

Table Detection for Visually Rich Document Images

Bin Xiao^a, Murat Simsek^a, Burak Kantarci^a, Ala Abu Alkheir^b

^a*School of Electrical Engineering and Computer Science, University of Ottawa, Ottawa, K1N 6N5, ON, Canada*

^b*Lytica Inc., 555 Legget Drive, Ottawa, K2K 2X3, ON, Canada*

Abstract

Table Detection (TD) is a fundamental task to enable visually rich document understanding, which requires the model to extract information without information loss. However, popular Intersection over Union (IoU) based evaluation metrics and IoU-based loss functions for the detection models cannot directly represent the degree of information loss for the prediction results. Therefore, we propose to decouple IoU into a ground truth coverage term and a prediction coverage term, in which the former can be used to measure the information loss of the prediction results. Besides, considering the sparse distribution of tables in document images, we use SparseR-CNN as the base model and further improve the model by using Gaussian Noise Augmented Image Size region proposals and many-to-one label assignments. Results under comprehensive experiments show that the proposed method can consistently outperform state-of-the-art methods with different IoU-based metrics under various datasets and demonstrate that the proposed decoupled IoU loss can enable the model to alleviate information loss.

Keywords: Object Detection, Table Detection, Tabular Data Extraction, Document Object Detection

1. Introduction

Table Detection (TD) is often a pre-processor step for information extraction and document understanding tasks [1, 2]. One typical formulation for the TD problem is transforming the electronic documents into images and then using object detection models to generate the bounding boxes of defined table objects. Current state-of-the-art methods [3, 4] for the TD problem usually employ two-stage object detectors, which require dense candidates and apply data augmentation and multiple-stage transfer learning techniques. However, tables in visually rich documents are usually well formatted and large so that human readers can easily interpret them. Besides, the number of tables in a single document image is typically small, which means their distribution in a single document is sparse. Based on these observations, we use SparseR-CNN [5] as the base model in this study, which is a competitive detector using sparse learnable regional proposals. It is worth mentioning that many state-of-the-art studies for the TD problem often

*This work was supported in part by Mathematics of Information Technology and Complex Systems (MITACS) Accelerate Program, Smart Computing for Innovation (SOSCIIP) Program and Lytica Inc.

Email addresses: bxiao103@uottawa.ca (Bin Xiao), murat.simsek@uottawa.ca (Murat Simsek), burak.kantarci@uottawa.ca (Burak Kantarci), ala_abualkheir@lytica.com (Ala Abu Alkheir)

employ two-stage detectors using dense candidates and multiple-stage transfer learning techniques, which are usually more complex than the proposed method. We also propose to use image size regional proposals to cover all the information of target tables in the proposal boxes and use the noise augmentation method to enrich the diversity of proposal boxes.

Since information extraction tasks often follow TD tasks, it is vital to avoid information loss. Therefore, a larger prediction box is preferable to a smaller box that can lose information even when the latter box has a larger Intersection over Union (IoU) score with the ground truth. Figure 1 uses a table as an example to further illustrate this observation. The green box in Figure 1 has an IoU score of 0.77 with the red ground truth box, while the blue box has an IoU score of 0.82. Even though the blue prediction has a larger IoU score, the green prediction is preferable for TD tasks. Motivated by these observations, we argue that the IoU score cannot directly reflect the information loss of a prediction box. Therefore, we propose to decouple the IoU score into two terms: a ground truth coverage term and a prediction coverage term, in which the former term can be used to measure the information loss for the prediction boxes. It is worth mentioning that the proposed decoupled IoU score termed the Information Coverage Score (ICS), can replace the IoU score in the IoU-based loss functions and evaluation metrics.

Division of Reproductive and Urologic Drug Products

ADMINISTRATIVE REVIEW OF APPLICATION

Application Number:	20-527/S-024
Name of Drug:	Prempro™/Premphase® (0.3 mg/1.5 mg)
Sponsor:	Wyeth-Ayerst Laboratories, Inc.
Material Reviewed:	Supplement-024
Submission Date:	November 5, 2001
Receipt Date:	November 7, 2001
Filing Date:	January 7, 2002
User-Fee Goal Date:	September 7, 2002
Proposed Indication:	Relief of moderate-to-severe vasomotor symptoms and treatment of vulvar and vaginal atrophy
Other Background Information:	NDA 20-303; NDA 20-527, NDA 4-782.

Review

PART I: OVERALL FORMATTING^a

Figure 1: Preferable prediction for TD tasks. The IoU scores of green and blue predictions are 0.77 and 0.82, respectively. The green prediction is preferable for TD tasks, even though its IoU score is smaller.

On the other hand, label assignment in object detection models is to define the classification and regression targets

for anchors [6]. Many studies [7, 6] have shown that label assignment plays a vital role in the success of a detector, and the one-to-one scheme used in SparseR-CNN [5] is not optimal. Besides, SparseR-CNN employs a cascade architecture that uses the outputs of i th Dynamic Head as the inputs of the $i + 1$ th Dynamic Head to refine the predictions, as shown in Figure 2, which means that the proposal quality of each Dynamic Head varies. Therefore, inspired by the studies [7, 8, 6, 9, 10], we leverage a SimOTA [10] based many-to-one label assignment approach, which further improves the SimOTA by adapting a dynamic scheduling scheme to adjust the number of positive assignments dynamically and integrating the proposed ICS loss to the cost function.

To sum up, the contributions of this study are three-fold:

1. We introduce a decoupled IoU score, termed Information Coverage Score (ICS), which can reflect the information loss of the prediction boxes when it is used as evaluation metrics and encourage the model to alleviate the information loss when it is used as loss functions.
2. We improve the SimOTA method by adapting a dynamic scheduling scheme and integrating the ICS loss, propose a Gaussian Noise Augmented Image Size region proposal method, and apply them to the SparseR-CNN model to further improve the performance of the proposed method.
3. To compare with state-of-the-art models fairly, we first conduct extensive experiments using IoU-based evaluation metrics and loss functions on various manually annotated datasets to demonstrate the efficiency and effectiveness of our proposed detection model. Then, we conduct further experiments to demonstrate the benefits of the proposed ICS score when it is applied to the TD problem. The experimental results show that the proposed method can consistently outperform the state-of-the-art benchmark models under different evaluation metrics.

The rest of this paper is organized as follows: Section 2 discusses related studies, including studies in Object Detection models and Table Detection models. Section 3 introduces the formal problem definition and our proposed SparseTableDet model. Section 4 presents the experiments and discusses the design aspects of the proposed method. At last, we draw our conclusion and possible directions in section 5.

2. Related Work

As discussed in section 1, we formulate the TD problem as an Object Detection problem and propose an ICS that can replace the IoU in loss functions and evaluation metrics. Therefore, this section first discusses popular object detection models and loss functions. Since the proposed method employs a Many-to-One label assignment, we include popular label assignment studies. At last, we focus on the studies specifically designed or optimized for the TD problem.

2.1. Object Detection Methods

Object detection problem has been widely discussed in recent years. Object detection models are often categorized into one-stage and two-stage models based on their number of regression steps. Popular one-stage models include

YOLO [11] and its variants [12, 13, 14, 10], SSD [15], FCOS [16], and many others. These one-stage models do not contain the step of generating region proposals and usually have faster training and inference speed compared with two-stage models. In contrast, two-stage models usually first use a region proposal network to generate a series of regional proposals, then further regress and classify the regional proposals by well-designed models. Typical two-stage models include FasterR-CNN [17], MaskR-CNN [18], CascadeR-CNN [19] and many others. Along with these typical one-stage and two-stage models, transformer-based approaches recently attracted a lot of attention. DETR [20] is the first study introducing transformer and self-attention mechanism [21] to the object detection. Following the design of DETR, many variants of DETR have been proposed to improve the performance further and accelerate the convergence, such as Deformable-DETR [22] and DAB-DETR [23]. SparseR-CNN [5] adopts the Set prediction loss and Hungarian matching from DETR and proposes to use learnable region proposals and learnable instance features to simplify the region proposal generation. One key characteristic of the learnable region proposals in SparseR-CNN is that they can be initialized with different methods, including random initialization and image size initialization. Our proposed method in this study is based on SparseR-CNN and uses the image size initialization considering the requirements of TD applications as discussed in section 1.

Loss functions used in object detectors’ regression and classification tasks have also been discussed widely. For the regression task, it is a natural choice to use L_n -norms and their variants, such as smooth-l1 [24], as the loss function. However, these functions are not aligned with the widely accepted evaluation metric IoU score, meaning that for some cases minimizing these loss functions cannot lead to better IoU scores [25, 26]. IoU-based loss functions can alleviate this issue and become the most popular choice. IoU loss [27] has the gradient vanishing issue when the prediction and ground truth boxes have no overlaps. GIoU loss [25] extends the IoU loss by adding an extra penalty term to alleviate the gradient vanishing problem when two boxes have no overlap. More specifically, assuming that A, B denote two arbitrary convex shapes and C is the smallest enclosing convex, then the term used in GIoU loss is defined as $\frac{|C-A \cup B|}{|C|}$, where $-$ means the complementary operation. A limitation of GIoU loss is that it can be degraded to IoU loss for enclosing bounding boxes. To address this limitation of GIoU loss, DIoU loss [28] proposes to use the distance between two boxes’ centers as the additional term, which can lead to faster convergence and alleviate the gradient vanishing problem. CIoU loss [29] also considers the geometric factors of bounding boxes and proposes an aspect ratio term, a distance term, and the IoU term. Besides these popular IoU-based loss functions, there are other loss functions without using IoU score, such as SCA Loss [26] and KLLoss [30]. SCA Loss defines two terms considering side overlap and corner distance. KLLoss requires the output to be a distribution instead of location coordinates.

Besides these object detection models, label assignment methods have also been widely discussed in recent studies. Typically, label assignment methods can be categorized into Fixed Label Assignment and Dynamic Label Assignment [6]. Fixed Label Assignments are methods defining fixed criterion to determine the positive and negative samples of each ground truth. For example, Region Proposal Network (RPN) in FasterR-CNN [17] defines two IoU scores as the thresholds to the positive and negative proposals. YOLO [11] uses the closest anchor points to the center of ground truths as positive anchor points. In contrast, Dynamic Label Assignment methods often formulate the

problem as an optimization problem and solve the problem more dynamically. For example, OTA [6] formulates the label assignment problem as an Optimal Transport (OT) problem, which can be optimized by the Sinkhorn-Knopp algorithm. SimOTA [9, 10] uses the top-k candidates whose centers are in the ground truth bounding boxes to avoid the time-consuming optimization process.

2.2. Table Detection

Many studies have discussed the TD problem recently. One of the most popular formulations for the TD problem is defining tables in visually rich documents as objects and then applying popular object detectors. Following this problem formulation, the object detection approaches discussed in section 2.1 can be easily adapted to the TD problem and widely used as benchmark models in many studies [31]. Considering the special characters and requirements for the TD problem, many studies further optimized the popular object detection methods to improve the model performance. Due to the limited number of training samples for the TD problem, transfer learning methods are widely used. CascadeTabNet [4] extends the Cascade Mask R-CNN [19] model and uses HRNet [32] as the backbone network. Besides, CascadeTabNet applies a two-stage transfer learning approach and various augmentation methods. Similarly, TableDet [3] is based on Cascade R-CNN [19], proposes a Table Aware Cutout augmentation method, and also leverages a two-step transfer learning approach to improve the model performance further. Besides two-stage detectors, one-stage methods, such as YOLO [11] and its variants, also have been adapted to the TD problem. YOLOv3-TD [33] employs the YOLOv3 [13] as the base model and proposes some adaptive adjustments, including a new anchor optimization method and a new post-processing process. Besides these one-stage and two-stage methods, transformer-based approaches such as DETR [20] and Deformable-Detr [22], also have been applied to the TD problem [34, 31]. There are many other studies discussing the TD problem, including DeepDeSRT [35], TableDet [3], and many others [36, 37]. All in all, these studies usually adopt popular object detection models to the TD problem and use some specifically designed methods to improve the model performance based on the characteristics of the TD problem.

3. Proposed Method

3.1. Architecture of the Proposed Method

Following the architecture of SparseR-CNN [5], our proposed method also consists of an Initialization Module, a Feature Pyramid Network, and a series of Dynamic Heads. The Initialization Module is used to initialize the learnable proposal boxes and the learnable proposal features. In this study, we use the Noise Augmented Image Size region proposals, which will be discussed in section 3.2. Feature Pyramid Network (FPN) [38] is the backbone network to generate image features for every Dynamic Head. The Dynamic Heads are used to do the regression and classification tasks. Dynamic Head $t + 1$ takes the image features generated by FPN and the outputs of Dynamic Head t , including the Refined Proposal Features and Refined Proposal Boxes, as the input to further refine the predictions of Dynamic

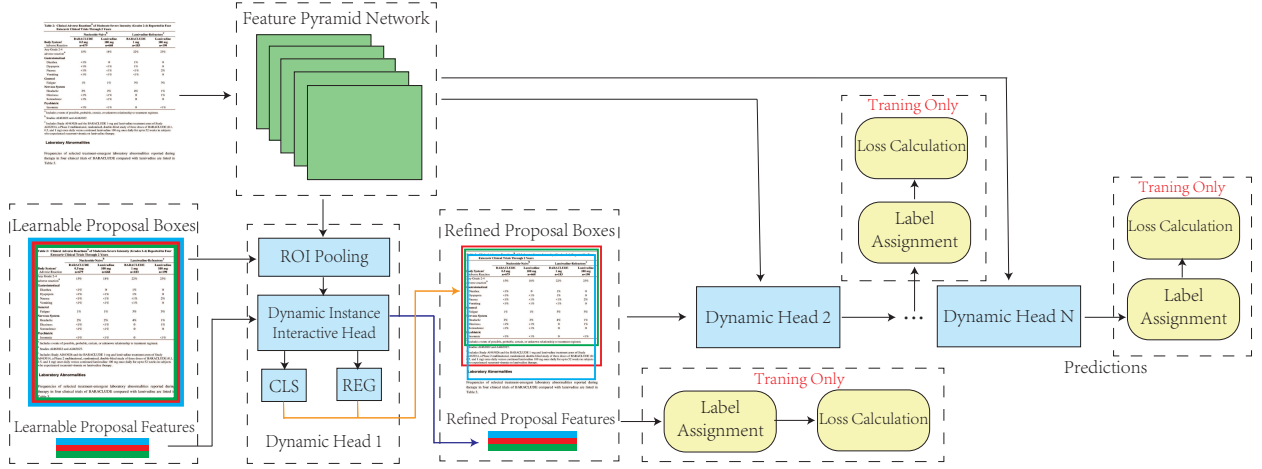


Figure 2: The overall architecture of the proposed method. Notably, all Dynamic Heads share an identity structure. We only show the details of Dynamic Head1 in this figure for simplicity.

Head t when $t > 1$. Since the predictions of each Dynamic Head are used to calculate the loss, a label assignment process is operated on these predictions to further calculate the losses. Since our refinements to the SparseR-CNN model are mainly on the proposal initialization, label assignment, and the loss functions, which will be discussed in section 3.2, 3.3 and 3.4, we keep the default implementations of SparseR-CNN for other parts which are detailed described in the study [5].

3.2. Noise Augmentation to Region Proposals

As discussed in section 1, TD applications typically require predictions to avoid information loss, and the tables in the documents are usually large and have no overlaps. Considering these characteristics of TD applications, using Image Size to initialize the region proposals becomes a good choice compared with other initialization methods, such as Random Initialization [5] and Grid Initialization [5], because it can avoid information loss at the first step of the detector. However, simply using a number of the same proposals may not be optimal. Therefore, we propose a simple but effective augmentation method to the region proposals by adding Gaussian Noise to enrich the proposals' diversity. More specifically, assuming that a proposal box is represented by its box center, width, and height, namely $b = \{c_x, c_y, w, h\}$, then the augmented proposal box can be defined as Equation 1, where \mathcal{N} means Gaussian Distribution. In our implementation, boxes are normalized, meaning that an image size box can be represented as $b = \{0.5, 0.5, 1, 1\}$, and μ, σ^2 are set as 0 and 0.01, respectively.

$$b_{aug} = f(\{c_x, c_y, w, h\}) = \{c_x + \epsilon_x, c_y + \epsilon_y, w - 2 \cdot |\epsilon_x|, h - 2 \cdot |\epsilon_y|\}, \epsilon_x \in \mathcal{N}(\mu, \sigma^2), \epsilon_y \in \mathcal{N}(\mu, \sigma^2) \quad (1)$$

It is worth mentioning that adding noise to the regional proposal boxes can be interpreted as a movement of these boxes. Since we set initial boxes to Image Size, any movement ϵ of the center leads to 2ϵ reduction of height or width,

as shown in Figure 3.

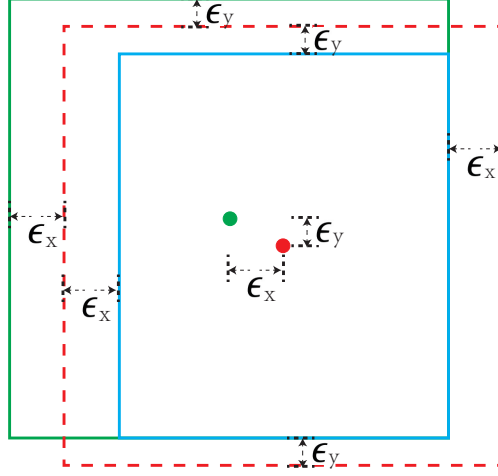


Figure 3: A sample of noise augmentation to a regional proposal box. The green box is the original box, the dashed red box is the result of center movement, and the blue box is the result box after augmentation.

3.3. Many-to-One Label Assignment

As discussed in section 1 and 2, label assignment plays a key role in the object detection models. SparseR-CNN employs Hungarian algorithm [20] to perform one-to-one label assignment so that the Non Maximum Suppression (NMS) can be removed from the processing pipeline. However, as aforementioned, tables in documents are usually large and have no overlaps, meaning that applying NMS to the TD problem doesn't necessarily lead to some drawbacks caused by the NMS, such as the performance degradation caused by the object overlaps [39]. Moreover, the cascaded Dynamic Heads take input proposal features and boxes with different qualities, making it necessary to determine the label assignment dynamically. Many studies [6, 7, 9] have demonstrated that Many-to-One label assignment can bring benefits to the model performance. Therefore, we adapt SimOTA [9] as the base label assignment method in this study.

SimOTA is a simplified version of OTA [6], which can avoid the complex optimization process of OTA. More specifically, SimOTA directly uses the top-k candidates whose centers are in the ground truth bounding boxes as the positive samples, as defined by Equation 2, which contains a classification cost, a regression cost, and a center cost. It is worth mentioning that the cost function of SparseR-CNN is the sum of the `cls_cost` and `regression_cost` in Equation 2.

$$cost_{SimOTA} = \underbrace{\lambda_{cls} \cdot cost_{cls}}_{cls_cost} + \underbrace{\lambda_{l1} \cdot cost_{l1} + \lambda_{giou} \cdot cost_{giou}}_{regression_cost} + \underbrace{\lambda_{center} \cdot cost_{center}}_{center_cost} \quad (2)$$

SimOTA employs a dynamic method to determine the number of positive samples assigned to each ground truth box using the sum of the top 10 IoU scores between a ground truth box and its corresponding prediction boxes without considering the difference of Dynamic Heads. Considering that the inputs of Dynamic Head $t+1$ should contain higher

quality boxes than that of t after the refinement of Dynamic Head t , the Dynamic Head $t + 1$ should have more positive samples. Therefore, we further extend this dynamic method of SimOTA by adding a scheduling scheme as defined by Equation 3, where N is the number of Dynamic Heads, IoU_i is the IoU matrix between the predictions and the i th ground truth, n is a hyper parameter and k_t^i means the number of positive samples assigned to the i th ground truth for Dynamic Head t .

$$k_t^i = SUM(TOPK(IoU_i, n - 0.5 * (N - t))), t \in [1, N] \quad (3)$$

At last, we can define the loss function as the sum of all the Dynamic Heads' loss, as defined by Equation 4. For the classification loss, we simply use cross entropy and focal loss [40] for binary and multi-class classification, respectively.

$$\mathcal{L} = \sum_{t=1}^N loss_t = \sum_{t=1}^N \lambda_{cls} loss_{cls}^t + \lambda_{l1} loss_{l1}^t + \lambda_{giou} loss_{giou}^t \quad (4)$$

It is worth mentioning that some studies, such as YOLOX [10] and Dynamic SparseR-CNN [7], use similar Label Assignment approaches. Dynamic SparseR-CNN introduces an assignment scheduling scheme to the OTA method [6] to dynamically adjust the number of positive label assignments. However, the OTA method requires a complex optimization procedure, which is significantly more time-consuming than SimOTA. Therefore, in this study, we leverage SimOTA and further improve it by adapting the assignment scheduling scheme and the proposed ICS loss function.

3.4. Information Coverage Score

As discussed in section 1, the IoU score cannot directly reflect the information coverage of the prediction boxes. This section discusses our proposed decoupled IoU, the Information Coverage Score (ICS). Assume that G and P are the ground truth and prediction boxes, respectively. IoU is the ratio of the intersection of G and P to the union of G and P , as defined by Equation 5. In contrast, ICS contains a ground truth coverage term (GT_Coverage) and a prediction coverage term (Pred_Coverage), as defined by Equation 6. The ground truth coverage term is the ratio of the intersection of G and P to the G , which can directly measure the information covered by the prediction box. Similarly, the prediction coverage score is defined as the ratio of the intersection of G and P to the P . Figure 4 shows three cases for the calculation of ICS, in which green boxes, red boxes, and yellow areas represent the ground truth boxes, prediction boxes, and their intersection areas. It is worth mentioning that the proposed ICS can be used to replace IoU in a variety of IoU-based loss functions, such as GIoU loss [25] and DIoU loss [28]. A simple ICS loss can be defined as Equation 7.

$$IoU = \frac{|G \cap P|}{|G \cup P|} \quad (5)$$

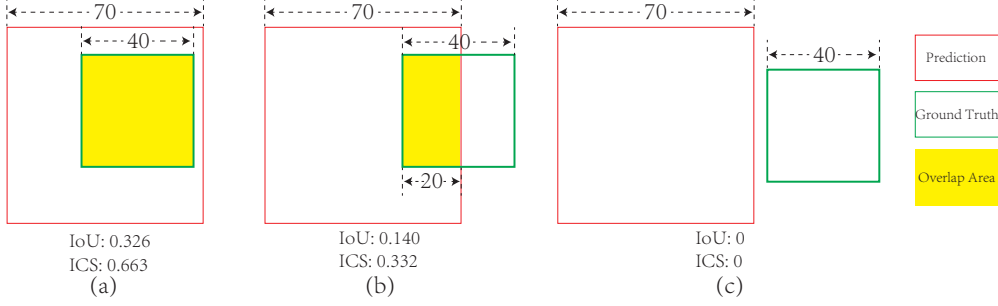


Figure 4: Three cases for the Information Coverage Score. λ is set to 0.5. All the boxes are squares.

$$ICS = \underbrace{\lambda \frac{|G \cap P|}{|G|}}_{GT_Coverage} + \underbrace{(1 - \lambda) \frac{|G \cap P|}{|P|}}_{Pred_Coverage} \quad (6)$$

$$ICS_Loss = 1 - ICS \quad (7)$$

4. Experimental Results and Analysis

This section compares the proposed method with state-of-the-art models using IoU-based evaluation metrics and loss functions. Then, we conduct experiments to demonstrate the benefits of using ICS as the evaluation metrics and loss functions. Lastly, an ablation study is conducted to demonstrate the effectiveness of the proposed many-to-one label assignment method and Gaussian Noise Augmented Image Size region proposal.

4.1. Experiment settings and Main results

Many datasets have been proposed for the TD problem. We can roughly categorize these datasets into two groups: human-annotated datasets and generated datasets by parsing meta-data. The former dataset type usually has higher-quality annotations, but the number of samples is usually limited. In contrast, the latter type can have a large number of samples but often contain much noise. In this study, we only consider the datasets with high-quality annotations, including ICDAR2017 [41], ICDAR2019 [42], TNCR [31] and ICT-TD [43] datasets.

The TD problem is often the pre-processor step of the information extraction tasks, which requires models to avoid missing tables or predicting other document components as tables. Therefore, the widely used evaluation mAP [44] for the detection models cannot fulfill this requirement per se. Hence, we use Precision, Recall, and F1 scores as evaluation metrics, which are also widely used in other studies [45, 41, 42, 31, 43]. However, the IoU thresholds for these metrics vary among studies, making it hard to compare these models directly. Therefore, in this study, we align the evaluation metric for all datasets and use Weighted Average F1, as defined in Equation 8, as evaluation metric whose thresholds are 60%, 70%, 80% and 90%, and attach detailed results containing Precision, Recall and F1 under the IoU thresholds from 50% to 95% with a 5% interval in Appendix A. A. Notably, the evaluation metric used here

is identical to the one used in ICDAR2019 competition [42], and other evaluation metrics, such as mAP. The metrics in other competitions [45, 41] can be found in the detailed experimental results in section Appendix A.

$$\text{Weighted Avg. F1} = \frac{\sum_{i=1}^4 IoU_i \cdot F1@IoU_i}{\sum_{i=1}^4 IoU_i} \quad (8)$$

Table 1: Key training parameters of the proposed model.

Parameter	Value	Description
IMS_PER_BATCH	16	number of training samples in an iteration
MAX_ITER	40,600	total number of mini-batch
STEPS	29,000	the mini-batch to apply the learning rate schedule
SCHEDULER	MultiStepLR	the scheduler to change the learning rate
BASE_LR	2.5e-05	the learning rate before applying the scheduler
WEIGHTS	r50_300pro_3x_model	initialization weight of the model
NUM_HEADS	6	the number of Dynamic Head
NUM_PROPOSALS	300	number of region proposals
OPTIMIZER	AdamW	the optimizer to train the model
LABEL_N	8	the hyper parameter N defined by Equation 3
NOISE_MEAN	0	mean value of the Gaussian Noise
NOISE_VAR	0.01	variance value of the Gaussian Noise
NMS_THRESH	0.9	non-maximum suppression threshold

For the benchmark models, we include all the types of popular object detection models discussed in section 2, including FasterR-CNN [17], MaskR-CNN [18], TableDet [3], DiffusionDet [8], Deformable-DETR [22], SparseR-CNN [5], RetinaNet [40], FCOS [16], YOLOX-X [10], YOLOR-X [46], YOLOv5-X [47], YOLOv7-X [48], YOLOv8-X [49]. We used the default settings of their implementations, trained FasterR-CNN, MaskR-CNN, TableDet, DiffusionDet, Deformable-DETR, and SparseR-CNN for 120 epochs, and other one-stage detectors for 300 epochs. The detailed settings of these benchmark models are included in section Appendix A. Our proposed SparseTableDet is built on the code base of SparseR-CNN, and the key training parameters are summarized in Table 1. It is worth mentioning that the parameter names in Table 1 are aligned with the names in Detectron2 [50]. More specifically, IMS_PER_BATCH is the total number of training samples in an iteration. MAX_ITER and STEPS refer to the total number of mini-batch used in the training and the mini-batch to apply the learning rate scheduler, respectively. WEIGHTS is the initialization weight of the model. In our implementation, we use the SparseR-CNN model pre-trained with COCO dataset [44] as the initialization weight. NUM_HEADS and LABEL_N are two custom parameters in the proposed SparseTableDet, which refer to the number of Dynamic Head and the hyper parameter N defined by Equation 3, as discussed in section 3. At last, the model is trained with AdamW [51] optimizer.

Table 2: Experimental results on ICDAR2017 dataset.

Model	F1				Weighted Average F1
	IoU(60%)	IoU(70%)	IoU(80%)	IoU(90%)	
RetinaNet	96.0	93.4	91.7	87.3	91.6
FCOS	96.0	93.8	92.1	87.6	91.9
YOLOX-X	97.7	95.5	92.3	80.0	90.4
YOLOR-X	97.1	95.3	93.4	89.7	93.5
YOLOV5-X	98.5	96.8	95.4	91.9	95.3
YOLOV7-X	97.6	96.3	94.6	91.5	94.6
YOLOV8-X	97.9	96.2	95.3	92.5	95.2
FasterR-CNN	97.1	96.0	93.8	89.6	93.7
MaskR-CNN	96.8	96.0	94.8	91.1	94.4
TableDet	98.8	97.1	95.0	90.4	94.9
DiffusionDet	98.3	97.0	94.9	90.3	94.7
Deformable-DETR	97.5	96.7	94.4	91.4	94.6
SparseR-CNN	98.3	97.9	96.1	94.0	96.3
SparseTableDet (Proposed)	99.5	99.4	98.2	94.8	97.7

The experimental results for the ICDAR2017, ICDAR2019, TNCR and ICT-TD datasets are shown in Table 2, 3, 4 and 5, respectively. The experimental results show that the proposed SparseTableDet can consistently outperform the state-of-the-art benchmark models regarding the Weighted Average F1 score. We also compare our proposed model with other state-of-the-art models optimized for the TD problem following the evaluation protocols of ICDAR2013, ICDAR2017, and ICT-TD datasets, and include the results in section Appendix A. With these competition evaluation protocols, our proposed method can still consistently outperform the benchmark models.

4.2. ICS for model training and evaluation

As discussed in section 3.4, GT_Coverage term in the ICS is a direct metric to measure whether the prediction box covers all the target content. This section uses the GT_Coverage as the evaluation metric to evaluate the model performance. More specifically, we replace the IoU score defined in Equation 8 with GT_Coverage to define a new Weighted Average F1 score as the evaluation metric, as defined by Equation 9. To demonstrate the effectiveness of ICS as the loss function, we replace the cost_giou in Equation 2 and GIoU loss in Equation 4 with their ICS-based counterparts. For simplicity, we use M_{giou} and M_{ics} to represent the model trained with GIoU loss and ICS loss, respectively. As shown in Table 6, when the IoU-based metrics are used, M_{giou} can perform better than M_{ics} . However, as aforementioned in section 1 and 3.4, GT_Coverage term defined in the ICS is a direct measure to

Table 3: Experimental results on ICDAR2019 dataset.

Model	F1				Weighted Average F1
	IoU(60%)	IoU(70%)	IoU(80%)	IoU(90%)	
RetinaNet	98.0	96.7	94.5	86.8	93.4
FCOS	97.6	96.5	93.6	85.7	92.7
YOLOX-X	97.1	96.0	94.6	89.2	93.8
YOLOR-X	98.6	98.2	97.2	93.4	96.6
YOLOV5-X	99.0	98.9	98.2	95.7	97.8
YOLOV7-X	99.2	98.6	98.0	94.1	97.2
YOLOV8-X	99.2	99.1	98.1	94.7	97.5
FasterR-CNN	97.4	96.2	95.0	90.4	94.4
MaskR-CNN	98.2	97.0	95.8	91.9	95.4
TableDet	98.1	96.8	94.9	91.5	94.9
DiffusionDet	98.9	97.4	95.8	91.3	95.5
Deformable-DETR	98.4	97.9	96.5	92.7	96.0
SparseR-CNN	98.6	98.1	97.5	94.9	97.1
SparseTableDet (Proposed)	99.3	99.1	98.9	96.3	98.3

evaluate the ground truth information covered by the prediction. GT_Coverage-based evaluation metric is used, M_{ics} can perform better. Figure 5 shows two prediction results of M_{giou} and M_{ics} . Using an ICS-based loss function can encourage the model to alleviate the information loss during the optimization process because the GT_Coverage term in the ICS is a direct measure of the information loss and is more sensitive to the information loss. We include some other prediction samples of these two models in Appendix A. Notably, bias might be introduced when using the GT Coverage score as the evaluation metric because the GT Coverage score cannot reflect the difference of prediction boxes once the predictions can cover the ground truth. However, the proposed ICS and GT Coverage scores can provide more insights regarding the quality of predictions and complement IoU-based metrics.

$$\text{Weighted Avg. F1} = \frac{\sum_{i=1}^4 GT_Coverage_i \cdot F1@GT_Coverage_i}{\sum_{i=1}^4 GT_Coverage_i} \quad (9)$$

4.3. Ablation Study

This section discusses the effectiveness of the proposed Image Size regional proposals, Noise Augmented Proposals, and Many-to-One label assignment method. We use SparseR-CNN as the baseline model in this section, which uses Hungarian Matching for the label assignment and random region proposals. We use the ICDAR2019 dataset to conduct experiments and use the Weighted Average F1 scores defined by Equation 8 and Equation 9 as the evaluation

Table 4: Experimental results on TNCR dataset.

Model	F1				Weighted Average F1
	IoU(60%)	IoU(70%)	IoU(80%)	IoU(90%)	
RetinaNet	92.7	92.0	90.6	84.8	89.6
FCOS	90.8	89.9	88.8	83.3	87.8
YOLOX-X	90.6	89.3	86.1	79.6	85.8
YOLOR-X	94.2	93.4	91.8	86.4	91.0
YOLOV5-X	95.8	95.5	94.	89.6	93.5
YOLOV7-X	95.2	95.0	93.7	89.3	93.0
YOLOV8-X	96.1	95.5	94.6	90.1	93.7
FasterR-CNN	91.5	91.0	90.3	84.4	88.9
MaskR-CNN	92.5	92.2	90.9	84.7	89.6
TableDet	94.7	94.4	93.3	87.7	92.2
Deformable-DETR	94.4	94.1	92.9	89.3	92.4
DiffusionDet	95.4	94.6	93.1	88.5	92.5
SparseR-CNN	95.1	94.9	94.4	90.9	93.6
SparseTableDet (Proposed)	96.3	96.2	95.8	92.7	95.1

metrics. It is worth mentioning that we choose 60%, 70%, 80% and 90% as thresholds to align with the metric in section 4.1. The experimental results are shown in Table 7 and 8, where SparseR-CNN(R), SparseR-CNN(I) are the SparseR-CNN initialized with the random proposals and image size proposals, respectively. ManytoOne and Noise represent the proposed many-to-one label assignment and the Noise Augmentation to regional proposals. The experimental results show that using image size region proposals, adding noise to the regional proposals, and Many-to-One label assignment can improve the performance of the base SparseR-CNN model.

5. Conclusion

In this study, we propose to use SparseR-CNN [5] as the base model and further improve the model by introducing Noise Augmented region proposal generation, Many-to-One label assignment, and a decoupled IoU. The experimental results show that the proposed method can consistently outperform benchmark models regarding the Weighted Average F1 score on various datasets. Furthermore, considering the requirement of TD applications, we propose to use GT_Coverage in ICS to replace IoU to act as the evaluation metric and use ICS to replace IoU to derive ICS-based loss functions. The experimental results demonstrate that the GT_Coverage can be a better metric reflecting the prediction’s information loss, and ICS-based loss can guide models to cover more information of the target objects. In this study, we assume that all the area in a ground truth box contains information without considering the inner

Table 5: Experimental results on the ICT-TD dataset.

Model	F1				Weighted Average F1
	IoU(60%)	IoU(70%)	IoU(80%)	IoU(90%)	
RetinaNet	95.8	93.6	91.0	83.7	90.4
FCOS	91.8	90.4	87.9	82.3	87.6
YOLOX-X	95.8	93.6	90.1	81.6	89.5
YOLOR-X	97.5	96.0	94.3	89.0	93.8
YOLOV5-X	98.0	97.2	95.8	91.7	95.3
YOLOV7-X	98.6	97.6	95.7	92.6	95.8
YOLOV8-X	97.9	97.2	95.6	92.3	95.4
FasterR-CNN	96.8	94.7	92.9	86.8	92.3
MaskR-CNN	96.2	94.8	92.8	87.9	92.5
TableDet	96.9	95.7	93.6	89.1	93.4
DiffusionDet	97.6	96.8	95.5	91.1	94.9
Deformable-DETR	97.4	96.5	95.0	91.2	94.7
SparseR-CNN	97.1	95.9	94.3	90.4	94.1
SparseTableDet (Proposed)	98.2	97.9	97.2	94.2	96.7

structure of tables. However, some tables contain extra spaces, meaning that some smaller prediction boxes than their ground truth boxes do not lead to any information loss. Therefore, it can be a direction to consider the inner structure of a table to build more reliable evaluation metrics for the TD applications. Besides, as aforementioned in section 4.2, the proposed GT_Coverage score cannot reflect the difference of box size once the prediction box can cover the whole ground truth. Therefore, it can be another direction to integrate the size of boxes to the GT_Coverage score to make it more versatile.

References

- [1] B. Xiao, M. Simsek, B. Kantarci, A. A. Alkheir, Handling big tabular data of ict supply chains: a multi-task, machine-interpretable approach, in: GLOBECOM 2022-2022 IEEE Global Communications Conference, IEEE, 2022, pp. 504–509.
- [2] B. Xiao, Y. Akkaya, M. Simsek, B. Kantarci, A. A. Alkheir, Efficient information sharing in ict supply chain social network via table structure recognition, in: GLOBECOM 2022-2022 IEEE Global Communications Conference, IEEE, 2022, pp. 4661–4666.
- [3] J. Fernandes, M. Simsek, B. Kantarci, S. Khan, Tabledet: An end-to-end deep learning approach for table detection and table image classification in data sheet images, *Neurocomputing* 468 (2022) 317–334.
- [4] D. Prasad, A. Gadpal, K. Kapadni, M. Visave, K. Sultanpure, Cascadetabnet: An approach for end to end table detection and structure recognition from image-based documents, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops, 2020, pp. 572–573.

表 1: 港股有色金属行业龙头业绩增速与估值

代码	名称	净利润同比增长		收入同比增长		市盈率 TTM
		Y17	Y16	Y17	Y16	
3993.HK	洛阳钼业	173.32%	31.12%	254.25%	69.92%	10.86
2899.HK	紫金矿业	90.66%	11.12%	19.57%	6.05%	13.39
2600.HK	中国铝业	274.16%	170.82%	25.67%	16.83%	21.39
1378.HK	中国宏桥	-25.27%	84.81%	51.99%	39.19%	5.99
0486.HK	赣锋	3.65%	111.29%	24.88%	-8.03%	2.70
0358.HK	江西铜业股份	96.19%	21.93%	1.24%	8.91%	13.28
0847.HK	哈钨矿业-S	152.54%	1575.00%	117.10%	15.19%	0.00
1208.HK	五矿资源	196.33%	85.12%	66.47%	27.57%	13.04
1333.HK	中国忠旺	23.06%	2.37%	16.55%	3.24%	4.16
1818.HK	招金矿业	82.26%	14.66%	0.14%	13.21%	39.78
1636.HK	中国金矿集团	159.39%	37.72%	178.53%	175.51%	60.04
0639.HK	贵研资源	866.63%	126.84%	91.83%	-9.35%	7.88
2303.HK	陕西黄金	18.28%	279.13%	31.03%	157.61%	19.45
2362.HK	金川国际	398.67%	102.86%	50.53%	-22.49%	10.67
1258.HK	中国有色矿业	1103.75%	104.23%	40.01%	10.44%	5.29
2099.HK	中国黄金国际	574.64%	-62.48%	21.64%	-0.40%	30.95
0976.HK	齐合环保	197.02%	61.37%	475.79%	2.38%	5.87
1021.HK	龙达新材-S	0.00%	76.29%	0.00%	-1.75%	30.67
2326.HK	新疆广汇能源	0.00%	-62.67%	0.00%	-17.49%	7.17
1164.HK	中广核矿业	-86.62%	30.59%	-47.33%	0.62%	19.42
	均值	214.93%	140.11%	71.00%	24.36%	16.10
	中位数	124.37%	68.83%	28.35%	7.48%	11.95
	最大值	1103.75%	1575.00%	475.79%	175.51%	60.04
	最小值	-86.62%	-62.67%	-47.33%	-22.49%	0.00

资料来源: wind, 中国银河证券研究院

(a) ICS-based loss

表 2: 港股有色金属行业龙头业绩增速与估值

代码	名称	净利润同比增长		收入同比增长		市盈率 TTM
		Y17	Y16	Y17	Y16	
3993.HK	洛阳钼业	173.32%	31.12%	254.25%	69.92%	10.86
2899.HK	紫金矿业	90.66%	11.12%	19.57%	6.05%	13.39
2600.HK	中国铝业	274.16%	170.82%	25.67%	16.83%	21.39
1378.HK	中国宏桥	-25.27%	84.81%	51.99%	39.19%	5.99
0486.HK	赣锋	3.65%	111.29%	24.88%	-8.03%	2.70
0358.HK	江西铜业股份	96.19%	21.93%	1.24%	8.91%	13.28
0847.HK	哈钨矿业-S	152.54%	1575.00%	117.10%	15.19%	0.00
1208.HK	五矿资源	196.33%	85.12%	66.47%	27.57%	13.04
1333.HK	中国忠旺	23.06%	2.37%	16.55%	3.24%	4.16
1818.HK	招金矿业	82.26%	14.66%	0.14%	13.21%	39.78
1636.HK	中国金矿集团	159.39%	37.72%	178.53%	175.51%	60.04
0639.HK	贵研资源	866.63%	126.84%	91.83%	-9.35%	7.88
2303.HK	陕西黄金	18.28%	279.13%	31.03%	157.61%	19.45
2362.HK	金川国际	398.67%	102.86%	50.53%	-22.49%	10.67
1258.HK	中国有色矿业	1103.75%	104.23%	40.01%	10.44%	5.29
2099.HK	中国黄金国际	574.64%	-62.48%	21.64%	-0.40%	30.95
0976.HK	齐合环保	197.02%	61.37%	475.79%	2.38%	5.87
1021.HK	龙达新材-S	0.00%	76.29%	0.00%	-1.75%	30.67
2326.HK	新疆广汇能源	0.00%	-62.67%	0.00%	-17.49%	7.17
1164.HK	中广核矿业	-86.62%	30.59%	-47.33%	0.62%	19.42
	均值	214.93%	140.11%	71.00%	24.36%	16.10
	中位数	124.37%	68.83%	28.35%	7.48%	11.95
	最大值	1103.75%	1575.00%	475.79%	175.51%	60.04
	最小值	-86.62%	-62.67%	-47.33%	-22.49%	0.00

资料来源: wind, 中国银河证券研究院

(b) IoU-based loss

Figure 5: Prediction samples of models trained with ICS-based loss and IoU-based loss.

Table 6: Experimental results on ICDAR2019 dataset evaluated by Weighted Average F1 scores using GT_Coverage and IoU as thresholds.

Loss Function	Metric	F1				Weighted Average F1
		60%	70%	80%	90%	
GIoU	IoU	99.3	99.1	98.9	96.3	98.3
ICS	IoU	99.4	98.8	98.1	92.9	97.0
GIoU	GT_C	99.6	99.5	99.1	98.3	99.1
ICS	GT_C	99.7	99.6	99.5	98.6	99.3

- [5] P. Sun, R. Zhang, Y. Jiang, T. Kong, C. Xu, W. Zhan, M. Tomizuka, L. Li, Z. Yuan, C. Wang, et al., Sparse r-cnn: End-to-end object detection with learnable proposals, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2021, pp. 14454–14463.
- [6] Z. Ge, S. Liu, Z. Li, O. Yoshie, J. Sun, Ota: Optimal transport assignment for object detection, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 303–312.
- [7] Q. Hong, F. Liu, D. Li, J. Liu, L. Tian, Y. Shan, Dynamic sparse r-cnn, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 4723–4732.
- [8] S. Chen, P. Sun, Y. Song, P. Luo, Diffusiondet: Diffusion model for object detection, arXiv preprint arXiv:2211.09788 (2022).
- [9] Y. Du, W. Guo, Y. Xiao, V. Lepetit, 1st place solution for the uvo challenge on image-based open-world segmentation 2021, arXiv preprint arXiv:2110.10239 (2021).
- [10] Z. Ge, S. Liu, F. Wang, Z. Li, J. Sun, Yolox: Exceeding yolo series in 2021, arXiv preprint arXiv:2107.08430 (2021).
- [11] J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You only look once: Unified, real-time object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 779–788.
- [12] M. J. Shafiee, B. Chywl, F. Li, A. Wong, Fast yolo: A fast you only look once system for real-time embedded object detection in video, arXiv preprint arXiv:1709.05943 (2017).
- [13] J. Redmon, A. Farhadi, Yolov3: An incremental improvement, arXiv preprint arXiv:1804.02767 (2018).

Table 7: The effectiveness of each component using IoU scores as thresholds.

	F1				Weighted Average F1
	IoU(60%)	IoU(70%)	IoU(80%)	IoU(90%)	
SparseR-CNN (R)	98.6	98.1	97.5	94.9	97.1
SparseR-CNN (I)	99.1	98.7	98.2	95.3	97.6
I+ManytoOne	99.4	99.1	98.8	95.3	97.9
I+Noise+ManytoOne	99.3	99.1	98.9	96.3	98.3

Table 8: The effectiveness of each component using GT Coverage scores as thresholds.

	F1				Weighted Average F1
	GT_C (60%)	GT_C(70%)	GT_C(80%)	GT_C(90%)	
SparseR-CNN (R)	98.7	98.7	97.9	96.7	97.9
SparseR-CNN (I)	99.1	99.0	98.3	97.6	98.4
I+ManytoOne	99.4	99.4	99.2	98.1	98.9
I+Noise+ManytoOne	99.6	99.5	99.1	98.3	99.1

- [14] A. Bochkovskiy, C.-Y. Wang, H.-Y. M. Liao, Yolov4: Optimal speed and accuracy of object detection, arXiv preprint arXiv:2004.10934 (2020).
- [15] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A. C. Berg, Ssd: Single shot multibox detector, in: European conference on computer vision, Springer, 2016, pp. 21–37.
- [16] Z. Tian, C. Shen, H. Chen, T. He, Fcos: Fully convolutional one-stage object detection, in: Proceedings of the IEEE/CVF international conference on computer vision, 2019, pp. 9627–9636.
- [17] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, Advances in neural information processing systems 28 (2015).
- [18] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask r-cnn, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2961–2969.
- [19] Z. Cai, N. Vasconcelos, Cascade r-cnn: Delving into high quality object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 6154–6162.
- [20] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, S. Zagoruyko, End-to-end object detection with transformers, in: European conference on computer vision, Springer, 2020, pp. 213–229.
- [21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, Advances in neural information processing systems 30 (2017).
- [22] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, J. Dai, Deformable detr: Deformable transformers for end-to-end object detection, arXiv preprint arXiv:2010.04159 (2020).
- [23] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, L. Zhang, Dab-detr: Dynamic anchor boxes are better queries for detr, arXiv preprint arXiv:2201.12329 (2022).
- [24] R. Girshick, Fast r-cnn, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 1440–1448.
- [25] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, S. Savarese, Generalized intersection over union: A metric and a loss for bounding

- box regression, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 658–666.
- [26] T. Zheng, S. Zhao, Y. Liu, Z. Liu, D. Cai, Scaloss: Side and corner aligned loss for bounding box regression, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 36, 2022, pp. 3535–3543.
 - [27] J. Yu, Y. Jiang, Z. Wang, Z. Cao, T. Huang, Unitbox: An advanced object detection network, in: Proceedings of the 24th ACM international conference on Multimedia, 2016, pp. 516–520.
 - [28] Z. Zheng, P. Wang, W. Liu, J. Li, R. Ye, D. Ren, Distance-iou loss: Faster and better learning for bounding box regression, in: Proceedings of the AAAI conference on artificial intelligence, Vol. 34, 2020, pp. 12993–13000.
 - [29] Z. Zheng, P. Wang, D. Ren, W. Liu, R. Ye, Q. Hu, W. Zuo, Enhancing geometric factors in model learning and inference for object detection and instance segmentation, *IEEE Transactions on Cybernetics* 52 (8) (2021) 8574–8586.
 - [30] Y. He, C. Zhu, J. Wang, M. Savvides, X. Zhang, Bounding box regression with uncertainty for accurate object detection, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 2888–2897.
 - [31] A. Abdallah, A. Berendeyev, I. Nuradin, D. Nurseitov, Tncr: Table net detection and classification dataset, *Neurocomputing* 473 (2022) 79–97.
 - [32] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang, et al., Deep high-resolution representation learning for visual recognition, *IEEE transactions on pattern analysis and machine intelligence* 43 (10) (2020) 3349–3364.
 - [33] Y. Huang, Q. Yan, Y. Li, Y. Chen, X. Wang, L. Gao, Z. Tang, A yolo-based table detection method, in: 2019 International Conference on Document Analysis and Recognition (ICDAR), IEEE, 2019, pp. 813–818.
 - [34] B. Smock, R. Pesala, R. Abraham, Pubtables-1m: Towards comprehensive table extraction from unstructured documents, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 4634–4642.
 - [35] S. Schreiber, S. Agne, I. Wolf, A. Dengel, S. Ahmed, Deepdesrt: Deep learning for detection and structure recognition of tables in document images, in: 2017 14th IAPR international conference on document analysis and recognition (ICDAR), Vol. 1, IEEE, 2017, pp. 1162–1167.
 - [36] S. A. Siddiqui, M. I. Malik, S. Agne, A. Dengel, S. Ahmed, Decnt: Deep deformable cnn for table detection, *IEEE access* 6 (2018) 74151–74161.
 - [37] E. Kara, M. Traquair, M. Simsek, B. Kantarci, S. Khan, Holistic design for deep learning-based discovery of tabular structures in datasheet images, *Engineering Applications of Artificial Intelligence* 90 (2020) 103551.
 - [38] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, S. Belongie, Feature pyramid networks for object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 2117–2125.
 - [39] N. Bodla, B. Singh, R. Chellappa, L. S. Davis, Soft-nms—improving object detection with one line of code, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 5561–5569.
 - [40] T.-Y. Lin, P. Goyal, R. Girshick, K. He, P. Dollár, Focal loss for dense object detection, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2980–2988.
 - [41] L. Gao, X. Yi, Z. Jiang, L. Hao, Z. Tang, Icdar2017 competition on page object detection, in: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR), Vol. 1, IEEE, 2017, pp. 1417–1422.
 - [42] L. Gao, Y. Huang, H. Déjean, J.-L. Meunier, Q. Yan, Y. Fang, F. Kleber, E. Lang, Icdar 2019 competition on table detection and recognition (ctdar), in: 2019 International Conference on Document Analysis and Recognition (ICDAR), IEEE, 2019, pp. 1510–1515.
 - [43] B. Xiao, M. Simsek, B. Kantarci, A. A. Alkheir, Revisiting table detection datasets for visually rich documents, *arXiv preprint arXiv:2305.04833* (2023).
 - [44] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V* 13, Springer, 2014, pp. 740–755.
 - [45] M. Göbel, T. Hassan, E. Oro, G. Orsi, Icdar 2013 table competition, in: 2013 12th International Conference on Document Analysis and Recognition, IEEE, 2013, pp. 1449–1453.
 - [46] C.-Y. Wang, I.-H. Yeh, H.-Y. M. Liao, You only learn one representation: Unified network for multiple tasks, *arXiv preprint arXiv:2105.04206*

(2021).

- [47] G. Jocher, Ultralytics yolov5 (2020). doi:10.5281/zenodo.3908559.
URL <https://github.com/ultralytics/yolov5>
- [48] C.-Y. Wang, A. Bochkovskiy, H.-Y. M. Liao, YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors, arXiv preprint arXiv:2207.02696 (2022).
- [49] G. Jocher, A. Chaurasia, J. Qiu, Ultralytics yolov8 (2023).
URL <https://github.com/ultralytics/ultralytics>
- [50] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, R. Girshick, Detectron2, <https://github.com/facebookresearch/detectron2> (2019).
- [51] I. Loshchilov, F. Hutter, Decoupled weight decay regularization, arXiv preprint arXiv:1711.05101 (2017).
- [52] S. S. Paliwal, D. Vishwanath, R. Rahul, M. Sharma, L. Vig, Tablenet: Deep learning model for end-to-end table detection and tabular data extraction from scanned document images, in: 2019 International Conference on Document Analysis and Recognition (ICDAR), IEEE, 2019, pp. 128–133.
- [53] Y. Li, L. Gao, Z. Tang, Q. Yan, Y. Huang, A gan-based feature generator for table detection, in: 2019 International Conference on Document Analysis and Recognition (ICDAR), IEEE, 2019, pp. 763–768.

Appendix A. Appendix

Appendix A.1. Model implementations and settings

In this section, we list the implementations and configuration files of the baseline models, including RetinaNet [40], FCOS [16], YOLOX-X [10], YOLOR-X [46], YOLOv5-X [47], YOLOv7-X [48], YOLOv8-X [49], FasterR-CNN [17], DiffusionDet [8], Deformable-DETR [22], and SparseR-CNN [5], as summarized in Table A.9. It is worth mentioning that we modified the training epochs of the listed configuration files, trained FasterR-CNN, MaskR-CNN, DiffusionDet, Deformable-DETR, and SparseR-CNN for 120 epochs, and trained other one-stage detectors for 300 epochs.

Table A.9: Summary of model implementations and settings

Model	Implementation	Setting File
RetinaNet	Detectron2	https://github.com/facebookresearch/detectron2/blob/main/configs/COCO-Keypoint.yaml
FCOS	Detectron2	https://github.com/facebookresearch/detectron2/blob/main/configs/COCO-Keypoint.yaml
YOLOX	Official codebase	https://github.com/Megvii-BaseDetection/YOLOX/blob/main/exps/default.yaml
YOLOR	Official codebase	https://github.com/WongKinYiu/yolor/blob/main/cfg/yolor_csp_x.cfg
YOLOv5	Official codebase	https://github.com/ultralytics/ultralytics/blob/main/ultralytics/cfg/default.yaml
YOLOv7	Official codebase	https://github.com/WongKinYiu/yolov7/blob/main/cfg/training/yolov7x.yaml
YOLOv8	Official codebase	https://github.com/ultralytics/ultralytics/blob/main/ultralytics/cfg/default.yaml
FasterR-CNN	Detectron2	https://github.com/facebookresearch/detectron2/blob/main/configs/COCO-Keypoint.yaml
MaskR-CNN	Detectron2	https://github.com/facebookresearch/detectron2/blob/main/configs/COCO-Keypoint.yaml
DiffusionDet	Official codebase	https://github.com/ShoufaChen/DiffusionDet/blob/main/configs/diffdet.yaml
Deformable-DETR	detrex	https://github.com/IDEA-Research/detrex/blob/main/projects/deformable_deetr/configs/deformable_deetr.yaml

Table A.10: Experimental results on ICDAR2013 dataset (IoU = 50%).

Model	Precision	Recall	F1
CascadeTabNet [4]	100	100	100
TableDet [3]	100	100	100
DeCNT [36]	99.6	99.6	99.6
YOLOv3-TD [33]	94.9	100	97.3
DeepDeSRT [35]	97.4	96.2	96.8
TableNet [52]	97.0	96.3	96.6
SparseTableDet (Proposed)	100	100	100

SparseR-CNN	Official codebase	https://github.com/PeizeSun/SparseR-CNN/blob/main/projects/SparseRCNN
-------------	-------------------	---

Appendix A.2. Compared with other Table Detection models

In this section, we include the experimental results using the evaluation protocols in ICDAR2013, ICDAR2017, and ICT-TD datasets. More specifically, for the ICDAR2013 dataset, the F1 score thresholded by 50% is the competition evaluation metric. For the ICDAR2017 dataset, Precision, Recall, and F1 scores thresholded by 60% and 80% are used as evaluation metrics. For the ICT-TD dataset, Weighted Average F1 score, as defined in Equation 8, is used as the evaluation metric whose thresholds are 80%, 85%, 90%, and 95%. The experimental results of them are shown in Table A.10, A.11 and A.12. It is worth mentioning that the experimental results of TableDet [3], DeCNT [36], YOLOv3-TD [33], DeepDeSRT [35], TableNet [52], GAN-TD [53], in Table A.10, A.11 are from study [3].

Appendix A.3. Detailed experimental results

In this appendix section, we include the detailed experimental results on the TNCR and ICT-TD datasets, as shown in Table A.15 and Table A.16. Besides, we also include some prediction results for the models trained with IoU-based and ICS-based losses, as discussed in section 4.2.

Table A.13: Detailed Experimental results on the ICDAR2017 dataset.

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
RetinaNet	Precision	96.4	96.4	95.2	94.9	92.1	90.7	90.3	88.1	85.5	75.0	90.4
	Recall	97.8	97.8	96.9	96.3	94.7	93.8	93.2	91.6	89.1	80.1	93.1
	F1	97.1	97.1	96.0	95.6	93.4	92.2	91.7	89.8	87.3	77.5	91.7

Continued on next page

Table A.13 Detailed Experimental results on the ICDAR2017 dataset (continued from previous page).

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
FCOS	Precision	97.3	96.3	95.2	94.9	92.4	90.8	90.2	87.6	84.9	72.7	90.2
	Recall	98.1	97.5	96.9	96.6	95.3	94.7	94.1	92.8	90.3	82.6	93.9
	F1	97.7	96.9	96.0	95.7	93.8	92.7	92.1	90.1	87.5	77.3	92.0
YOLOX-X	Precision	97.0	97.0	96.3	95.5	94.1	93.3	91.2	87.1	78.8	54.2	88.5
	Recall	99.4	99.4	99.1	98.4	96.9	95.6	93.5	89.4	81.3	59.5	91.3
	F1	98.2	98.2	97.7	96.9	95.5	94.4	92.3	88.2	80.0	56.7	89.9
YOLOR-X	Precision	97.5	97.4	96.4	95.1	94.3	93.4	92.2	92.2	88.1	79.9	92.6
	Recall	98.8	98.8	97.8	96.9	96.3	95.3	94.7	94.4	91.3	84.1	94.8
	F1	98.1	98.1	97.1	96.0	95.3	94.3	93.4	93.3	89.7	81.9	93.7
YOLOV5-X	Precision	98.3	98.3	97.9	97.9	95.8	95.4	94.2	93.8	90.2	82.8	94.4
	Recall	99.7	99.7	99.1	99.1	97.8	97.2	96.6	96.3	93.8	86.6	96.6
	F1	99.0	99.0	98.5	98.5	96.8	96.3	95.4	95.0	91.9	84.7	95.5
YOLOV7-X	Precision	98.0	97.1	97.0	95.9	95.3	94.1	93.3	93.3	89.8	80.7	93.5
	Recall	99.1	98.4	98.1	97.5	97.2	96.6	96.0	95.6	93.2	85.1	95.7
	F1	98.5	97.7	97.5	96.7	96.2	95.3	94.6	94.4	91.5	82.8	94.6
YOLOV8-X	Precision	99.0	97.9	97.1	96.5	95.2	94.6	94.1	93.9	91.0	80.6	94.0
	Recall	100.0	99.1	98.8	98.1	97.2	96.9	96.6	96.3	94.1	85.1	96.2
	F1	99.5	98.5	97.9	97.3	96.2	95.7	95.3	95.1	92.5	82.8	95.1
FasterR-CNN	Precision	97.8	96.9	96.7	96.6	95.4	94.0	92.9	91.3	88.2	79.4	92.9
	Recall	98.1	97.8	97.5	97.2	96.6	95.3	94.7	93.8	91.0	83.2	94.5
	F1	97.9	97.3	97.1	96.9	96.0	94.6	93.8	92.5	89.6	81.3	93.7
MaskR-CNN	Precision	97.5	96.5	96.2	96.2	95.1	94.0	93.9	92.6	89.8	68.5	92.0
	Recall	98.1	97.8	97.5	97.2	96.9	96.0	95.6	94.7	92.5	76.6	94.3
	F1	97.8	97.1	96.8	96.7	96.0	95.0	94.7	93.6	91.1	72.3	93.1
TableDet	Precision	99.4	98.6	98.5	98.5	96.4	95.1	94.1	93.5	89.0	78.8	94.2
	Recall	100.0	99.7	99.1	99.1	97.8	96.6	96.0	95.0	91.9	84.1	95.9
	F1	99.7	99.1	98.8	98.8	97.1	95.8	95.0	94.2	90.4	81.4	95.0
DiffusionDet	Precision	98.1	98.0	97.6	97.1	96.1	94.6	93.9	92.5	89.1	77.8	93.5
	Recall	99.7	99.7	99.1	98.8	97.8	96.3	96.0	94.4	91.6	83.5	95.7

Continued on next page

Table A.13 Detailed Experimental results on the ICDAR2017 dataset (continued from previous page).

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
	F1	98.9	98.8	98.3	97.9	96.9	95.4	94.9	93.4	90.3	80.5	94.6
Deformable-DETR	Precision	97.2	97.2	96.8	96.8	95.8	93.7	92.6	91.3	89.7	78.6	93.0
	Recall	98.8	98.4	98.1	98.1	97.5	96.6	96.3	95.0	93.2	83.2	95.5
	F1	98.0	97.8	97.4	97.4	96.6	95.1	94.4	93.1	87.1	80.8	94.2
SparseR-CNN	Precision	99.3	98.3	97.8	97.8	97.0	95.9	94.8	94.8	92.7	84.6	95.3
	Recall	100.0	99.7	99.1	99.1	98.8	98.1	97.5	97.2	95.3	89.1	97.4
	F1	99.6	99.0	98.4	98.4	97.9	97.0	96.1	96.0	94.0	86.8	96.3
SparseTableDet (Proposed)	Precision	99.7	99.7	99.1	99.1	98.9	97.5	96.7	95.6	93.1	83.6	96.3
	Recall	100.0	100.0	100.0	100.0	100.0	100.0	99.7	99.1	96.6	89.1	98.4
	F1	99.8	99.8	99.5	99.5	99.4	98.7	98.2	97.3	94.8	86.3	97.3

Table A.14: Detailed Experimental results on the ICDAR2019 dataset.

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
RetinaNet	Precision	98.6	98.4	97.4	97.2	96.2	94.3	93.9	90.5	85.3	72.8	92.5
	Recall	99.6	99.1	98.7	98.2	97.3	96.0	95.1	92.7	88.4	78.6	94.4
	F1	99.1	98.7	98.0	97.7	96.7	95.1	94.5	91.6	86.8	75.6	93.4
FCOS	Precision	97.1	96.7	96.7	95.7	95.5	94.4	92.2	90.2	83.1	63.7	90.5
	Recall	98.9	98.4	98.4	97.8	97.6	96.9	95.1	93.5	88.4	75.7	94.1
	F1	98.0	97.5	97.5	96.7	96.5	95.6	93.6	91.8	85.7	69.2	92.3
YOLOX-X	Precision	97.5	97.1	96.3	95.7	95.3	94.3	94.1	92.2	88.3	70.4	92.1
	Recall	98.7	98.4	98.0	97.1	96.7	96.0	95.1	93.3	90.2	74.4	93.8
	F1	98.1	97.7	97.1	96.4	96.0	95.1	94.6	92.7	89.2	72.3	92.9
YOLOR-X	Precision	98.7	98.7	98.2	98.2	97.5	97.5	96.6	95.6	92.5	82.1	95.5
	Recall	99.6	99.6	99.1	99.1	98.9	98.7	97.8	96.4	94.2	85.3	96.9
	F1	99.1	99.1	98.6	98.6	98.2	98.1	97.2	96.0	93.3	83.7	96.2
YOLOV5-X	Precision	98.7	98.6	98.5	98.5	98.5	98.4	97.5	96.5	95.1	83.7	96.4
	Recall	99.8	99.8	99.6	99.3	99.3	99.3	98.9	97.8	96.4	86.6	97.7

Continued on next page

Table A.14 Detailed Experimental results on the ICDAR2019 dataset (continued from previous page).

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
YOLOV7-X	F1	99.2	99.2	99.0	98.9	98.9	98.8	98.2	97.1	95.7	85.1	97.0
	Precision	99.5	98.6	98.6	98.4	98.1	97.7	97.6	96.7	93.3	81.0	95.9
	Recall	100.0	99.8	99.8	99.1	99.1	98.7	98.4	97.8	94.9	84.9	97.2
	F1	99.7	99.2	99.2	98.7	98.6	98.2	98.0	97.2	94.1	82.9	96.5
YOLOV8-X	Precision	99.0	99.0	98.8	98.8	98.8	98.7	97.7	96.8	94.0	89.5	97.1
	Recall	99.8	99.8	99.6	99.3	99.3	99.1	98.4	98.0	95.3	92.0	98.1
	F1	99.4	99.4	99.2	99.0	99.0	98.9	98.0	97.4	94.6	90.7	97.6
FasterR-CNN	Precision	97.9	97.8	96.8	96.8	95.6	94.5	94.5	93.5	89.7	76.0	93.3
	Recall	98.7	98.4	98.0	97.6	96.9	95.8	95.6	94.7	91.1	80.2	94.7
	F1	98.3	98.1	97.4	97.2	96.2	95.1	95.0	94.1	90.4	78.0	94.0
MaskR-CNN	Precision	98.9	97.8	97.7	96.5	96.4	95.4	95.4	94.5	91.2	71.2	93.5
	Recall	99.3	98.9	98.7	98.0	97.6	96.9	96.2	95.8	92.7	76.6	95.1
	F1	99.1	98.3	98.2	97.2	97.0	96.1	95.8	95.1	91.9	73.8	94.3
TableDet	Precision	98.5	97.5	97.5	97.4	96.3	95.3	94.4	93.5	90.7	77.2	93.8
	Recall	99.1	98.9	98.7	98.4	97.3	96.4	95.3	94.4	92.2	83.7	95.5
	F1	98.8	98.2	98.1	97.9	96.8	95.8	94.8	93.9	91.4	80.3	94.6
DiffusionDet	Precision	98.5	98.3	98.3	97.6	96.4	96.0	95.0	92.6	90.3	74.4	93.7
	Recall	99.8	99.6	99.6	99.3	98.4	97.6	96.7	94.4	92.4	82.2	96.0
	F1	99.1	98.9	98.9	98.4	97.4	96.8	95.8	93.5	91.3	78.1	94.8
Deformable-DETR	Precision	98.7	97.9	97.6	97.1	96.9	96.5	95.3	94.1	91.0	80.2	94.5
	Recall	99.6	99.3	99.1	98.9	98.9	98.7	97.8	96.7	94.4	86.6	97.0
	F1	99.1	98.6	98.3	98.0	97.9	97.6	96.5	95.4	92.7	83.3	95.7
SparseR-CNN	Precision	98.4	97.9	97.9	97.9	97.1	96.5	96.4	96.1	93.5	86.0	95.8
	Recall	99.8	99.6	99.3	99.3	99.1	98.9	98.7	98.2	96.2	91.8	98.1
	F1	99.1	98.7	98.6	98.6	98.1	97.7	97.5	97.1	94.8	88.8	96.9
SparseTableDet (Proposed)	Precision	99.5	98.8	98.8	98.8	98.6	98.5	98.3	97.5	95.2	88.6	97.3
	Recall	100.0	99.8	99.6	99.6	99.6	99.6	99.6	98.9	97.3	92.0	98.6
	F1	99.7	99.3	99.2	99.2	99.1	99.0	98.9	98.2	96.3	90.3	97.9

Table 9-7. No-Decompression Limits and Repetitive Group Designators for No-Decompression Air Dives.

0 100% (fsw)	No-Stop Limit	Repetitive Group Designation															
		A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	Z
10	Unlimited	57	101	158	245	426	*										
15	Unlimited	36	60	88	121	163	217	297	449	*							
20	Unlimited	26	43	61	82	106	133	165	205	256	330	461	*				
25	595	20	33	47	62	78	97	117	140	166	198	236	285	354	469	595	
30	371	17	27	38	50	62	76	91	107	125	145	167	193	223	260	307	371
35	232	14	23	32	42	52	63	74	87	100	115	131	148	168	190	215	232
40	163	12	20	27	36	44	53	63	73	84	95	108	121	135	151	163	
45	125	11	17	24	31	39	46	55	63	72	82	92	102	114	125		
50	92	9	15	21	28	34	41	48	56	63	71	80	89	92			
55	74	8	14	19	25	31	37	43	50	56	63	71	74				
60	60	7	12	17	22	28	33	39	45	51	57	60					
70	48	6	10	14	19	23	28	32	37	42	47	48					
80	39	5	9	12	16	20	24	28	32	36	39						
90	30	4	7	11	14	17	21	24	28	30							
100	25	4	6	9	12	15	18	21	25								
110	20	3	6	8	11	14	16	19	20								
120	15	3	5	7	10	12	15										
130	10	2	4	6	9	10											
140	10	2	4	6	8	10											
150	5	2	3	5													
160	5	3	5														
170	5		4	5													
180	5		4	5													
190	5		3	5													

* Highest repetitive group that can be achieved at this depth regardless of bottom time.

(a) ICS-based loss

Table 9-7. No-Decompression Limits and Repetitive Group Designators for No-Decompression Air Dives.

0 100% (fsw)	No-Stop Limit	Repetitive Group Designation															
		A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	Z
10	Unlimited	57	101	158	245	426	*										
15	Unlimited	36	60	88	121	163	217	297	449	*							
20	Unlimited	26	43	61	82	106	133	165	205	256	330	461	*				
25	595	20	33	47	62	78	97	117	140	166	198	236	285	354	469	595	
30	371	17	27	38	50	62	76	91	107	125	145	167	193	223	260	307	371
35	232	14	23	32	42	52	63	74	87	100	115	131	148	168	190	215	232
40	163	12	20	27	36	44	53	63	73	84	95	108	121	135	151	163	
45	125	11	17	24	31	39	46	55	63	72	82	92	102	114	125		
50	92	9	15	21	28	34	41	48	56	63	71	80	89	92			
55	74	8	14	19	25	31	37	43	50	56	63	71	74				
60	60	7	12	17	22	28	33	39	45	51	57	60					
70	48	6	10	14	19	23	28	32	37	42	47	48					
80	39	5	9	12	16	20	24	28	32	36	39						
90	30	4	7	11	14	17	21	24	28	30							
100	25	4	6	9	12	15	18	21	25								
110	20	3	6	8	11	14	16	19	20								
120	15	3	5	7	10	12	15										
130	10	2	4	6	9	10											
140	10	2	4	6	8	10											
150	5	2	3	5													
160	5	3	5														
170	5		4	5													
180	5		4	5													
190	5		3	5													

* Highest repetitive group that can be achieved at this depth regardless of bottom time.

(b) IoU-based loss

Figure A.6: Prediction samples of models trained with ICS-based loss and IoU-based loss.

0 100%	Annual cap for the year ended December 31,		
	2019	2020	2021
	(RMB in thousands)		
(i) Pre-delivery property management services	18,200	23,650	35,500
(ii) Management and related services to the display units, sales offices and common area of our property projects	27,000	27,350	30,400
Total	45,200	51,000	65,900

(a) ICS-based loss

0 100%	Annual cap for the year ended December 31,		
	2019	2020	2021
	(RMB in thousands)		
(i) Pre-delivery property management services	18,200	23,650	35,500
(ii) Management and related services to the display units, sales offices and common area of our property projects	27,000	27,350	30,400
Total	45,200	51,000	65,900

(b) IoU-based loss

Figure A.7: Prediction samples of models trained with ICS-based loss and IoU-based loss.

Table A.11: Experimental results on ICDAR2017 dataset.

IoU Threshold	Model	Precision	Recall	F1
60%	TableDet [3]	98.8	99.7	99.3
	YOLOv3-TD [33]	97.2	97.8	97.5
	DeCNT [36]	96.5	97.1	96.8
	GAN-TD [53]	94.4	94.4	94.4
	SparseTableDet (Proposed)	99.1	100.0	99.5
80%	TableDet [3]	97.4	98.4	97.9
	YOLOv3-TD [33]	96.8	97.5	97.1
	DeCNT [36]	96.7	93.7	95.2
	GAN-TD [53]	90.3	90.3	90.3
	SparseTableDet (Proposed)	96.7	99.7	98.2

Table A.12: Experimental results on the ICT-TD dataset.

Model	F1				Weighted Average F1
	IoU(80%)	IoU(85%)	IoU(90%)	IoU(95%)	
TableDet [43]	93.6	91.6	89.1	75.7	87.1
DiffusionDet [43]	95.5	94.2	91.1	76.4	88.9
Deformable-DETR [43]	95.0	93.9	91.2	83.0	90.5
SparseR-CNN [43]	94.3	93.0	90.4	78.8	88.8
SparseTableDet (Proposed)	97.2	96.4	94.2	81.8	92.1

Table A.15: Detailed Experimental results on the TNCR dataset.

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
RetinaNet	Precision	89.7	89.7	89.6	89.4	88.9	88.6	87.5	85.6	81.5	69.8	86.0
	Recall	96.2	96.2	96.1	96.0	95.3	95.0	94.1	92.3	88.2	78.2	92.8
	F1	92.8	92.8	92.7	92.6	92.0	91.7	90.6	88.8	84.8	73.8	89.3
FCOS	Precision	87.8	87.7	87.6	87.3	86.7	86.3	85.6	83.1	79.4	68.5	84.0
	Recall	94.4	94.3	94.2	94.0	93.4	93.1	92.3	90.4	87.5	78.4	91.2
	F1	91.0	90.9	90.8	90.5	89.9	89.6	88.8	86.6	83.3	73.1	87.5

Continued on next page

Table A.15 Detailed Experimental results on the TNCR dataset (continued from previous page).

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
YOLOX-X	Precision	87.1	86.8	86.5	86.1	85.5	84.4	83.0	80.6	76.7	57.9	81.5
	Recall	96.0	95.6	95.0	94.2	93.4	91.8	89.5	86.7	82.7	65.9	89.1
	F1	91.3	91.0	90.6	90.0	89.3	87.9	86.1	83.5	79.6	61.6	85.1
YOLOR-X	Precision	90.4	90.3	90.2	90.1	89.6	89.2	88.2	86.0	82.9	75.1	87.2
	Recall	98.7	98.6	98.5	98.3	97.6	97.1	95.6	93.3	90.1	83.6	95.1
	F1	94.4	94.3	94.2	94.0	93.4	93.0	91.8	89.5	86.4	79.1	91.0
YOLOV5-X	Precision	93.0	92.8	92.7	92.4	92.2	91.8	91.1	88.9	86.0	79.3	90.0
	Recall	99.3	99.2	99.1	99.0	98.9	98.5	97.9	96.1	93.4	88.2	97.0
	F1	96.0	95.9	95.8	95.6	95.4	95.0	94.4	92.4	89.5	83.5	93.4
YOLOV7-X	Precision	92.0	91.9	91.8	91.7	91.6	91.3	90.3	88.6	85.8	77.3	89.2
	Recall	99.1	99.0	98.9	98.9	98.8	98.4	97.3	96.0	93.1	86.4	96.6
	F1	95.4	95.3	95.2	95.2	95.1	94.7	93.7	92.2	89.3	81.6	92.8
YOLOV8-X	Precision	93.1	93.1	93.0	92.8	92.4	92.2	91.4	89.5	86.6	79.2	90.3
	Recall	99.3	99.3	99.3	99.2	98.9	98.6	98.0	96.3	93.9	87.6	97.0
	F1	96.1	96.1	96.0	95.9	95.5	95.3	94.6	92.8	90.1	83.2	93.5
FasterR-CNN	Precision	89.1	89.1	88.9	88.7	88.5	88.2	87.8	86.5	81.5	66.6	85.5
	Recall	94.3	94.3	94.2	94.1	93.7	93.4	93.0	91.7	87.5	76.1	91.2
	F1	91.6	91.6	91.5	91.3	91.0	90.7	90.3	89.0	84.4	71.0	88.3
MaskR-CNN	Precision	90.2	90.0	90.0	89.8	89.7	89.1	88.2	86.3	81.7	64.4	85.9
	Recall	95.3	95.2	95.1	95.0	94.8	94.3	93.7	92.0	88.0	75.5	91.9
	F1	92.7	92.5	92.5	92.3	92.2	91.6	90.9	89.1	84.7	69.5	88.8
TableDet	Precision	91.7	91.7	91.7	91.5	91.3	90.7	90.2	88.4	84.0	72.9	88.4
	Recall	98.1	98.1	98.0	97.9	97.7	97.0	96.6	95.2	91.7	83.6	95.4
	F1	94.8	94.8	94.7	94.6	94.4	93.7	93.3	91.7	87.7	77.9	91.8
Deformable-DETR	Precision	90.3	90.3	90.1	90.1	89.8	89.2	88.5	86.8	84.4	77.2	87.7
	Recall	99.2	99.1	99.0	98.9	98.8	98.6	97.9	97.0	94.7	89.4	97.3
	F1	94.5	94.5	94.3	94.3	94.1	93.7	93.0	91.6	89.3	82.9	92.3
DiffusionDet	Precision	91.7	91.6	91.6	91.3	90.8	90.2	89.4	87.7	84.6	74.0	88.3
	Recall	99.6	99.6	99.6	99.4	98.7	97.9	97.2	95.5	92.9	85.2	96.6

Continued on next page

Table A.15 Detailed Experimental results on the TNCR dataset (continued from previous page).

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
SparseR-CNN	F1	95.5	95.4	95.4	95.2	94.6	93.9	93.1	91.4	88.5	79.2	92.3
	Precision	90.9	90.9	90.7	90.7	90.5	90.3	89.8	89.0	85.5	77.1	88.6
	Recall	99.9	99.8	99.8	99.7	99.7	99.7	99.5	98.7	97.1	89.9	98.4
	F1	95.2	95.1	95.0	95.0	94.9	94.8	94.4	93.6	90.9	83.0	93.2
SparseTableDet (Proposed)	Precision	93.0	93.0	92.9	92.9	92.7	92.5	92.2	91.3	88.8	77.2	90.6
	Recall	100.0	100.0	100.0	99.9	99.9	99.8	99.6	99.0	96.9	87.9	98.3
	F1	96.4	96.4	96.3	96.3	96.2	96.0	95.8	95.0	92.7	82.2	94.3

Table A.16: Detailed Experimental results on the ICT-TD dataset.

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
RetinaNet	Precision	95.8	95.7	94.7	94.4	92.4	91.1	90.0	87.9	82.1	68.0	89.2
	Recall	97.3	97.2	96.9	96.2	94.8	93.3	92.1	90.7	85.4	72.1	91.6
	F1	96.5	96.4	95.8	95.3	93.6	92.2	91.0	89.3	83.7	70.0	90.4
FCOS	Precision	92.0	91.8	90.8	90.5	89.4	88.1	86.8	84.4	80.5	66.8	86.1
	Recall	93.6	93.1	92.8	92.1	91.5	90.6	89.0	87.4	84.1	73.7	88.8
	F1	92.8	92.4	91.8	91.3	90.4	89.3	87.9	85.9	82.3	70.1	87.4
YOLOX-X	Precision	95.8	94.9	94.4	93.6	92.2	90.7	88.9	86.2	80.4	64.5	88.2
	Recall	98.7	97.8	97.3	96.7	95.1	93.5	91.4	88.3	82.9	67.9	91.0
	F1	97.2	96.3	95.8	95.1	93.6	92.1	90.1	87.2	81.6	66.2	89.6
YOLOR-X	Precision	97.6	97.3	96.4	95.5	95.0	94.2	93.3	90.7	88.3	79.1	92.7
	Recall	99.4	99.1	98.6	97.9	97.1	96.5	95.3	92.1	89.8	81.1	94.7
	F1	98.5	98.2	97.5	96.7	96.0	95.3	94.3	91.4	89.0	80.1	93.7
YOLOV5-X	Precision	97.4	97.2	97.2	97.0	96.4	95.6	94.8	93.7	90.7	81.2	94.1
	Recall	98.9	98.8	98.8	98.6	98.0	97.4	96.8	95.8	92.8	83.8	96.0
	F1	98.1	98.0	98.0	97.8	97.2	96.5	95.8	94.7	91.7	82.5	95.0
YOLOV7-X	Precision	98.2	98.2	98.0	97.9	96.8	95.8	94.8	93.7	91.8	80.9	94.6
	Recall	99.5	99.4	99.3	99.1	98.5	97.8	96.7	95.6	93.4	83.5	96.3

Continued on next page

Table A.16 Detailed Experimental results on the ICT-TD dataset (continued from previous page).

Method	Metric	IoU										
		50%	55%	60%	65%	70%	75%	80%	85%	90%	95%	50%:95%
YOLOV8-X	F1	98.8	98.8	98.6	98.5	97.6	96.8	95.7	94.6	92.6	82.2	95.4
	Precision	97.9	97.1	97.0	96.9	96.4	95.6	94.7	93.6	91.4	82.4	94.3
	Recall	99.1	98.8	98.7	98.6	98.1	97.5	96.6	95.7	93.2	84.9	96.1
	F1	98.5	97.9	97.8	97.7	97.2	96.5	95.6	94.6	92.3	83.6	95.2
FasterR-CNN	Precision	96.6	96.5	96.5	95.4	94.2	93.2	92.1	90.0	86.0	73.6	91.4
	Recall	97.4	97.2	97.1	96.4	95.3	94.8	93.7	91.7	87.6	76.5	92.8
	F1	97.0	96.8	96.8	95.9	94.7	94.0	92.9	90.8	86.8	75.0	92.1
MaskR-CNN	Precision	96.6	96.5	95.5	95.3	94.2	93.1	92.1	90.0	86.9	74.4	91.5
	Recall	97.4	97.1	96.9	96.3	95.5	94.4	93.4	91.7	88.9	77.8	92.9
	F1	97.0	96.8	96.2	95.8	94.8	93.7	92.7	90.8	87.9	76.1	92.2
TableDet [43]	Precision	97.4	96.4	96.3	96.3	95.1	94.0	92.9	90.5	88.2	72.5	92.0
	Recall	98.2	97.9	97.5	97.2	96.3	95.5	94.4	92.7	90.1	79.3	93.9
	F1	97.8	97.1	96.9	96.7	95.7	94.7	93.6	91.6	89.1	75.7	92.9
DiffusionDet [43]	Precision	96.5	96.4	96.3	95.8	95.2	94.5	93.9	92.5	89.2	73.6	92.4
	Recall	99.3	99.2	99.0	98.8	98.4	97.8	97.2	96.0	93.1	79.5	95.8
	F1	97.9	97.8	97.6	97.3	96.8	96.1	95.5	94.2	91.1	76.4	94.1
Deformable-DETR [43]	Precision	97.0	96.6	96.3	96.0	95.2	94.4	93.7	92.6	89.7	80.3	93.2
	Recall	98.9	98.6	98.5	98.3	97.8	96.9	96.3	95.2	92.8	85.8	95.9
	F1	97.9	97.6	97.4	97.1	96.5	95.6	95.0	93.9	91.2	83.0	94.5
SparseR-CNN [43]	Precision	96.2	96.0	95.5	95.3	94.2	93.3	92.6	91.1	88.3	75.6	91.8
	Recall	99.0	98.9	98.7	98.4	97.7	97.0	96.2	94.9	92.5	82.4	95.6
	F1	97.6	97.4	97.1	96.8	95.9	95.1	94.3	93.0	90.4	78.8	93.7
SparseTableDet (Proposed)	Precision	97.5	97.4	97.1	96.9	96.7	96.2	96.0	95.1	92.8	79.6	94.5
	Recall	99.3	99.3	99.3	99.2	99.2	98.7	98.5	97.8	95.6	84.1	97.1
	F1	98.4	98.3	98.2	98.0	97.9	97.4	97.2	96.4	94.2	81.8	95.8