	27.0 14.5 21.6 41.7 45.3 -	≥ 41.9 31.0 19.2 52.5 53.4 - -	15.9 14.3 28.0 21.1 27.7 -	26.1 16.3 28.6 34.4 42.6 -	14.1 4.2 22.9 15.5 23.9 - -	6.5 2.1 1.6 7.8 11.2 -	12.9 14.2 1.0 29.2 58.8 - -	57.6 0.2 9.6 60.7 53.2 –	30.1 27.2 30.6 42.2 54.1 -	tion 61.3 55.1 48.4 62.7 80.4 - -	<u>+yii</u> 44.8 37.5 41.8 47.4 59.2 – –	32.1 34.8 28.1 38.6 45.5 -	qnuuture 39.2 38.2 27.6 28.5 52.6 - -	4.8 0.2 0 7.3 15.9 -	15.2 7.1 9.8 18.8 12.7 –	15.1 16.4 -	30.0 23.1 24.5 31.4 29.3 - -	8 4 65.4 60.3 63.8 70.1 71.9 - -	Doju 34.0 28.6 31.5 40.1 45.9 43.9 51.6	90000000000000000000000000000000000000	Dol
Long [7] Gupta [6] Deng [36] He [37] Cheng <sup>‡</sup> [10] Wang [35] Daniel [38] Gu [39]	27.3 6.4 8.9 29.8 34.2 - -	27.0 14.5 21.6 41.7 45.3 - -	≥ 41.9 31.0 19.2 52.5 53.4 - -	15.9 14.3 28.0 21.1 27.7 - -	26.1 16.3 28.6 34.4 42.6 - -	14.1 4.2 22.9 15.5 23.9 - -	6.5 2.1 1.6 7.8 11.2 - -	12.9 14.2 1.0 29.2 58.8 - - -	57.6 0.2 9.6 60.7 53.2 - -	30.1 27.2 30.6 42.2 54.1 - -	tion 61.3 55.1 48.4 62.7 80.4 - - -	¥ 44.8 37.5 41.8 47.4 59.2 - - -	32.1 34.8 28.1 38.6 45.5 - - -	qnuutura 39.2 38.2 27.6 28.5 52.6 - - -	4.8 0.2 0 7.3 15.9 - -	15.2 7.1 9.8 18.8 12.7 - -	7.7 6.1 7.6 15.1 16.4 - -	30.0 23.1 24.5 31.4 29.3 - -	8 4 65.4 60.3 63.8 70.1 71.9 - - -	Doju           34.0           28.6           31.5           40.1           45.9           43.9           51.6           50.3	46.1 - 53.8 60.7 53.5 -	Pol i i i i i i i i i i i i i
Long [7] Gupta [6] Deng [36] He [37] Cheng <sup>‡</sup> [10] Wang [35] Daniel [38] Gu [39] Zhou [40]	27.3 6.4 8.9 29.8 34.2 - - 33.5	27.0 14.5 21.6 41.7 45.3 - - 56.0	≥ 41.9 31.0 19.2 52.5 53.4 − − 60.8	15.9 14.3 28.0 21.1 27.7 - - 31.7	26.1 16.3 28.6 34.4 42.6 - - 47.7	14.1 4.2 22.9 15.5 23.9 - - 25.3	6.5 2.1 1.6 7.8 11.2 - - 14.8	12.9 14.2 1.0 29.2 58.8 - - 83.7	U05394 57.6 0.2 9.6 60.7 53.2 - - 77.6	30.1 27.2 30.6 42.2 54.1 - 40.2	61.3 55.1 48.4 62.7 80.4 - - 83.8	¥U 44.8 37.5 41.8 47.4 59.2 - - 67.3	32.1 34.8 28.1 38.6 45.5 - - 48.2	900 900 900 900 900 900 900 900	4.8 0.2 0 7.3 15.9 - - 11.0	15.2 7.1 9.8 18.8 12.7 - - 30.6	7.7 6.1 7.6 15.1 16.4 - - 21.2	30.0 23.1 24.5 31.4 29.3 - - 39.2	8 4 65.4 60.3 63.8 70.1 71.9 - - 76.6	Doju           34.0           28.6           31.5           40.1           45.9           43.9           51.6           50.3           51.2	46.1 - 53.8 60.7 53.5 - 63.8	Pol × 49.5 47.0 48.5 55.7 59.3 - - -
Long [7] Gupta [6] Deng [36] He [37] Cheng <sup>‡</sup> [10] Wang [35] Daniel [38] Gu [39] Zhou [40] Lin [41]	27.3 6.4 8.9 29.8 34.2 - - 33.5 <b>37.8</b>	27.0 14.5 21.6 41.7 45.3 - 56.0 56.2	≥ 41.9 31.0 19.2 52.5 53.4 - - 60.8 67.1	15.9 14.3 28.0 21.1 27.7 - - 31.7 32.5	26.1 16.3 28.6 34.4 42.6 - - 47.7 44.2	14.1 4.2 22.9 15.5 23.9 - - 25.3 <b>39.1</b>	6.5 2.1 1.6 7.8 11.2 - - 14.8 12.5	12.9 14.2 1.0 29.2 58.8 - - 83.7 52.6	E05394 57.6 0.2 9.6 60.7 53.2 - 77.6 82.6	90000000000000000000000000000000000000	ting 61.3 55.1 48.4 62.7 80.4 - - 83.8 68.2	¥ 44.8 37.5 41.8 47.4 59.2 - - 67.3 63.8	32.1 34.8 28.1 38.6 45.5 - - 48.2 45.2	99.2 39.2 38.2 27.6 28.5 52.6 - - 66.2 61.4	4.8 0.2 0 7.3 15.9 - - 11.0 21.5	15.2 7.1 9.8 18.8 12.7 - - 30.6 34.7	7.7 6.1 7.6 15.1 16.4 - 21.2 18.3	30.0 23.1 24.5 31.4 29.3 - - 39.2 44.8	8 4 65.4 60.3 63.8 70.1 71.9 - - 76.6 77.0	Doju 34.0 28.6 31.5 40.1 45.9 43.9 51.6 50.3 51.2 51.2	¥ 46.1 - 53.8 60.7 53.5 - 63.8 64.0	Pol *J 49.5 47.0 48.5 55.7 59.3 - - - -
Long [7] Gupta [6] Deng [36] He [37] Cheng <sup>‡</sup> [10] Wang [35] Daniel [38] Gu [39] Zhou [40] Lin [41] Cao [42]	27.3 6.4 8.9 29.8 34.2 - - 33.5 37.8 -	27.0 14.5 21.6 41.7 45.3 - 56.0 56.2 -	≥ 41.9 31.0 19.2 52.5 53.4 - - 60.8 67.1 -	15.9 14.3 28.0 21.1 27.7 - - 31.7 32.5 -	26.1 16.3 28.6 34.4 42.6 - - 47.7 44.2 -	14.1 4.2 22.9 15.5 23.9 - - 25.3 39.1 -	6.5 2.1 1.6 7.8 11.2 - - 14.8 12.5 -	12.9 14.2 1.0 29.2 58.8 - - 83.7 52.6 -	57.6 0.2 9.6 60.7 53.2 - 77.6 82.6 -	pure status 30.1 27.2 30.6 42.2 54.1 - 40.2 47.1 -	tine 61.3 55.1 48.4 62.7 80.4 - - 83.8 68.2 -	<u>¥</u> 44.8 37.5 41.8 47.4 59.2 - - 67.3 63.8 -	32.1 34.8 28.1 38.6 45.5 - - 48.2 45.2 -	99.2 39.2 38.2 27.6 28.5 52.6 - - 66.2 61.4 -	500 4.8 0.2 0 7.3 15.9 - - 11.0 21.5 -	15.2 7.1 9.8 18.8 12.7 - - 30.6 34.7 -	100 7.7 6.1 7.6 15.1 16.4 - - 21.2 18.3 -	30.0 23.1 24.5 31.4 29.3 - - 39.2 44.8 -	<del>8</del> 65.4 60.3 63.8 70.1 71.9 - - 76.6 <b>77.0</b> 76.4	Doll 34.0 28.6 31.5 40.1 45.9 43.9 51.6 50.3 51.2 51.2 51.2 51.3	¥ 46.1 - 53.8 60.7 53.5 - 63.8 64.0 63.5	Por *; 49.5 47.0 48.5 55.7 59.3 - - - 63.0
Long [7] Gupta [6] Deng [36] He [37] Cheng <sup>‡</sup> [10] Wang [35] Daniel [38] Gu [39] Zhou [40] Lin [41] Cao [42] Baseline	<u>sy</u> ooq 27.3 6.4 8.9 29.8 34.2 - - 33.5 <b>37.8</b> - 24.9	27.0 14.5 21.6 41.7 45.3 - 56.0 56.2 - 45.8	≥ 41.9 31.0 19.2 52.5 53.4 - - 60.8 67.1 - 53.2	15.9 14.3 28.0 21.1 27.7 - - 31.7 32.5 - 23.2	26.1 16.3 28.6 34.4 42.6 - - 47.7 44.2 - 39.8	14.1 4.2 22.9 15.5 23.9 - - 25.3 39.1 - 27.1	<u>če</u> 6.5 2.1 1.6 7.8 11.2 - - 14.8 12.5 - 5.1	12.9 14.2 1.0 29.2 58.8 - - 83.7 52.6 - 73.5	57.6 0.2 9.6 60.7 53.2 - 77.6 82.6 - 64.9	pure status 30.1 27.2 30.6 42.2 54.1 - 40.2 47.1 - 38.4	tine 61.3 55.1 48.4 62.7 80.4 - - 83.8 68.2 - 86.2	<u>viji</u> 44.8 37.5 41.8 47.4 59.2 - - 67.3 63.8 - 67.5	32.1 34.8 28.1 38.6 45.5 - - 48.2 45.2 - 48.2 45.2 - 43.3	-01 -01 -01 -02 -02 -02 -02 -02 -02 -02 -02	500 4.8 0.2 0 7.3 15.9 - - 11.0 21.5 - 5.1	15.2 7.1 9.8 18.8 12.7 - - 30.6 34.7 - 28.0	19.99 19.99 19.90 19.90 19.90 19.90	30.0 23.1 24.5 31.4 29.3 - - 39.2 44.8 - 38.6	8 65.4 60.3 63.8 70.1 71.9 - - 76.6 77.0 76.4 73.4	Definition           34.0           28.6           31.5           40.1           45.9           43.9           51.6           50.3           51.2           51.3           46.7	97 46.1 - 53.8 60.7 53.5 - 63.8 64.0 63.5 61.7	Por *; 49.5 47.0 48.5 55.7 59.3 - - - 63.0 61.2

Figure 7: Performance Analysis. Depth estimation on NYUDepthv2 dataset.



depend on feature types, we utilize two methods, atrous spatial pyramid pooling (ASPP) and pyramid pooling (PSP) to extract multi-scale feature representations. In both experiments, we keep all settings exactly the same and extract four kinds of scale features in each modal network. We find that the use of ASPP (52.3) to extract multi-scale features is slightly better than PSP (52.1).

To demonstrate the effectiveness of the scale-aware module, we compare the results of using and not using this module. The qualitative analysis is shown in Figure 4. From the comparison result, with the scaleaware module, our model learns more proper scale fea-392 ture representations for pixels. A pixel itself does not 393 have enough contextual information for semantic seg-394 mentation, so it has to look around to check which class 395 it belongs to. Whether it is from texture (RGB) or depth 396 values (depth), the "board" object is very similar to sur-397 rounding pixels. The method of the forcible concatena-398 tion of the multi-scale feature would make some pixels 399 confused when determining which category they belong 400 to. Without the scale-aware module, the confidence 401 map on the board region is low as shown in Figure 4. 402 Also, the "wire" object (categorized into "otherprops") 403 is too thin to be classified. With the scale-aware module, 404 the model learns an appropriate feature representation. 405 When comparing with the baseline model that does not 406 have this module, the performance of our model is su-407 perior. To discover the importance of the self-attention 408 module, we provide the comparison results with and 409 without the module, as shown in Figure 6. From the 410 results, we can find that, through the self-attention mod-411 ule, our model can model long-range dependencies. 412

The feature extracted from the scale-aware module, we can find that the feature extracted from the different atrous rate ai which is focused on the different region on the input images, as shown in Figure 3. From the

	wall	floor	cabinet	bed	chair	sofa	table	door	window	bookshelf	picture	counter	blinds	desk	shelves	curtain	dresser	pillow	mirror
Song [43]	36.4	45.8	15.4	23.3	19.9	11.6	19.3	6.0	7.9	12.8	3.6	5.2	2.2	7.0	1.7	4.4	5.4	3.1	5.6
Liu [44]	37.8	48.3	17.2	23.6	20.8	12.1	20.9	6.8	9.0	13.1	4.4	6.2	2.4	6.8	1.0	7.8	4.8	3.2	6.4
Ren [45]	43.2	78.6	26.2	42.5	33.2	40.6	34.3	33.2	43.6	23.1	57.2	31.8	42.3	12.1	18.4	59.1	31.4	49.5	24.8
Li [ <mark>46</mark> ]	74.9	82.3	47.3	62.1	67.7	55.5	57.8	45.6	52.8	43.1	56.7	39.4	48.6	37.3	9.6	63.4	35.0	45.8	44.5
Cheng [10]	91.9	94.7	61.6	82.2	87.5	62.8	68.3	47.9	68.0	48.4	69.1	49.4	51.3	35.0	24.0	68.7	60.5	66.5	57.6
Wang [35]	-	-	-	-	-	-	_	-	-	-	-	_	-	_	_	-	-	_	_
Zhou [40]	-	-	-	-	-	-	_	-	-	-	-	_	-	_	_	-	-	_	_
Cao [42]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
SAMD	93.4	96.9	79.2	84.6	87.4	79.1	57.6	49.9	76.3	55.1	72.5	83.8	71.9	29.3	36.4	65.3	60.2	65.9	59.2
	1																		1
	floormat	clothes	ceiling	books	fridge	tv	paper	towel	shower	box	board	person	nightstanc	toilet	sink	lamp	bathhub	bag	mAcc
Song [43]	0.0	clothes	ceiling 35.8	pooks 6.1	6.6 fridge	<u>는</u> 0.7	baper 1.4	lowel 0.2	shower 0.0	xoq 0.6	poard 7.6	berson 0.7	1.7 nightstand	toilet 12.0	yus 15.2	lamp 6.0	1.1 bathhub	pag 0.0	mAcc 0.6
Song [43] Liu [44]	0.0 0.0	clothes 1.4 1.6	م تا 35.8 49.2	syooq 6.1 8.7	9.5 10.1	≥ 0.7 0.6	baber.	[] 0.2 0.2	0.0 0.0	ход 0.6 0.8	poard 5.6	berson 0.7 0.8	nightstane 1.7	toilet 12.0 14.9	чі <u>ў</u> 15.2 16.8	dung 0.9 1.2	pathhub 1:1	бар 0.6 1.3	Зүш 9.0 10.1
Song [43] Liu [44] Ren [45]	0.0 0.0 0.0	clothes 1.4 1.6 27.0	35.8 49.2 84.5	syooq 6.1 8.7 35.7	9.5 10.1 24.2	≥ 0.7 0.6 36.5	Laber 1.4 1.4 26.8	0.2 0.2 19.2	0.0 0.0 0.0 0.0	0.6 0.8 11.7	p. p. 7.6 8.6 51.4	U 0.7 0.8 35.7	nightstanc 1.7 1.8 25.0	tojet 12.0 14.9 64.1	<u>ущ</u> 15.2 16.8 53.0	due 0.9 1.2 44.2	qnµµnp 1.1 1.1 47.0	0.6 1.3 18.6	Э Чш 9.0 10.1 36.3
Song [43] Liu [44] Ren [45] Li [46]	U0000000000000000000000000000000000000	solution 1.4 1.6 27.0 28.4	35.8 49.2 84.5 68.0	syooq 6.1 8.7 35.7 47.9	9.5 10.1 24.2 61.5	≥ 0.7 0.6 36.5 52.1	Laber 1.4 1.4 26.8 36.4	0.2 0.2 19.2 36.7	0.0 shower 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.	xoq 0.6 0.8 11.7 38.1	poard 7.6 8.6 51.4 48.1	U 0.7 0.8 35.7 72.6	1.7 1.8 25.0 36.4	to 12.0 14.9 64.1 68.8	чц 15.2 16.8 53.0 67.9	duue 0.9 1.2 44.2 58.0	qnuutpaq 1.1 1.1 47.0 65.6	0.6 1.3 18.6 23.6	Э Чш 9.0 10.1 36.3 48.1
Song [43] Liu [44] Ren [45] Li [46] Cheng [10]	U 0.0 0.0 0.0 0.0 0.0 0.0	Sclothes 1.4 1.6 27.0 28.4 44.7	35.8 49.2 84.5 68.0 88.8	syooq 6.1 8.7 35.7 47.9 61.5	9.5 10.1 24.2 61.5 51.4	≥ 0.7 0.6 36.5 52.1 71.7	Laber 1.4 1.4 26.8 36.4 37.3	0.2 0.2 19.2 36.7 51.4	unit constraints and a constraints and a constraint constraints and a constraint constraint constraints and a constraint constraint constraint constraints and a constraint cons	0.6 0.8 11.7 38.1 46.0	p	Uosied 0.7 0.8 35.7 72.6 49.1	1.7 1.8 25.0 36.4 44.6	tion 12.0 14.9 64.1 68.8 82.2	<u>хи</u> 15.2 16.8 53.0 67.9 74.2	duer 0.9 1.2 44.2 58.0 64.7	qnuuteq 1.1 1.1 47.0 65.6 77.0	0.6 1.3 18.6 23.6 47.6	8 9.0 10.1 36.3 48.1 58.0
Song [43] Liu [44] Ren [45] Li [46] Cheng [10] Wang [35]	U0000000000000000000000000000000000000	solution 1.4 1.6 27.0 28.4 44.7	35.8 49.2 84.5 68.0 <b>88.8</b>	syooq 6.1 8.7 35.7 47.9 61.5	9.5 10.1 24.2 61.5 51.4	≥ 0.7 0.6 36.5 52.1 71.7 -	1.4 1.4 26.8 36.4 37.3	0.2 0.2 19.2 36.7 51.4	shower 0.0 9.0 0.0 2.9	0.6 0.8 11.7 38.1 46.0	7.6 8.6 51.4 48.1 54.2	0.7 0.8 35.7 72.6 49.1	1.7 1.8 25.0 36.4 44.6	12.0 14.9 64.1 68.8 82.2 -	× 15.2 16.8 53.0 67.9 74.2	0.9 1.2 44.2 58.0 64.7	qnythed 1.1 1.1 47.0 65.6 77.0 -	0.6 1.3 18.6 23.6 47.6	9.0 10.1 36.3 48.1 58.0 53.5
Song [43] Liu [44] Ren [45] Li [46] Cheng [10] Wang [35] Zhou [40]	0.0 0.0 5.6 0.0 0.0 - -	1.4 1.6 27.0 28.4 44.7 -	<sup>50</sup> 35.8 49.2 84.5 68.0 <b>88.8</b> -	6.1 8.7 35.7 47.9 61.5 –	9.5 10.1 24.2 61.5 51.4	≥ 0.7 0.6 36.5 52.1 71.7 -	1.4 1.4 26.8 36.4 37.3 -	0.2 0.2 19.2 36.7 51.4 -	0.0 0.0 9.0 0.0 2.9 -	0.6 0.8 11.7 38.1 46.0 -	7.6 8.6 51.4 48.1 54.2 -	0.7 0.8 35.7 72.6 49.1 –	1.7 1.8 25.0 36.4 44.6 -	12.0 14.9 64.1 68.8 82.2 –	<u>×iii</u> 15.2 16.8 53.0 67.9 74.2 –	0.9 1.2 44.2 58.0 64.7 -	qnyypeq 1.1 1.1 47.0 65.6 77.0 - -	0.6 1.3 18.6 23.6 47.6 –	З Ч 9.0 10.1 36.3 48.1 58.0 53.5 60.5
Song [43] Liu [44] Ren [45] Li [46] Cheng [10] Wang [35] Zhou [40] Cao [42]	U COUNTRE CONTREMENTE U COUNTRE CONTREMENTE U COUNTRE CONTREMENTE U COUNTRE CONTREMENTE U COUNTRE CONTREMENTE U COUNTRE CONTREMENTE U COUNTRE CONTRE	sequence of the second	<sup>50</sup> 35.8 49.2 84.5 68.0 <b>88.8</b> - -	6.1 8.7 35.7 47.9 61.5 - -	9.5 10.1 24.2 61.5 51.4 - -	≥ 0.7 0.6 36.5 52.1 71.7 - -	1.4 1.4 26.8 36.4 37.3 - -	0.2 0.2 19.2 36.7 51.4 - -	0.0 0.0 9.0 0.0 2.9 - -	0.6 0.8 11.7 38.1 46.0 - -	7.6 8.6 51.4 48.1 54.2 - -	0.7 0.8 35.7 72.6 49.1 - -	1.7 1.8 25.0 36.4 44.6 - -	12.0 14.9 64.1 68.8 82.2 - - -	ту 15.2 16.8 53.0 67.9 74.2 – – –	0.9 1.2 44.2 58.0 64.7 -	qnyypeq 1.1 1.1 47.0 65.6 77.0 - - - -	0.6 1.3 18.6 23.6 47.6 - -	9.0 10.1 36.3 48.1 58.0 53.5 <b>60.5</b> 58.5

Table 2: Performance on the SUN RGB-D dataset. The SAMD row shows the results of our SAMD model. The class of background is ignored during performance evaluation.

Table 3: Performance on different feature extractor encoder backbone network of our model.

Backbone	w/o SAMD	w/ SAMD
ResNet50	45.1	48.1
ResNet101	46.7	52.3

Table 4: Performance on different modality fusion methods of our model.

Methods	mIoU
Late fusion [7]	48.9
Gated fusion [10]	51.3
Modality-aware fusion	51.9

figure, we also find that the scale feature extracted from

the scale-aware module, has different levels of attention on each modality. This phenomenon spurs us to design

the next modality-aware module.

To demonstrate the effectiveness of the modality-421 aware module, we provide three results on the NYU-422 Depth v2 dataset with different modality fusion meth-423 ods as shown in Table 4. In all three experiments, they 424 all include the scale-aware and self-attention modules. 425 All parameter settings in the experiments are the same 426 427 except for the fusion method used. The late fusion approach follows the instruction in [7], which fuses RGB 428 and depth networks by equal-weight score. [10] pro-429 poses a gated fusion way to fuse RGB and depth by 430

Table 5: Performance on the NYU-Depth v2 test dataset (4-class).

	Acc	mAcc
Courprie [5]	64.5	63.5
Hermans [47]	69.0	68.1
Stuckler [48]	70.6	66.8
Wang [49]	-	74.7
Eigen [50]	83.2	82.0
He [37]	83.6	82.5
SAMD	86.9	85.7

To demonstrate that the depth estimation module is workable and useful, we provide the depth estimation results of input images, as shown in Figure 7. From the results, we can find that the depth estimation can provide a plausible depth value for the input image. 439

To have a better understanding of how the proposed440SAMD model outperforms the baseline method, we pro-441vide the visualization results of the improvement of IoU442for each semantic category in Figure 8. As can be seen443from the statistics result in Figure 8, our SAMD is superior to the baseline in most classes.444

In Table 6, we give the quantitative comparisons of 446 with and without our SAMD components on the NYU-

Table 6: Ablation study of the proposed SAMD model on NYU-Depth v2 dataset. S, A, M and D denote scale-aware module, self-attention module, modality-aware module and depth estimation module, respectively.

	Methods	mIoU
а.	Baseline	46.7
b.	Baseline + $S$	47.8
с.	Baseline + $A$	48.6
d.	Baseline + $M$	48.4
е.	Baseline + $S + A$	49.8
f.	Baseline + $S + M$	49.9
g.	Baseline $+ A + M$	49.5
h.	Baseline + $S + A + M$	51.9
i.	Baseline + $S + A + M + D$	52.3

Figure 8: Performance Analysis. Per-class IoU improvement of our SAMD model over baseline on NYU Depth-v2 test dataset.



Depth v2 dataset. From the comparison results ( $b \sim$ 448 i), each component in the proposed SAMD module will 449 benefit the performance of the indoor semantic segmen-450 tation. The qualitative results are illustrated in Figure 5, 451 it gives the visualized comparisons with and without our 452 SAMD module on the NYU-Depth v2 dataset. In Table 453 454 1, we give the results of the comparison between our model and state-of-the-art methods on the NYU-Depth 455 v2 dataset. From the results, we can find that our model 456 is better than state-of-the-art methods in many classes. 457 We also test our model on the SUN RGBD dataset, and 458 we obtain a state-of-the-art comparable result, 63.4% 459 mean accuracy, more detail please refer to Table 2. 460

We compare SAMD to other state-of-the-art methods on the 4-class of the NYU-Depth v2 dataset, and 462 the quantitative results are shown in Table 5. We also 463 check the model size saved by PyTorch for both base-464 465 line and our SAMD to demonstrate the proposed module wouldn't increase the parameter size of the baseline 466 too much. The size of baseline is 127.28M, and our 467 SAMD is 136.19M. 468

## 5. Conclusion

In this paper, we propose SAMD to tackle the chal-470 lenging problems for indoor semantic segmentation 471 with RGB-D data. SAMD is composed of three main 472 parts: (1) the scale-aware module which is designed for 473 generating a spatial-sampled and scale-sampled feature 474 representation, (2) the modality-aware module which 475 can weigh the varying contributions of the two comple-476 mentary modalities for better fusion, and (3) the self-477 attention module and depth estimation module, which 478 can produce long-range dependencies for better model-479 ing and push the RGB network to extract more plausible 480 features. Theoretical analysis, qualitative and quantita-481 tive experimental results on NYU-Depth v2 and SUN 482 RGB-D dataset demonstrate that SAMD can achieve 483 significant performance gains for indoor semantic seg-484 mentation. 485

## Acknowledgements

The authors thank the editor and all the reviewers 487 for their very helpful comments to improve this paper. 488 The authors also thank Professor Ou wu and Professor 489 Xukun Shen for their suggestions and to polish this pa-490 per. 491

## References

- [1] Q. Xie, O. Remil, Y. Guo, M. Wang, M. Wei, J. Wang, Object 493 detection and tracking under occlusion for object-level RGB-D 494 video segmentation, IEEE Transactions on Multimedia 20 (3) 495 (2017) 580-592. 496
- [2] J. Huang, Z. Liu, Y. Wang, Joint scene classification and seg-497 mentation based on hidden Markov model. IEEE Transactions 498 on Multimedia 7 (3) (2005) 538-550.
- X. Liang, L. Lin, W. Yang, P. Luo, J. Huang, S. Yan, Clothes 500 co-parsing via joint image segmentation and labeling with ap-501 plication to clothing retrieval, IEEE Transactions on Multimedia 502 18 (6) (2016) 1175-1186. 503
- [4] I. Ahn, C. Kim, Face and hair region labeling using semisupervised spectral clustering-based multiple segmentations, IEEE Transactions on Multimedia 18 (7) (2016) 1414-1421.
- [5] C. Couprie, C. Farabet, L. Najman, Y. LeCun, Indoor semantic segmentation using depth information, arXiv preprint arXiv:1301.3572.
- [6] S. Gupta, R. Girshick, P. Arbeláez, J. Malik, Learning rich fea-510 tures from RGB-D images for object detection and segmenta-511 tion, in: ECCV, 2014. 512
- J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: CVPR, 2015.
- [8] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs, PAMI.
- [9] X. Q. X. W. J. J. Hengshuang Zhao, Jianping Shi, Pyramid scene 519 parsing network, in: CVPR, 2017. 520

10

492

499

504

505

506

507

508

509

513

514

515

516

517

518

486

- [10] Y. Cheng, R. Cai, Z. Li, X. Zhao, K. Huang, Locality-sensitive 58621 deconvolution networks with gated fusion for RGB-D indoor semantic segmentation, in: CVPR, 2017. 58**5**23
- X. Li, H. Zhao, L. Han, Y. Tong, S. Tan, K. Yang, Gated fully [11] fusion for semantic segmentation, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 11418-11425
- [12] X. Li, H. He, X. Li, D. Li, G. Cheng, J. Shi, L. Weng, Y. Tong, 528 Z. Lin. Pointflow: Flowing semantics through points for aerial 529 image segmentation, arXiv preprint arXiv:2103.06564. 530

525

526

527

534

536

537

- [13] L.-C. Chen, G. Papandreou, F. Schroff, H. Adam, Rethink-531 ing atrous convolution for semantic image segmentation, arXiv 532 preprint arXiv:1706.05587. 533
- X. Chen, K.-Y. Lin, J. Wang, W. Wu, C. Qian, H. Li, [14] G. Zeng, Bi-directional cross-modality feature propagation with 535 separation-and-aggregation gate for RGB-D semantic segmentation, arXiv preprint arXiv:2007.09183.
- 538 [15] D. Lin, Y. Ji, D. Lischinski, D. Cohen-Or, H. Huang, Multiscale context intertwining for semantic segmentation, in: ECCV, 539 540 2018
- [16] C. Liang-Chieh, G. Papandreou, I. Kokkinos, K. Murphy, 541 A. Yuille, Semantic image segmentation with deep convolu-542 tional nets and fully connected CRFs, in: ICLR, 2015. 543
- 544 L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-decoder with atrous separable convolution for semantic 545 image segmentation, arXiv preprint arXiv:1802.02611. 546
- 547 [18] X. Li, Z. Jie, W. Wang, C. Liu, J. Yang, X. Shen, Z. Lin, Q. Chen, S. Yan, J. Feng, Foveanet: Perspective-aware urban 548 scene parsing, arXiv preprint arXiv:1708.02421. 549
- [19] S. Kong, C. Fowlkes, Recurrent scene parsing with perspective 550 understanding in the loop, arXiv preprint arXiv:1705.07238. 551
- W. Wang, T. Zhou, F. Yu, J. Dai, E. Konukoglu, L. Van Gool, 552 [20] Exploring cross-image pixel contrast for semantic segmentation, 553 554 in: ICCV. 2021.
- [21] J. Hu, L. Shen, G. Sun, Squeeze-and-excitation networks, arXiv 555 preprint arXiv:1709.01507. 556
- 557 [22] D. G. Lowe, Distinctive image features from scale-invariant keypoints, International journal of computer vision 60 (2) (2004) 558 559 91-110.
- [23] D. Bahdanau, K. Cho, Y. Bengio, Neural machine transla-560 tion by jointly learning to align and translate, arXiv preprint 561 arXiv:1409.0473. 562
- 563 [24] K. Gregor, I. Danihelka, A. Graves, D. J. Rezende, D. Wierstra, Draw: A recurrent neural network for image generation, arXiv 564 preprint arXiv:1502.04623. 565
- [25] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, 566 R. Zemel, Y. Bengio, Show, attend and tell: Neural image cap-567 tion generation with visual attention, in: ICML, 2015. 568
- [26] Z. Yang, X. He, J. Gao, L. Deng, A. Smola, Stacked attention 569 570 networks for image question answering, in: CVPR, 2016.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. 571 [27] Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, in: 572 NIPS 2017 573
- [28] H. Zhang, I. Goodfellow, D. Metaxas, A. Odena, Self-574 attention generative adversarial networks, arXiv preprint 575 arXiv:1805.08318. 576
- [29] H. Hu, J. Gu, Z. Zhang, J. Dai, Y. Wei, Relation networks for 577 object detection, in: CVPR, 2018. 578
- T. Zhou, S. Wang, Y. Zhou, Y. Yao, J. Li, L. Shao, Motion-579 [30] attentive transition for zero-shot video object segmentation, in: 580 Proceedings of the AAAI Conference on Artificial Intelligence. 581 Vol. 34, 2020, pp. 13066-13073. 582
- L. Hoyer, D. Dai, Y. Chen, A. Koring, S. Saha, L. Van Gool, [31] 583 584 Three ways to improve semantic segmentation with selfsupervised depth estimation, in: CVPR, 2021, pp. 11130-585

11140.

- 58822 [32] V. Guizilini, R. Hou, J. Li, R. Ambrus, A. Gaidon, Semanticallyguided representation learning for self-supervised monocular depth, arXiv preprint arXiv:2002.12319. 58\$24
  - [33] X. Wang, R. Girshick, A. Gupta, K. He, Non-local neural networks, in: CVPR, 2018.

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

- [34] C. Godard, O. Mac Aodha, M. Firman, G. J. Brostow, Digging into self-supervised monocular depth estimation, in: ICCV, 2019, pp. 3828-3838.
- [35] W. Wang, U. Neumann, Depth-aware CNN for RGB-D segmentation, in: ECCV, 2018.
- [36] Z. Deng, S. Todorovic, L. Jan Latecki, Semantic segmentation of RGBD images with mutex constraints, in: ICCV, 2015.
- [37] Y. He, W.-C. Chiu, M. Keuper, M. Fritz, STD2P: RGBD semantic segmentation using spatio-temporal data-driven pooling, in: CVPR, 2017.
- [38] D. Seichter, M. Köhler, B. Lewandowski, T. Wengefeld, H.-M. Gross, Efficient RGB-D semantic segmentation for indoor scene analysis, arXiv e-prints (2020) arXiv-2011.
- [39] Z. Gu, L. Niu, H. Zhao, L. Zhang, Hard pixel mining for depth privileged semantic segmentation, IEEE Transactions on Multimedia.
- [40] H. Zhou, L. Qi, Z. Wan, H. Huang, X. Yang, RGB-D coattention network for semantic segmentation, in: ACCV, 2020.
- D. Lin, H. Huang, Zig-zag network for semantic segmentation of RGB-D images, IEEE Transactions on Pattern Analysis and Machine Intelligence 2019 42 (10) (2020) 2642-2655
- [42] J. Cao, H. Leng, D. Lischinski, D. Cohen-Or, C. Tu, Y. Li, Shapeconv: Shape-aware convolutional layer for indoor RGB-D semantic segmentation, in: ICCV, 2021, pp. 7088-7097
- [43] S. Song, S. P. Lichtenberg, J. Xiao, SUN RGB-D: A RGB-D scene understanding benchmark suite, in: CVPR, 2015.
- [44] C. Liu, J. Yuen, A. Torralba, Sift flow: Dense correspondence across scenes and its applications, IEEE transactions on pattern analysis and machine intelligence 33 (5) (2010) 978-994.
- [45] X. Ren, L. Bo, D. Fox, RGB-(D) scene labeling: Features and algorithms, in: CVPR, IEEE, 2012, pp. 2759-2766.
- [46] Z. Li, Y. Gan, X. Liang, Y. Yu, H. Cheng, L. Lin, LSTM-CF: Unifying context modeling and fusion with LSTMs for RGB-D scene labeling, in: ECCV, 2016.
- [47] A. Hermans, G. Floros, B. Leibe, Dense 3d semantic mapping of indoor scenes from RGB-D images, in: ICRA, IEEE, 2014, pp. 2631-2638.
- J. Stückler, B. Waldvogel, H. Schulz, S. Behnke, Dense real-[48] time mapping of object-class semantics from RGB-D video, Journal of Real-Time Image Processing 10 (4) (2015) 599-609.
- [49] J. Wang, Z. Wang, D. Tao, S. See, G. Wang, Learning com-632 mon and specific features for RGB-D semantic segmentation 633 with deconvolutional networks, in: ECCV, Springer, 2016, pp. 634 664-679 635
- [50] D. Eigen, R. Fergus, Predicting depth, surface normals and se-636 mantic labels with a common multi-scale convolutional archi-637 tecture, in: ICCV, 2015, pp. 2650-2658. 638