

Article

Mask-Aware Semi-Supervised Object Detection in Floor Plans

Tahira Shehzadi^{1,2,3}, Khurram Azeem Hashmi^{1,2,3}, Alain Pagani³, Marcus Liwicki⁴, Didier Stricker^{1,3} and Muhammad Zeshan Afzal^{1,2,3}

¹ Department of Computer Science, Technical University of Kaiserslautern, 67663 Kaiserslautern, Germany; Tahira.Shehzadi@dfki.de (T.S.); khurram_azeem.hashmi@dfki.de (K.A.H.); muhammad_zeshan.afzal@dfki.de (M.Z.A.); didier.stricker@dfki.de (D.S.)

² Mindgarage, Technical University of Kaiserslautern, 67663 Kaiserslautern, Germany

³ German Research Institute for Artificial Intelligence (DFKI), 67663 Kaiserslautern, Germany; alain.pagani@dfki.de

⁴ Department of Computer Science, Luleå University of Technology, 971 87 Luleå, Sweden; marcus.liwicki@ltu.se

* Correspondence: Tahira.Shehzadi@dfki.de

Abstract: Research has been growing on object detection using semi-supervised methods in past few years. We examine the intersection of these two areas for floor-plan objects to promote the research objective of detecting more accurate objects with less labelled data. The floor-plan objects include different furniture items with multiple types of the same class, and this high inter-class similarity impacts the performance of prior methods. In this paper, we present Mask R-CNN based semi-supervised approach that provides pixel-to-pixel alignment to generate individual annotation masks for each class to mine the inter-class similarity. The semi-supervised approach has a student-teacher network that pulls information from the teacher network and feeds it to the student network. The teacher network uses unlabeled data to form pseudo-boxes, and the student network uses both unlabeled data with the pseudo boxes and labelled data as ground truth for training. It learns representations of furniture items by combining labelled and unlabeled data. On the Mask R-CNN detector with ResNet-101 backbone network, the proposed approach achieves mAP of 98.8%, 99.7%, and 99.8% with only 1%, 5% and 10% labelled data, respectively. Our experiment affirms the efficiency of the proposed approach as it outperforms the fully supervised counterpart using only 10% of the labels.

Keywords: object detection; semi-supervised learning; Mask R-CNN; floor-plan images; computer vision.

1. Introduction

Semi-supervised learning-based research getting more attention in the past few years as it can use unlabeled data to increase model performance when it's impossible to annotate large datasets. The first layout of semi-supervised approach-based learning uses consistency-based self-learning [1–8] approaches. The main idea is to create artificial labels and then predict those self-generated labels by training the model on unlabeled data with stochastic augmentations. Those self-generated labels can be the network's predictive distribution or one-hot prediction. The second point of improvement in semi-supervised approach-based learning is the variety of available data augmentation techniques. Data augmentation techniques boost the performance of the training network [9,10] and are also efficient for consistency-based learning [5–8]. The augmentation approaches progress from image transformation such as cropping, flipping, scaling, brightness, colour augmentation, contrast, saturation, translation, and rotation to image generation [11–13] and model training by reinforcement-learning [14,15]. Previously, the researchers applied supervised learning techniques for floor-plan object detection. We use the semi-supervised approach for floor-plan analysis, which matches the fully supervised counterpart using only 10% of the label data.

The floor-plan object detection problem has high value because of its usage in tremendous applications such as property value estimation, furniture setting and designing, etc.

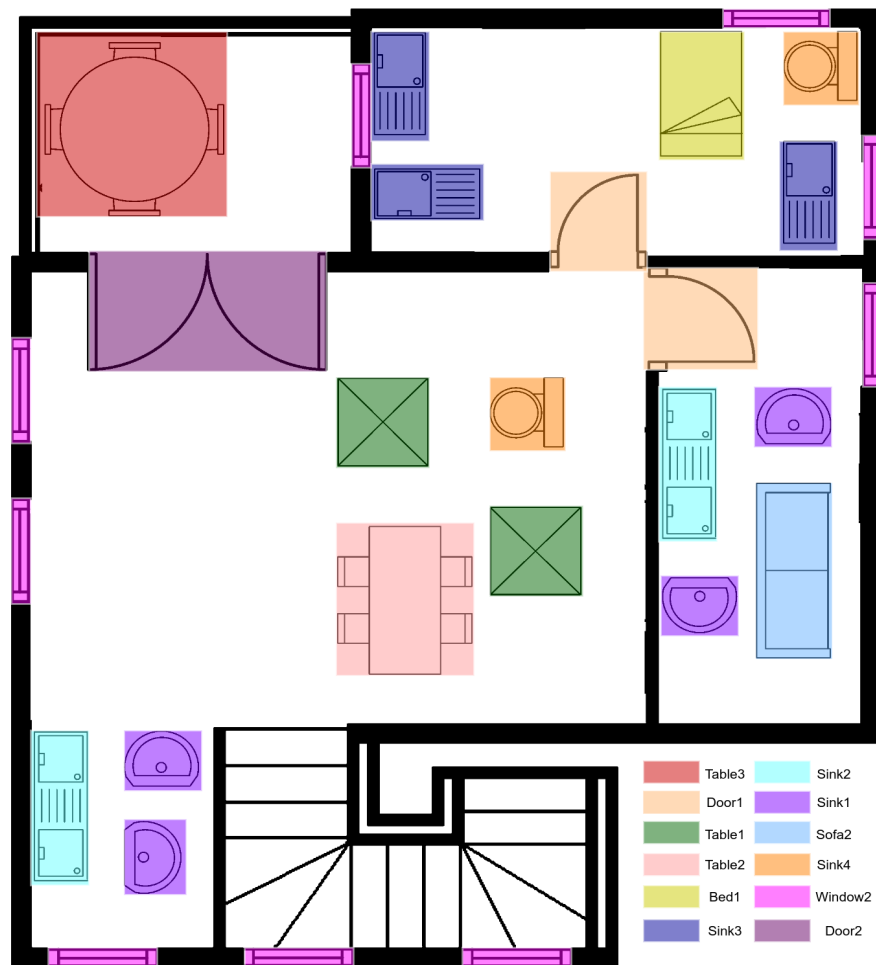


Figure 1. The sample image of the floor-plan dataset contains furniture items such as sink and bed items in the top-right room, while the top-left room contains a dining table. The bottom right corner of the image shows labels of furniture items.

The floor-plan objects include furniture items, windows, doors, and walls. Humans can readily recognize floor-plan objects, but to automatically recognize and detect floor-plan objects is challenging because of the similarity between room types and furniture items. For example, the Drawing room contains a limited number of furniture items, and the furniture category of the kitchen and dining room is almost similar. There are many applications of floor-plan object detection, such as 3d reconstruction of floor-plan [16] and similarity search [17]. Floor-plan object detection is necessary for the floor-plan analysis applications. Figure 1 is an overview of the floor-plan layout with different furniture items that explains room size and furniture categories. The top left room is the dining room, where a single round table is present. The top right room contains a kitchen with a bathroom. The next room is a living area where different sofa items are present. Thus, all other rooms have names according to their furniture items. This floor-plan category can help furniture installation [16].

Semi-supervised approach-based object detection needs a small amount of labelled data with unlabeled data. There are some multi-stage approaches [18,19] that use label data for training in the first stage, followed by unlabeled data for generating pseudo labels, then retraining on unannotated data. The model performance depends on the accuracy of the generated pseudo label, but the available training data is small, which reduces model efficiency. To increase label data, we generate pseudo labels using a semi-supervised approach and then use these pseudo labels and small portions like 1% of label data to

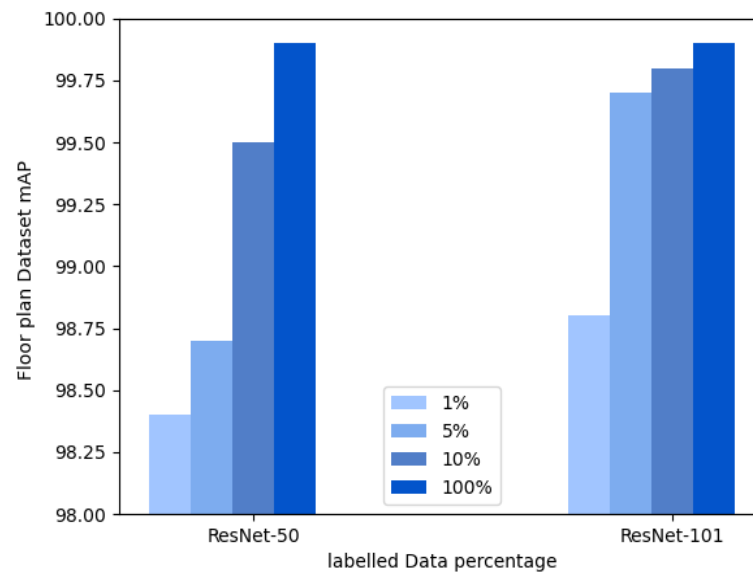


Figure 2. Compares the performance of ResNet-50 [24] with the ResNet-101 backbone network using Mask R-CNN [22] framework under different label data settings. The increasing colour of bars indicates the increase in the percentage of label data.

train the model. We randomly sample unlabeled and labelled data in which both portions include all classes present in available data. We used two models for our experiment on the floor-plan dataset, the first is for detector training, and the second is for generating pseudo labels for unlabeled data. This approach provides simplified multi-stage training. Further, it uses the flywheel effect [20], in which the pseudo label generator and training detector can boost each other to improve model performance with increasing training iterations. Another important benefit of this approach is that more weight is provided to the pseudo-label generator model rather than the training detector model as it guides the training model, instead of providing hard category labels as in earlier techniques [18,19]. This approach is also proposed in Soft-Teacher model [21]. In this network, the teacher model uses a pseudo-label generator, and the student model uses the training detector.

We detected objects on pixel level like different furniture items in floor-plan data using a semi-supervised approach. We used Mask R-CNN [22] with Feature Pyramid Network (FPN) [23] as a detection method, ResNet-50 [24] and ResNet-101 network pre-trained on the ImageNet dataset [25] as a backbone with student-teacher approach. We used 1%, 5% and 10% floor-plan images as labelled data for training and the rest of the images as unlabeled data. We used 5 data folds for each percentage level and calculated the final performance by taking the average of these five folds. Figure 2 compares the performance of both these backbones under the different percentage of label data settings. We obtain 98.8(%) mAP, 99.7(%) mAP and 99.8(%) mAP on Mask R-CNN [22] with ResNet-101 [24] backbone network for 1%, 5% and 10% floor-plan images, respectively. This paper provides an end-to-end Semi-supervised approach-based object detection in the floor-plan domain. The main contribution of this work is as follows:

- We present Mask R-CNN [22] based semi-supervised trainable network with ResNet-50 [24] and ResNet-101 backbone network for object detection in floor-plan domain.
- The Mask R-CNN [22] based semi-supervised approach improves state-of-the-art performance on the publicly available floor-plan dataset named SESYD [26] and SFPI [27] using only 10% of the labels.

The remaining paper is arranged as follows. Section 2 talks about previous research on semi-supervised approaches based-learning and floor-plan datasets. Section 3 explain the

methodology and Section 4 discusses the dataset briefly. Section 5 is about experimental setup. Section 6 discusses the evaluation matrices. In Section 7, we analyze the experimental results. Finally, Section 8 summarizes the experimental work and gives an idea about future directions.

2. Related Work

Object detection and semi-supervised learning are essential steps toward floor-plan image analysis. This section gives an overview of previous work in these domains and contains three parts. The first section describes the literature about object detection. The second section explains previous semi-supervised approaches. Finally, we explain the literature on the floor-plan domain.

2.1. Object Detection and its Applications

Object detection is the main computer-vision domain in which extensive work has been done in the past few years [22,28–36]. There are two main types: single-stage detectors [33,35,37] and two-stage detectors [30,32,38–40]. The two-stage detectors extract the object regions in the first stage and then classify and localize the object in the second stage. These detectors, such as Faster R-CNN [30], firstly generate region proposals, making a separate prediction for every object in the image. In contrast, single-stage detectors perform classification and localization in one pass through the neural network. The basic difference between these detectors is the cascade filter for object proposals. These detectors provide good results on a large amount of label data and used in different applications in many elds, such as document image analysis [41–46] face detection [47] and pedestrian detection [48].

2.2. Semi-Supervised Learning

Semi-supervised based image classification has two types: pseudo-label-based learning and consistency-based learning. The consistency-based learning [49–51] examine the similarity between original and augmented images. It provides more weight to unlabeled data than labelled data, which helps in perturbations of the same image for producing similar labels. There are different methods to apply perturbations using noise [49], augmentation [51], adversarial training [50]. In [52], the author predicted the training steps to assemble the training object. In [2], the author takes the weighted average by ensembling rather than predicting the model, called exponential mean-average (EMA). In [7,53], the authors annotated the unlabeled images with pseudo labels using the classification model and then retrained the detector using this pseudo label data. They analyzed the effect of data augmentation for semi-supervised learning [6,8,54,55].

Semi-supervised based object detection has two types as pseudo-label based learning [18,19,56–58] and consistency-based learning [59,60]. In [19,56], labels generated from different augmented images are ensembled to predict labels of unlabeled images. In [57], pseudo-labels are generated by training SelectiveNet [61]. In [58], the labelled image contains the detected box of the unlabeled image, and the author calculated the localization consistency estimation for the attached label image. It needs a deep detection procedure [58] as the image itself is changed. Recently, intricate augmentation approaches, including CTAugment [5], RandAugment [62] are proven to be very effective for Semi-supervised learning on object detection [6,8].

2.3. Floor-Plan Analysis

Research on object detection in floor-plan data is growing because of its usage in tremendous applications such as property value estimation, furniture setting and designing, etc. Ghorbel et al.[63] proposed a handwritten floor-plan recognition model. This network provides a CAD model for floor-plans data. In [64], the author proposed room detection model for floor-plan dataset. Moreover, [65] proposed a model for understanding floor-plan using Hough-transform and subgraph-isomorphism. Several graphic recognition

methods are applied to identify the basic structure and also consider human feedback during the analysis phase.

In [66], the author used a deep learning network to parse floor-plan images. The author applied Cascade Mask R-CNN [32] to get floor-plan information and keypoint-CNN for segmentation to extract accurate corner locations and obtained the final segmentation results after post-processing. In [67], textural information is extracted from floor-plan images. This work is helpful for visually impaired people to analyze house design and for customers to buy a house online. The morphological closure is applied to detect the walls of the floor-plan image, the flood fill method to detect corners, and scale-invariant features for door identification. After extracting all this information, the author applied text synthesis techniques.

In [68], the author proposed an object recognition method for floor-plan images. The main target is to recognize floor-plan items like windows, walls, rooms, doors and furniture items. To extract features, the VGG network [69] is used. It recognizes room types based on furniture items present in the room. But room type identification is not showing good results as variation in furniture items is less. It also detects room boundaries for doors, windows and walls, which gives good results.

Liu et al. [70] detected edges in the floor-plan dataset using the deep network and then used Integer programming to detect walls of different rooms by combining those corner points. However, this approach can only recognize walls of rectangular rooms with uniform thickness; it works on the Manhattan assumption that aligns the walls with two main axes in a floor-plan image. Yamasaki et al. [71] applied a fully convolutional network (FCN) to label pixels for detecting similar structure houses by forming a graph model of the floor-plan dataset with different classes. Their method ignores spatial relation between different classes as it detects pixels of different classes separately by using a simple segmentation network.

In [72], Faster R-CNN [30] is used to detect kitchen items like stoves, sliding doors, simple doors and bathtubs, and then adopted fully convolutional network (FCN) to detect walls pixels. They also estimated the size of different rooms by recognizing text using a library tool. Maće et al. [64] used Hough transform to identify doors and walls in floor-plan images. In [73], the author used a pixel-based segmentation approach to detect doors, walls, windows, and the bag-of-words (BOW) network to classify image patches. They trained these patches to generate graphs for detecting walls. The author detected the walls in [16] by recognizing parallel lines, determined the room size by calculating the distance between parallel lines and estimated the wall thickness by clustering distance value.

3. Method

The experiment is performed on Mask R-CNN [22] with ResNet-50 [24] and ResNet-101 backbone. We used this model with convolutional networks (CNN) and a student-teacher network. Figure 4 shows the whole pipeline of the used approach. In this section, we explain individual modules of the experiment.

3.1. Mask R-CNN

Mask R-CNN [22] is an extended version of Faster R-CNN [30] with a new branch for providing masks to the detected objects with the two already present branches for the classification and regression layer. This branch is applied on RoIs to deal with detection on the pixel level to segment each stance accurately. The basic architecture of Mask R-CNN is identical to Faster R-CNN as it uses a similar architecture to generate object proposals. The major difference is that Mask R-CNN uses an RoI-align layer rather than an RoI-Pooling layer to reduce misalignment on pixel level because of spatial quantization. Generally, training of Mask R-CNN [22] and Faster R-CNN [30] is identical. For accuracy and speed, we prefer ResNeXt-101 [74] as backbone with feature Pyramid Network (FPN) [75]. We create the mask for each class for pixel-level classification to reduce interclass similarity. We created ground truth for the mask using object width, height and bounding

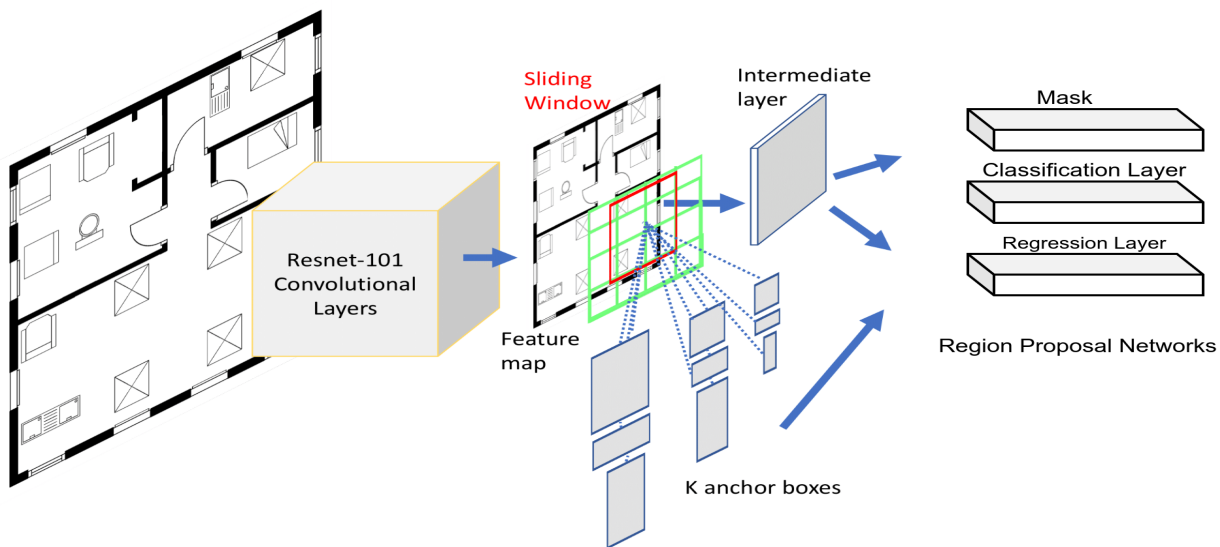


Figure 3. The overview of Mask R-CNN [22] framework with ResNet-101 [24] backbone for floor-plan object detection. This network gets a convolution feature map from the backbone layer, provides anchors generated by a sliding window and predicts the regions by Region-Proposal Network (RPN). Then, we implement a pooling process to resize and a Fully connected layer to produce three nodes as a mask, softmax classification and bounding-box regression.

box coordinates. The masks of all classes are in square boxes having four corner points as $(x_{min}y_{min}, x_{max}y_{min}, x_{max}y_{max}, x_{min}y_{max})$. Where (x_{min}, y_{min}) is the first corner point of the mask and obtained other corner points by adding the width and height of the bounding box in the first corner point. The model learns the mask of each class separately and is defined as average binary cross-entropy loss, as shown in the following equation.

$$L_{mask} = -\frac{1}{M^2} \sum_{1 \leq l, m \leq M} y_{lm} \log y_{lm}^n + (1 - y_{lm}) \log(1 - y_{lm}^n) \quad (1)$$

Where y_{lm} is the label of pixel (l, m) in true mask area $M * M$ and y_{lm}^n is the estimated value of the same pixel for the ground-truth class n . The loss function of Mask R-CNN [22] is the combination of localization, classification and segmentation mask loss, where classification and localization loss is the same as in Faster R-CNN [30].

3.2. Backbone Network

The model performance drops both for train and test data. This reduction is not because of overfitting. Instead, network initialization, exploding or vanishing gradients can also cause this problem. These can be easily optimized compared to the plain network, whose training error increases with adding more layers. The ResNet-50 [24] network is formed by replacing 2-layer block of resnet-34 with 3-layer block. This network has high accuracy than the resnet-34 network. The ResNet-101 contains three more layers. We used ResNet-50 [24] and ResNet-101 backbone network for this semi-supervised experiment. Figure 3 explains the Mask R-CNN [22] framework with ResNet-101 backbone.

3.3. Semi-Supervised Model

The performance of the model is dependent on the quality of pseudo labels. Setting a high threshold on foreground value to get more student-created boxes can provide better results than a low threshold. We get the best results when the threshold value is 0.9. However, a high threshold value provides good foreground precision, and the recall of box-candidate decreases quickly. Suppose we apply IoU between teacher-created pseudo-boxes and student-created box-candidate to provide background and foreground labels

as an ordinary object detection model does. In that case, we incorrectly classified some foreground boxes as negative, which reduces performance.

To eliminate this problem, we use the student-teacher network to generate pseudo-labels using a semi-supervised approach and then use these pseudo labels as well as a small portion like 1% of label data to train the model. This unlabeled and labelled data sampling includes all classes present in available data. The random samples of labelled and unlabelled images are selected using sampling ratio s_r to make training batches. The teacher model uses unlabeled data to form pseudo-boxes, and the student model uses both unlabeled data with the pseudo boxes and labelled data as ground truth for training. We assessed the reliability of student-created box candidates of a real background and used it to weigh background-class loss. The total loss is the combination of unsupervised and supervised loss:

$$L = L_{sup} + \alpha L_{un} \quad (2)$$

Where L_{sup} represents the supervised loss of labelled data while L_{un} represents the unsupervised loss of unlabeled data, α is the controlling factor of unsupervised loss. We normalized these losses by their respective amount of floor-plan images in the training batch. The supervised and unsupervised loss is the combination of classification, localization and segmentation mask loss. Equation 1 explains the mask loss while classification and localization loss is the same as in Faster R-CNN [30].

$$L_{sup} = \frac{1}{N_b} \sum_{n=1}^{N_b} (L_{class}(I_b^n) + L_{rg}(I_b^n) + L_{mask}(I_b^n)) \quad (3)$$

$$L_{un} = \frac{1}{N_u} \sum_{n=1}^{N_u} (L_{class}(I_u^n) + L_{rg}(I_u^n) + L_{mask}(I_u^n)) \quad (4)$$

where I_b^n represents n-th labeled-image, I_u^n represents n-th unlabeled-image, N_b indicates total labeled-images, N_u indicates total unlabeled-images, L_{class} , L_{rg} and L_{mask} is the classification, regression and mask loss, respectively.

Figure 4 explains the overall architecture of the student-teacher approach. We initialized the teacher and student model randomly to start training; then, the student model updates the teacher model just like [8,76] using the exponential moving average(EMA) approach. Generating pseudo-labels for detecting objects is more challenging than classifying objects, as an image typically has multiple objects. To annotate those objects, we need location and category. The teacher model gets unlabeled images to detect objects and generate many bounding boxes. The non-maximum suppression(NMS) is applied to minimize redundant boxes generated on the image objects. Even though we eliminated most iterating boxes, some non-foreground boxes remain.

The FixMatch [8] is a supervised learning-based image classification approach used to get better pseudo boxes and speed up the student network training. We applied weak-augmentation for generating pseudo-labels by the teacher network and strong-augmentation for training the student network. To calculate the reliability score is a little bit difficult. So, we used the background value generated by the teacher model using weak augmentation as a signal for the student model. This approach is just like simple negative-mining, not like OHEM [36,77] or Focal Loss [36], known as hard negative-mining.

To measure the consistency of regression boxes, we used a box jittering approach that calculates the localization reliability of pseudo boxes. We calculated jittered box around box candidate b_k and fed it into the teacher model to get refined box \hat{b}_k as follows:

$$\hat{b}_k = filtered(jitter(b_k)) \quad (5)$$

We repeated this process many times to get N_{jitter} filtered jitter boxes. The location probability of an object as regression-variance is determined as follows:

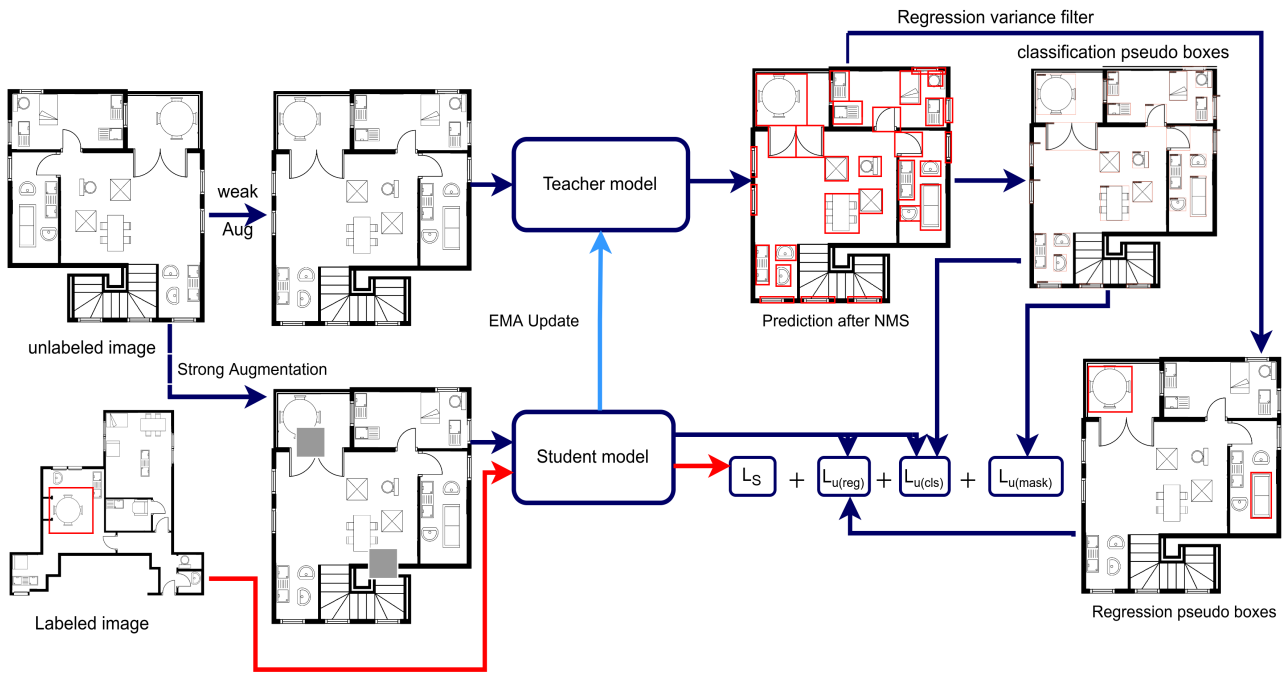


Figure 4. The complete architecture of the semi-supervised approach. We randomly initialized the teacher and student model to start training; then updated these models with Exponential Moving Average (EMA) approach represented with a light blue arrow. The teacher model has weak augmented unlabeled images, while the student model has strong augmented unlabeled images as input to detect objects, generating many bounding boxes. The student model also takes a small percentage of label images as additional input. The red arrows show supervised, and dark blue arrows show unsupervised training. The total loss is the combination of unsupervised and supervised loss.

$$\bar{\sigma}_k = \frac{1}{4} \sum_{n=1}^4 (\hat{\sigma}_n) \quad (6)$$

$$\hat{\sigma}_n = \frac{\sigma_n}{0.5(h(b_k) + w(b_k))} \quad (7)$$

where $\hat{\sigma}_n$ is the normalization of σ_n , σ_n is standard-derivation of n th coordinate of filtered jittered boxes, $w(b_k)$ is the width and $h(b_k)$ is the height of jittered box b_k .

The localization accuracy will be more when the regression variance of the box is smaller. But it is not feasible to assess the regression-variance of box candidates during the training process. So we compute reliability only for those boxes whose foreground value is above 0.5, reducing the number of boxes from hundreds to 16 per image, minimizing the computational cost.

4. Dataset

We need a large dataset with various floor-plan layouts for deep neural network training, and there should be enough classes to analyze variation in furniture items. The dataset is created from SESYD [26] named SFPI (Synthetic Floor-Plan Images). It contains 16 furniture classes as window1, sofa1, sink1, table1, door1, window2, sofa2, sink2, table2, door2, tub, sink3, table3, sink3, armchair, sink4, and bed placed in various room, which helps in generating more realistic results. We have 10,000 floor-plan dataset images containing 1,000 floor-plan layouts and around 300,000 furniture items of 16 classes. Figure 5 shows the sample image of this floor-plan dataset in which all furniture items are nearly the same size. Further, we notice that some furniture items are present in particular rooms, which helps recognize room categories for different furniture items.

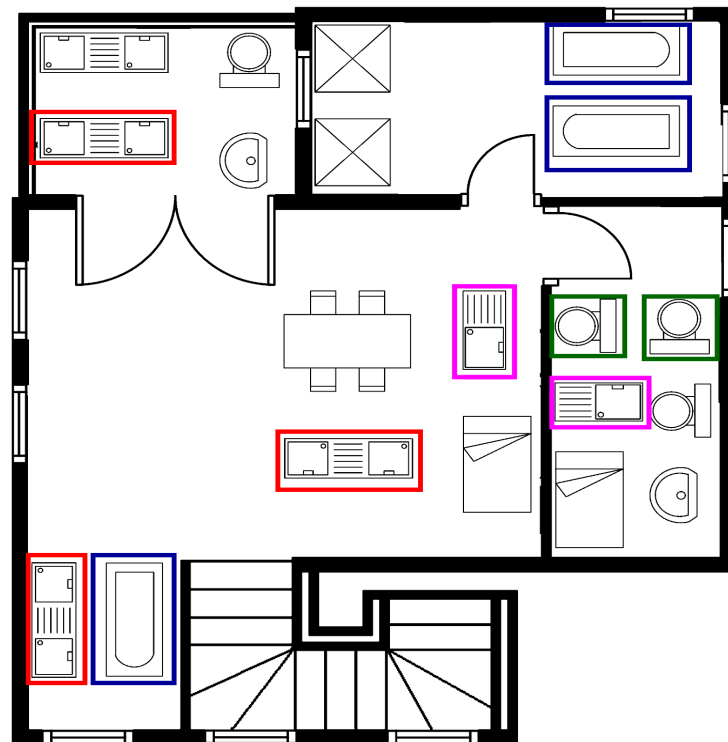


Figure 5. Sample image of the floor-plan dataset with different types of augmentation to create variation. The size of all furniture items is almost the same. The red and blue rectangular box objects show that sink and tub classes can vary in orientation.

We have different types of augmentation to create variation in our dataset. The first type of augmentation is rotation. We rotate with random angle between $[0, 30, 60, 90, 120, 180, 210, 270, 330]$. Figure 5 shows that tub and sink furniture class items have different directions based on the provided angle. Another augmentation method is scaling with a random scaling factor between $[20, 40, 75, 100, 145, 185, 200]$. During scaling, we keep the same aspect ratio for all furniture classes.

5. Experiments

5.1. Implementation Details

We used Mask R-CNN [22] based semi-supervised approach with ResNet-50[24] and ResNet-101 backbone pre-trained on ImageNet [25] as a detection method. The training data contains 1%, 5% and 10% floor-plan images as labelled data and the remaining images as unlabeled training data. We have 5 data folds for each type and calculated the final performance as the average of all folds. Our methodology and hyper-parameters are formulated from MMDetection [78]. For training, we used anchors with a three-aspect ratio and five-scale value and formed 1k and 2k region proposals with a 0.7 non-maximum suppression threshold. We selected a total of 512 proposals from 2k as box candidates for the training of RCNN. The IoU threshold value is set to 0.5 for mask bounding boxes.

5.1.1. Partially Labeled Data

We performed training for 80k iterations on 8 GPUs (A100) using eight images per GPU. For initial training, the learning rate has the value of 0.01 and then reduced to 0.001 at 30k iteration and 0.0001 at 40k iteration. The momentum and weight decay values are 0.9 and 0.0001, respectively. The data sampling ratio has an initial value of 0.2 and then decreases to 0 for the last 5k iterations, and the foreground threshold has a 0.9 value.

For selecting box regression pseudo-labels, we set a threshold value of 0.02, and the Njitter value is set as 10 to calculate the reliability of box localization. The jitter boxes are

sampled by setting offset values for all coordinates and selecting the offsets from -6% to 6% width or height of pseudo-box candidates. Moreover, different augmentations are used like FixMatch [8] to generate pseudo-label and train the labelled and unlabeled data.

5.1.2. Fully Labeled Data

We have 150k training iterations on 4 GPUs (A100) using eight images per GPU. For initial training, the learning rate has a value of 0.01 and then reduced to 0.001 at 30k iteration and 0.0001 at 40k iteration. The momentum has a value of 0.9. The data sampling ratio has an initial value of 0.2 and then decreases to 0 for the last 15k iterations, and the foreground threshold has a 0.9 value. The weight decay has a value of 0.0001. We assigned the N_{jitter} value of 10 to estimate box localization probability and the threshold value of 0.02 for selecting box regression pseudo labels.

6. Evaluation Criteria

We used some detection evaluation metrics to evaluate the performance of semi-supervised based floor-plan object detection approach. This section explains the used evaluation metrics.

6.1. Intersection over Union

We calculated the intersection over union (IoU) by taking the intersection divided by the union for the area of the ground-truth box A_g and the generated bounding box A_p .

$$IoU = \frac{area(A_g \cap A_p)}{area(A_g \cup A_p)} \quad (8)$$

IoU is used to estimate whether a detected object is false positive or true positive.

6.2. Average Precision

We calculated the average precision (AP) using a precision-recall curve. It is the area under the precision-recall curve and can be determined using this equation:

$$IoU = \sum_{k=1}^N (R_{k+1} - R_k) P_{intr}(R_{k+1}) \quad (9)$$

Where R_1, R_2, \dots, R_k are the values of the recall parameter.

6.3. mAP

The mean average precision (mAP) is the most common metric for evaluating the performance of object detection methods. We calculate it by taking the mean of average precision for all classes s of the dataset. While working with a floor-plan dataset, it is preferred to calculate mAP to lower 16 classes to a set of classes s . The overall performance of mAP depends on class mapping, where a slight change in the performance of one class can affect overall mAP; that is the only drawback of mAP. We set the IoU threshold value of 0.5 and 0.75 to calculate the mAP as:

$$mAP = \frac{1}{S} \sum_{s=1}^S AP_s \quad (10)$$

where S is the total number of classes. For our floor-plan dataset S , its value is 16.

7. Results and Discussion

We use Mask R-CNN [22] based semi-supervised network on the floor-plan dataset. This section will explain the qualitative as well as quantitative results of the student-teacher network. For our experiment, We take 1%, 5% and 10% floor-plan images as labelled data and the rest of the floor-plan images as unlabeled data. We have 5 data folds for each type and calculated the final performance as the average of all folds. We train and evaluate the

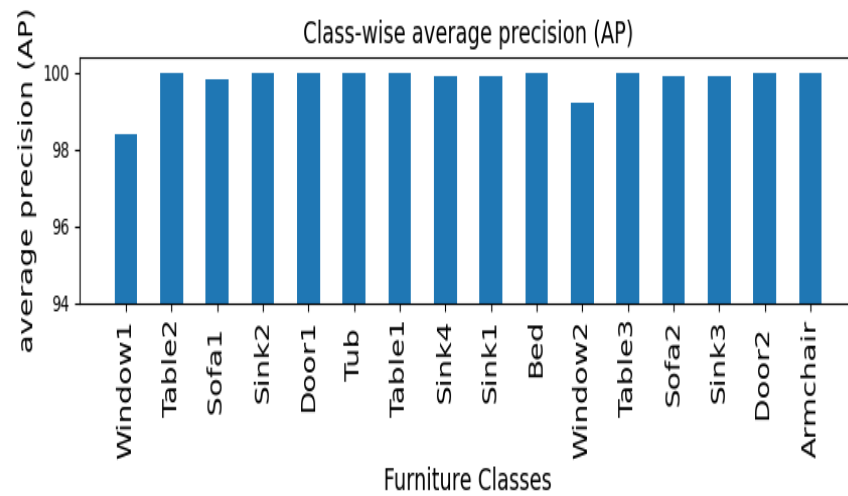


Figure 6. Class-wise average precision (AP) results of 5 data folds with 10% label data using Mask R-CNN [22] with ResNet-101 [24] backbone. Some furniture classes like armchairs, door2, table3, bed, table1, tub, door1, sink2 and table2 show one average precision while other classes like window1 and window2 need further improvements.

approach on Faster-RCNN [30] and Mask R-CNN [22]. Furthermore, we also compare the algorithm's performance on ResNet-50 [24] and ResNet-101 backbone with Mask R-CNN [22].

Table 1. Compared with different supervised train detectors on floor-plan under the semi-supervised setting. We used Faster R-CNN [30] and Mask R-CNN [22] with ResNet-50 [24] and ResNet-101 backbone for our experiment and then determined the relative error between these detectors represented with down arrows. It shows that Mask R-CNN [22] with 10% label data decreases the error by 37.5% and 50% with ResNet-50 and ResNet-101 backbone, respectively.

Detector	Backbone	1%	5%	10%
Faster R-CNN	ResNet-50	98.1	98.5	99.2
Faster R-CNN	ResNet-101	98.3	99.4	99.6
Mask R-CNN	ResNet-50	98.27 ↓ 8.94	98.74 ↓ 16	99.5 ↓ 37.5
Mask R-CNN	ResNet-101	98.8 ↓ 29.4	99.7 ↓ 50	99.8 ↓ 50

Figure 6 shows the average precision of every class separately. It is evident that some classes like armchairs, door2, table3, bed, table1, tub, door1, sink2 and table2 show one average precision while all other classes show average precision above 0.95 except window1 class. We can observe for which classes our model performs well and where we need further improvements. Figure 7 shows the furniture items detection and localization on the floor-plan test dataset. The final result where furniture items are detected and labelled in different colours accurately detects all 16 classes.

Table 2. The performance comparison by setting different values of jittered boxes using Mask R-CNN [22] with ResNet-101 [24] backbone on 10% label data.

N_{jitter}	mAP	mAP@0.5	mAP@0.75
5	0.996	0.998	0.997
10	0.998	1.0	1.0
15	0.997	1.0	1.0

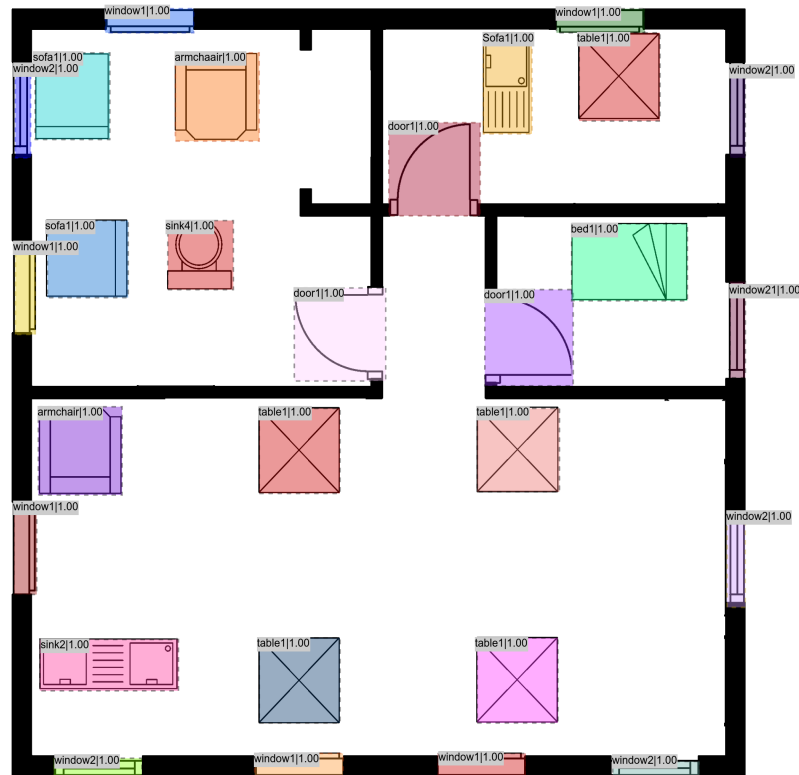


Figure 7. Test images where furniture items are detected and labelled using Mask R-CNN [22] with ResNet-101 [24] backbone on 10% label data.

We determine the relative error between Mask R-CNN [22] and Faster R-CNN [30] detectors using different backbone networks. Table 1 shows comparison of these detectors with ResNet-50 [24] and ResNet-101 backbone on floor-plan dataset under semi-supervised setting. It shows that Mask R-CNN decreases the error by 8.94%, 16%, 37.5% with ResNet-50 backbone and 29.4%, 50%, 50% with ResNet-101 backbone for 1%, 5% and 50% label data, respectively. Using 1% labelled data, we are getting 98.8% mAP on Mask R-CNN with ResNet-101 backbone, which shows that this approach gives the best results using a small amount of labelled data. This comparison also shows that the ResNet-101 backbone gives better results than the ResNet-50 [24] backbone for both detectors.

We also study the behaviour of hyper-parameters on model performance. The first hyperparameter is the jittered-box value that calculates the localization reliability of pseudo boxes. Table 2 compares the performance under different values of jittered boxes. By setting jittered box value of 10 it gives mAP 99.6% while $AP_{0.5}$ and $AP_{0.75}$ are 99.8% and 99.7%, respectively. We can observe from Table 2 that the model gives the highest accuracy when N_{jitter} has a value of 10.

We apply IoU between the teacher-created pseudo-boxes and student-created box-candidate to provide background and foreground labels as an ordinary object detection model does. In that case, some foreground boxes are incorrectly classified as negative, reducing performance. Table 3 shows the box regression-variance threshold. we are getting the best results by setting the threshold value to 0.02. However, a high threshold value provides good foreground precision, and the recall of box-candidate decreases quickly.

Figure 8 shows a test image where some furniture items are miss-classified. The network confuses between window1 and window2. The green box wrongly detects two windows as one window named window2. The size of window1 and window2 objects is small compared to all other floor-plan objects. The detection performance of such small objects can be improved further where the background occupies 95% area of the image.

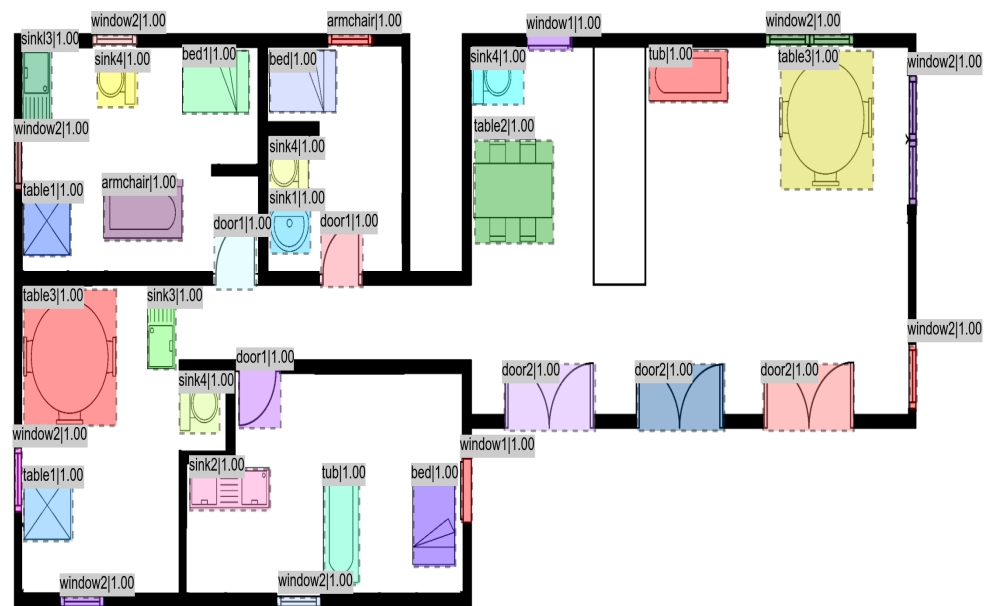


Figure 8. Test images where furniture items are detected and labelled using Mask R-CNN [22] with ResNet-101 [24] backbone on 10% label data. Some furniture items are miss-classified in the image, such as the green and purple window boxes wrongly detecting two windows as one named window2.

Table 3. The performance comparison by setting different box regression-variance thresholds to select pseudo-boxes for box regression using Mask R-CNN [22] with ResNet-101 [24] backbone on 10% label data. We set the threshold value of 0.02 for our experiment as it gives the best performance.

Threshold	mAP	mAP@0.5	mAP@0.75
0.04	0.998	0.998	1.0
0.03	0.997	1.0	1.0
0.02	0.998	1.0	1.0
0.01	0.992	1.0	1.0

7.1. Comparison with prior SOTA approaches

Table 4 shows the comparison of our semi-supervised network performance for 5 data folds with 10% label data on Faster R-CNN [30] and Mask R-CNN [22] with previously presented supervised approaches. We can not directly compare Ziran et al. [79] results because of the different datasets. It is observed from table 4 that the semi-supervised approach outperforms the previous supervised approaches using just 10% of label data.

8. Conclusion and Future Work

We examine the capabilities of the semi-supervised approach to detect objects in floor-plan data. It pulls information from the teacher network and feeds it to the student network. The teacher model uses unlabeled data to form pseudo-boxes, and the student model uses both unlabeled data (with the pseudo boxes) and labelled data as ground truth for training. On Mask R-CNN [22] detector with ResNet-101 backbone, the proposed approach achieves 98.8(%) mAP, 99.7(%) mAP, 99.8(%) mAP with 1%, 5% and 10% labelled data, respectively. We can see from the results that we can obtain the best performance by just using 10% labelled data. Furthermore, this experiment can be implemented in various floor-plan applications such as floor-plan text generation, furniture fitting, helping impaired people to analyze house design and for customers to buy a house online. Earlier, all these applications used supervised learning approaches [79,80] for floor-plan object

Table 4. Compared previously supervised detectors with our semi-supervised approach trained on the floor-plan dataset using Mask R-CNN [22] and Faster R-CNN [30] with ResNet-101 [24] backbone on 10% label data.

Method	Learning Approach	Dataset	Objects	Detector	mAP
Ziran et al. [79]	supervised	d1	1111	Faster R-CNN	0.31
Ziran et al. [79]	supervised	d2	1111	Faster R-CNN	0.39
Singh et al. [80]	supervised	SESYD+ROBIN	20,670+510	Faster R-CNN	0.756
Singh et al. [80]	supervised	SESYD+ROBIN	20,670+510	YOLO	0.857
Mishra et al. [27]	supervised	SFPI	316,160	Cascade Mask R-CNN	0.995
Our	semi-supervised(10%)	SFPI	316,160	Faster R-CNN	0.996
Our	semi-supervised(10%)	SFPI	316,160	Mask R-CNN	0.998

detection. However, now with our experiment, it is clear that the semi-supervised [21] approach gives better results for these applications.

In future, we can improve Mask R-CNN [22] based semi-supervised floor-plan detection system in different ways. We can add text information to detect room types, especially rooms that are not physically separated, like the dining hall attached to the kitchen. We can also label rooms according to their functionality. Further research using noisy labels in training and uncertainty estimation are also a few important topics to boost the efficiency of semi-supervised based object detection.

Author Contributions: writing—original draft preparation, T.S., K.A.H., M.Z.A.; writing—review and editing, T.S., K.A.H., M.Z.A.; supervision and project administration, A.P., D.S. All authors have read and agreed to the submitted version of the manuscript.

Funding: The work leading to this publication has been partially funded by the European project INFINITY under Grant Agreement ID 883293.

Informed Consent Statement: Not applicable.

Data Availability Statement: Not applicable.

Acknowledgments: Not applicable.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Sajjadi, M.; Javanmardi, M.; Tasdizen, T. Mutual exclusivity loss for semi-supervised deep learning. 2016 IEEE International Conference on Image Processing (ICIP), 2016, pp. 1908–1912. doi:10.1109/ICIP.2016.7532690.
2. Tarvainen, A.; Valpola, H. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in Neural Information Processing Systems*; Guyon, I.; Luxburg, U.V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; Garnett, R., Eds. Curran Associates, Inc., 2017, Vol. 30.
3. Miyato, T.; Maeda, S.I.; Koyama, M.; Ishii, S. Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2019**, *41*, 1979–1993. doi:10.1109/TPAMI.2018.2858821.
4. Berthelot, D.; Carlini, N.; Goodfellow, I.; Papernot, N.; Oliver, A.; Raffel, C.A. MixMatch: A Holistic Approach to Semi-Supervised Learning. *Advances in Neural Information Processing Systems*; Wallach, H.; Larochelle, H.; Beygelzimer, A.; dAlché-Buc, F.; Fox, E.; Garnett, R., Eds. Curran Associates, Inc., 2019, Vol. 32.
5. Xie, Q.; Dai, Z.; Hovy, E.H.; Luong, M.; Le, Q.V. Unsupervised Data Augmentation. *CoRR* **2019**, *abs/1904.12848*, [1904.12848].
6. Berthelot, D.; Carlini, N.; Cubuk, E.D.; Kurakin, A.; Sohn, K.; Zhang, H.; Raffel, C. ReMixMatch: Semi-Supervised Learning with Distribution Alignment and Augmentation Anchoring. *CoRR* **2019**, *abs/1911.09785*, [1911.09785].
7. Xie, Q.; Hovy, E.H.; Luong, M.; Le, Q.V. Self-training with Noisy Student improves ImageNet classification. *CoRR* **2019**, *abs/1911.04252*, [1911.04252].

8. Sohn, K.; Berthelot, D.; Li, C.; Zhang, Z.; Carlini, N.; Cubuk, E.D.; Kurakin, A.; Zhang, H.; Raffel, C. FixMatch: Simplifying Semi-Supervised Learning with Consistency and Confidence. *CoRR* **2020**, *abs/2001.07685*, [[2001.07685](#)].
9. Simard, P.; Steinkraus, D.; Platt, J. Best practices for convolutional neural networks applied to visual document analysis. Seventh International Conference on Document Analysis and Recognition, 2003. *Proceedings.*, 2003, pp. 958–963. doi:10.1109/ICDAR.2003.1227801.
10. Krizhevsky, A.; Sutskever, I.; Hinton, G.E. ImageNet Classification with Deep Convolutional Neural Networks. *Commun. ACM* **2017**, *60*, 84–90. doi:10.1145/3065386.
11. Zhu, J.; Park, T.; Isola, P.; Efros, A.A. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. *CoRR* **2017**, *abs/1703.10593*, [[1703.10593](#)].
12. Hoffman, J.; Tzeng, E.; Park, T.; Zhu, J.; Isola, P.; Saenko, K.; Efros, A.A.; Darrell, T. CyCADA: Cycle-Consistent Adversarial Domain Adaptation. *CoRR* **2017**, *abs/1711.03213*, [[1711.03213](#)].
13. Zakharov, S.; Kehl, W.; Ilic, S. DeceptionNet: Network-Driven Domain Randomization. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019, pp. 532–541. doi:10.1109/ICCV.2019.00062.
14. Cubuk, E.D.; Zoph, B.; Mané, D.; Vasudevan, V.; Le, Q.V. AutoAugment: Learning Augmentation Strategies From Data. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 113–123. doi:10.1109/CVPR.2019.00020.
15. Zoph, B.; Cubuk, E.D.; Ghiasi, G.; Lin, T.; Shlens, J.; Le, Q.V. Learning Data Augmentation Strategies for Object Detection. *CoRR* **2019**, *abs/1906.11172*, [[1906.11172](#)].
16. Gimenez, L.; Hippolyte, J.L.; Robert, S.; Suard, F.; Zreik, K. Review: reconstruction of 3D building information models from 2D scanned plans. *Journal of Building Engineering* **2015**, *2*, 24–35. doi:https://doi.org/10.1016/j.jobe.2015.04.002.
17. Ahmed, S.; Liwicki, M.; Weber, M.; Dengel, A. Automatic Room Detection and Room Labeling from Architectural Floor Plans. 2012 10th IAPR International Workshop on Document Analysis Systems, 2012, pp. 339–343. doi:10.1109/DAS.2012.22.
18. Sohn, K.; Zhang, Z.; Li, C.; Zhang, H.; Lee, C.; Pfister, T. A Simple Semi-Supervised Learning Framework for Object Detection. *CoRR* **2020**, *abs/2005.04757*, [[2005.04757](#)].
19. Zoph, B.; Ghiasi, G.; Lin, T.; Cui, Y.; Liu, H.; Cubuk, E.D.; Le, Q.V. Rethinking Pre-training and Self-training. *CoRR* **2020**, *abs/2006.06882*, [[2006.06882](#)].
20. Gurkan, H.; de Véricourt, F. Contracting, Pricing, and Data Collection Under the AI Flywheel Effect. *Microeconomics: Asymmetric & Private Information eJournal* **2020**.
21. Xu, M.; Zhang, Z.; Hu, H.; Wang, J.; Wang, L.; Wei, F.; Bai, X.; Liu, Z. End-to-End Semi-Supervised Object Detection with Soft Teacher. *CoRR* **2021**, *abs/2106.09018*, [[2106.09018](#)].
22. He, K.; Gkioxari, G.; Dollár, P.; Girshick, R. Mask R-CNN. 2017 IEEE International Conference on Computer Vision (ICCV), 2017, pp. 2980–2988. doi:10.1109/ICCV.2017.322.
23. Lin, T.Y.; Dollár, P.; Girshick, R.; He, K.; Hariharan, B.; Belongie, S. Feature Pyramid Networks for Object Detection. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 936–944. doi:10.1109/CVPR.2017.106.
24. Szegedy, C.; Ioffe, S.; Vanhoucke, V. Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning. *CoRR* **2016**, *abs/1602.07261*, [[1602.07261](#)].
25. Krizhevsky, A.; Sutskever, I.; Hinton, G.E. ImageNet classification with deep convolutional neural networks. *Communications of the ACM* **2012**, *60*, 84–90.
26. Delalandre, M.; Valveny, E.; Pridmore, T.; Karatzas, D. Generation of Synthetic Documents for Performance Evaluation of Symbol Recognition & Spotting Systems. *Int. J. Doc. Anal. Recognit.* **2010**, *13*, 187–207. doi:10.1007/s10032-010-0120-x.
27. Mishra, S.; Hashmi, K.A.; Pagani, A.; Liwicki, M.; Stricker, D.; Afzal, M.Z. Towards Robust Object Detection in Floor Plan Images: A Data Augmentation Approach. *Applied Sciences* **2021**, *11*. doi:10.3390/app112311174.
28. Girshick, R.; Donahue, J.; Darrell, T.; Malik, J. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. 2014 IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 580–587. doi:10.1109/CVPR.2014.81.
29. Girshick, R. Fast R-CNN. 2015 IEEE International Conference on Computer Vision (ICCV), 2015, pp. 1440–1448. doi:10.1109/ICCV.2015.169.
30. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2017**, *39*, 1137–1149. doi:10.1109/TPAMI.2016.2577031.
31. Lin, T.Y.; Dollár, P.; Girshick, R.; He, K.; Hariharan, B.; Belongie, S. Feature Pyramid Networks for Object Detection. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 936–944. doi:10.1109/CVPR.2017.106.
32. Cai, Z.; Vasconcelos, N. Cascade R-CNN: Delving Into High Quality Object Detection. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 6154–6162. doi:10.1109/CVPR.2018.00644.
33. Redmon, J.; Divvala, S.; Girshick, R.; Farhadi, A. You Only Look Once: Unified, Real-Time Object Detection. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 779–788. doi:10.1109/CVPR.2016.91.
34. Redmon, J.; Farhadi, A. YOLO9000: Better, Faster, Stronger. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 6517–6525. doi:10.1109/CVPR.2017.690.
35. Liu, W.; Anguelov, D.; Erhan, D.; Szegedy, C.; Reed, S.; Fu, C.Y.; Berg, A.C. SSD: Single Shot MultiBox Detector. *Computer Vision – ECCV 2016*; Leibe, B.; Matas, J.; Sebe, N.; Welling, M., Eds.; Springer International Publishing: Cham, 2016; pp. 21–37.
36. Lin, T.Y.; Goyal, P.; Girshick, R.; He, K.; Dollár, P. Focal Loss for Dense Object Detection. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.

37. Tian, Z.; Shen, C.; Chen, H.; He, T. FCOS: Fully Convolutional One-Stage Object Detection. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019, pp. 9626–9635. doi:10.1109/ICCV.2019.00972.
38. Deng, J.; Dong, W.; Socher, R.; Li, L.J.; Li, K.; Fei-Fei, L. ImageNet: A large-scale hierarchical image database. 2009 IEEE Conference on Computer Vision and Pattern Recognition, 2009, pp. 248–255. doi:10.1109/CVPR.2009.5206848.
39. Misra, I.; Shrivastava, A.; Hebert, M. Watch and Learn: Semi-Supervised Learning of Object Detectors from Videos. *CoRR* **2015**, *abs/1505.05769*, [1505.05769].
40. McLachlan, G.J. Iterative Reclassification Procedure for Constructing an Asymptotically Optimal Rule of Allocation in Discriminant Analysis. *Journal of the American Statistical Association* **1975**, *70*, 365–369, [https://www.tandfonline.com/doi/pdf/10.1080/01621459.1975.10479874]. doi:10.1080/01621459.1975.10479874.
41. Bhatt, J.; Hashmi, K.A.A.; Afzal, M.Z.; Stricker, D. A Survey of Graphical Page Object Detection with Deep Neural Networks. *Applied Sciences* **2021**, *11*. doi:10.3390/app11125344.
42. Ahmed, M.; Hashmi, K.A.; Pagani, A.; Liwicki, M.; Stricker, D.; Afzal, M.Z. Survey and Performance Analysis of Deep Learning Based Object Detection in Challenging Environments. *Sensors* **2021**, *21*. doi:10.3390/s21155116.
43. Nazir, D.; Hashmi, K.A.; Pagani, A.; Liwicki, M.; Stricker, D.; Afzal, M.Z. HybridTabNet: Towards Better Table Detection in Scanned Document Images. *Applied Sciences* **2021**, *11*. doi:10.3390/app11188396.
44. Hashmi, K.A.; Pagani, A.; Liwicki, M.; Stricker, D.; Afzal, M.Z. Cascade Network with Deformable Composite Backbone for Formula Detection in Scanned Document Images. *Applied Sciences* **2021**, *11*. doi:10.3390/app11167610.
45. Hashmi, K.A.; Pagani, A.; Liwicki, M.; Stricker, D.; Afzal, M.Z. CasTabDetectorRS: Cascade Network for Table Detection in Document Images with Recursive Feature Pyramid and Switchable Atrous Convolution. *Journal of Imaging* **2021**, *7*. doi:10.3390/jimaging7100214.
46. Naik, S.; Hashmi, K.A.; Pagani, A.; Liwicki, M.; Stricker, D.; Afzal, M.Z. Investigating Attention Mechanism for Page Object Detection in Document Images. *Applied Sciences* **2022**, *12*. doi:10.3390/app12157486.
47. Zhang, F.; Fan, X.; Ai, G.; Song, J.; Qin, Y.; Wu, J. Accurate Face Detection for High Performance. *CoRR* **2019**, *abs/1905.01585*, [1905.01585].
48. Khan, A.H.; Munir, M.; van Elst, L.; Dengel, A. F2DNet: Fast Focal Detection Network for Pedestrian Detection, 2022. doi:10.48550/ARXIV.2203.02331.
49. Bachman, P.; Alsharif, O.; Precup, D. Learning with Pseudo-Ensembles, 2014. doi:10.48550/ARXIV.1412.4864.
50. Miyato, T.; Maeda, S.I.; Koyama, M.; Ishii, S. Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2019**, *41*, 1979–1993. doi:10.1109/TPAMI.2018.2858821.
51. Sajjadi, M.; Javanmardi, M.; Tasdizen, T. Regularization with Stochastic Transformations and Perturbations for Deep Semi-Supervised Learning. Proceedings of the 30th International Conference on Neural Information Processing Systems; Curran Associates Inc.: Red Hook, NY, USA, 2016; NIPS'16, p. 1171–1179.
52. Laine, S.; Aila, T. Temporal Ensembling for Semi-Supervised Learning. *CoRR* **2016**, *abs/1610.02242*, [1610.02242].
53. Grandvalet, Y.; Bengio, Y. Semi-Supervised Learning by Entropy Minimization. Proceedings of the 17th International Conference on Neural Information Processing Systems; MIT Press: Cambridge, MA, USA, 2004; NIPS'04, p. 529–536.
54. Xie, Q.; Dai, Z.; Hovy, E.; Luong, T.; Le, Q. Unsupervised Data Augmentation for Consistency Training. Advances in Neural Information Processing Systems; Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M.; Lin, H., Eds. Curran Associates, Inc., 2020, Vol. 33, pp. 6256–6268.
55. Berthelot, D.; Carlini, N.; Goodfellow, I.J.; Papernot, N.; Oliver, A.; Raffel, C. MixMatch: A Holistic Approach to Semi-Supervised Learning. *CoRR* **2019**, *abs/1905.02249*, [1905.02249].
56. Radosavovic, I.; Dollár, P.; Girshick, R.; Gkioxari, G.; He, K. Data Distillation: Towards Omni-Supervised Learning. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 4119–4128. doi:10.1109/CVPR.2018.00433.
57. Li, Y.; Huang, D.; Qin, D.; Wang, L.; Gong, B. Improving Object Detection with Selective Self-Supervised Self-Training. Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX; Springer-Verlag: Berlin, Heidelberg, 2020; p. 589–607. doi:10.1007/978-3-030-58526-6_35.
58. Wang, K.; Yan, X.; Zhang, D.; Zhang, L.; Lin, L. Towards Human-Machine Cooperation: Self-Supervised Sample Mining for Object Detection. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 1605–1613. doi:10.1109/CVPR.2018.00173.
59. Jeong, J.; Lee, S.; Kim, J.; Kwak, N. Consistency-based Semi-supervised Learning for Object detection. Advances in Neural Information Processing Systems; Wallach, H.; Larochelle, H.; Beygelzimer, A.; d'Elché-Buc, F.; Fox, E.; Garnett, R., Eds. Curran Associates, Inc., 2019, Vol. 32.
60. Tang, P.; Ramaiah, C.; Xu, R.; Xiong, C. Proposal Learning for Semi-Supervised Object Detection. *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)* **2021**, pp. 2290–2300.
61. Geifman, Y.; El-Yaniv, R. SelectiveNet: A Deep Neural Network with an Integrated Reject Option. *CoRR* **2019**, *abs/1901.09192*, [1901.09192].
62. Cubuk, E.D.; Zoph, B.; Shlens, J.; Le, Q.V. RandAugment: Practical data augmentation with no separate search. *CoRR* **2019**, *abs/1909.13719*, [1909.13719].

63. Ghorbel, A.; Lemaitre, A.; Anquetil, E.; Fleury, S.; Jamet, E. Interactive interpretation of structured documents: Application to the recognition of handwritten architectural plans. *Pattern Recognition* **2015**, *48*, 2446–2458. doi:<https://doi.org/10.1016/j.patcog.2015.01.028>.
64. Macé, S.; Locteau, H.; Valveny, E.; Tabbone, S. A system to detect rooms in architectural floor plan images. DAS '10, 2010.
65. Lladós, J.; López-Krahe, J.; Martí, E. A System to Understand Hand-Drawn Floor Plans Using Subgraph Isomorphism and Hough Transform. *Mach. Vision Appl.* **1997**, *10*, 150–158. doi:10.1007/s001380050068.
66. Eklund, A. Cascade Mask R-CNN and Keypoint Detection used in Floorplan Parsing, 2020.
67. Goyal, S.; Chattopadhyay, C.; Bhatnagar, G. Plan2Text: A framework for describing building floor plan images from first person perspective. 2018 IEEE 14th International Colloquium on Signal Processing Its Applications (CSPA), 2018, pp. 35–40. doi:10.1109/CSPA.2018.8368681.
68. Zeng, Z.; Li, X.; Yu, Y.K.; Fu, C.W. Deep Floor Plan Recognition Using a Multi-Task Network With Room-Boundary-Guided Attention. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019, pp. 9095–9103. doi:10.1109/ICCV.2019.00919.
69. Simonyan, K.; Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image Recognition. *CoRR* **2015**, *abs/1409.1556*.
70. Liu, C.; Wu, J.; Kohli, P.; Furukawa, Y. Raster-to-Vector: Revisiting Floorplan Transformation. 2017 IEEE International Conference on Computer Vision (ICCV), 2017, pp. 2214–2222. doi:10.1109/ICCV.2017.241.
71. Yamasaki, T.; Zhang, J.; Takada, Y. Apartment Structure Estimation Using Fully Convolutional Networks and Graph Model. Proceedings of the 2018 ACM Workshop on Multimedia for Real Estate Tech; Association for Computing Machinery: New York, NY, USA, 2018; RETech'18, p. 1–6. doi:10.1145/3210499.3210528.
72. Dodge, S.; Xu, J.; Stenger, B. Parsing floor plan images. 2017 Fifteenth IAPR International Conference on Machine Vision Applications (MVA), 2017, pp. 358–361. doi:10.23919/MVA.2017.7986875.
73. de las Heras, L.P.; Ahmed, S.; Liwicki, M.; Valveny, E.; Sánchez, G. Statistical segmentation and structural recognition for floor plan interpretation. *International Journal on Document Analysis and Recognition (IJDAR)* **2014**, *17*, 221–237.
74. Xie, S.; Girshick, R.B.; Dollár, P.; Tu, Z.; He, K. Aggregated Residual Transformations for Deep Neural Networks. *CoRR* **2016**, *abs/1611.05431*, [[1611.05431](#)].
75. Lin, T.; Dollár, P.; Girshick, R.B.; He, K.; Hariharan, B.; Belongie, S.J. Feature Pyramid Networks for Object Detection. *CoRR* **2016**, *abs/1612.03144*, [[1612.03144](#)].
76. Tarvainen, A.; Valpola, H. Weight-averaged consistency targets improve semi-supervised deep learning results. *CoRR* **2017**, *abs/1703.01780*, [[1703.01780](#)].
77. Shrivastava, A.; Gupta, A.; Girshick, R. Training Region-Based Object Detectors with Online Hard Example Mining. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 761–769. doi:10.1109/CVPR.2016.89.
78. Chen, K.; Wang, J.; Pang, J.; Cao, Y.; Xiong, Y.; Li, X.; Sun, S.; Feng, W.; Liu, Z.; Xu, J.; Zhang, Z.; Cheng, D.; Zhu, C.; Cheng, T.; Zhao, Q.; Li, B.; Lu, X.; Zhu, R.; Wu, Y.; Dai, J.; Wang, J.; Shi, J.; Ouyang, W.; Loy, C.C.; Lin, D. MMDetection: Open MMLab Detection Toolbox and Benchmark. *CoRR* **2019**, *abs/1906.07155*, [[1906.07155](#)].
79. Ziran, Z.; Marinai, S. Object Detection in Floor Plan Images. *Artificial Neural Networks in Pattern Recognition*; Pancioni, L.; Schwenker, F.; Trentin, E., Eds.; Springer International Publishing: Cham, 2018; pp. 383–394.
80. GitHub, I. Open Source Survey. <https://github.com/dwmsingh/Object-Detection-in-Floor-Plan-Images>, 2019.