# Towards End-to-End Semi-Supervised Table Detection with Semantic Aligned Matching Transformer

Tahira Shehzadi[*1,2,3][0000−0002−7052−979X], Shalini Sarode[1,3][0009−0007−9968−4068],
Didier Stricker[1,2,3], and Muhammad Zeshan Afzal[1,2,3][0000−0002−0536−6867]

[1] Department of Computer Science, Technical University of Kaiserslautern, 67663, Germany
[2] Mindgarage, Technical University of Kaiserslautern, 67663, Germany
[3] German Research Institute for Artificial Intelligence (DFKI), 67663 Kaiserslautern, Germany
{tahira.shehzadi@dfki.de}

**Abstract.** Table detection within document images is a crucial task in document processing, involving the identification and localization of tables. Recent strides in deep learning have substantially improved the accuracy of this task, but it still heavily relies on large labeled datasets for effective training. Several semi-supervised approaches have emerged to overcome this challenge, often employing CNN-based detectors with anchor proposals and post-processing techniques like non-maximal suppression (NMS). However, recent advancements in the field have shifted the focus towards transformer-based techniques, eliminating the need for NMS and emphasizing object queries and attention mechanisms. Previous research has focused on two key areas to improve transformer-based detectors: refining the quality of object queries and optimizing attention mechanisms. However, increasing object queries can introduce redundancy, while adjustments to the attention mechanism can increase complexity. To address these challenges, we introduce a semi-supervised approach employing SAM-DETR, a novel approach for precise alignment between object queries and target features. Our approach demonstrates remarkable reductions in false positives and substantial enhancements in table detection performance, particularly in complex documents characterized by diverse table structures. This work provides more efficient and accurate table detection in semi-supervised settings.

**Keywords:** Semi-Supervised Learning · Detection Transformer · SAM-DETR · Table Analysis · Table Detection.

## 1 Introduction

Document analysis has been the fundamental task in various workflow pipelines[1,2], with document summarization as its core task. The essential task in document analysis is identifying graphical objects like tables, figures, and text paragraphs. Previously, this task was carried out manually by analyzing the documents, understanding their contents, and summarizing them. However, the number of documents that need to be analyzed has drastically increased, and manual inspection is impossible. The growing number of documents led businesses to use more efficient and reliable automated methods. Optical character recognition(OCR) [3,4] and rule-based table detection approaches[5,6,7] are classical approaches for visual summarization. These methods perform well for documents with highly structured layouts because they are rule-based[5,6,7]. However, they struggle to adapt to varying and newer table designs, such as borderless tables. These limitations has shifted the research focus to developing techniques using deep learning [8,9,10,11]. These methods show significant improvements over traditional approaches [12], precisely detecting tables in documents irrespective of their structure. This advancement provides a notable improvement in document analysis and visual summarization.

Deep learning methods [13,14,15,16,17,18] eliminate handcrafted rules and excel at generalizing problems. However, their reliance on large amounts of labeled data for training

counteracts the aim of reducing manual work. Generating these labels is time-consuming and prone to errors [19]. Although these supervised deep learning approaches achieve state-of-the-art results on public benchmarks, their usage in industries is limited without similarly large annotated datasets in specific domains. Semi-supervised learning methods [20] have emerged as a solution to insufficient labeled data for deep learning applications. Recent advancements [21,22,23] utilize two detectors: one generates pseudo-labels for unlabeled data, and the other refines predictions using these pseudo-labels and a smaller set of labeled data. These detectors update each other throughout training [24,25,26,27]. However, it's important to note that the initial pseudo-label generator is often not robust, potentially leading to inaccurate labels and affecting overall performance.

Additionally, there are two major drawbacks in the earlier CNN-based semi-supervised methods[28,21,22]: First, they rely on anchor points for region proposals that require manual tuning. Second, they use post-processing techniques like Non-Maximal Suppression(NMS) to limit the number of overlapping predictions. The emergence of transformer-based methods [29,30,31,32] make the network end-to-end without NMS and anchor-free. This is possible due to their dependence on the attention mechanism and object queries. Consequently, there has been research mainly to improve the quality of object queries and improve the attention mechanism[33]. For example, Deformable DETR [30], AdaMixer [31] and REGO [34] focus on advancing the attention mechanism. Meanwhile, models like DN DETR [35], DAB DETR [36], and DINO DETR [29] are dedicated to improving the quality of object queries, and H-DETR [37], Co-DETR [32], and FANet [38] aim to increase the quantity of object queries. However, this increase leads to redundant predictions, adversely affecting performance. To counter this, a dual-stage object query approach has been proposed, combining one-to-one and one-to-many matching strategies. Despite its effectiveness, this method still impacts performance [37]. Addressing these challenges, we employ SAM-DETR [39], a novel model designed to optimize the matching process between object queries and corresponding target features in a semi-supervised setting. This approach effectively reduces false positives and improves table detection performance in complex documents.

In this paper, we introduce a novel semi-supervised approach for table detection, employing SAM-DETR [39] detector. Our main objective is to solve the non-robustness of the pseudo-label generation process. The training procedure consists of two modules: the teacher and the student. The teacher module consists of a pseudo-labeling framework, and the student uses these pseudo-labels along with a smaller set of labeled data to produce the final predictions. The pseudo-labeling process is optimized by iteratively refining the labels and the detector. The teacher module is updated by an Exponential Moving Average (EMA) from the student to improve the pseudo-label generation and detection modules. Our approach differs from conventional pseudo-labeling methods by incorporating a SAM-DETR detector without object proposal generation and post-processing steps like NMS. We enhance the ability to accurately match object queries with corresponding target features in complex documents, particularly excelling in the detection and handling of tables in semi-supervised settings. The intrinsic flexibility of this method enables consistent and reliable performance in various scenarios, including diverse table sizes and scales, within a semi-supervised learning context. Furthermore, this framework creates a reinforcing loop where the Teacher model consistently guides and improves the Student model. Our evaluation results demonstrate that our semi-supervised table detection approach achieves superior results compared to both CNN-based and other transformer-based semi-supervised methods without needing object proposals and post-processing steps such as NMS.

We summarize the primary contributions of this paper as follows:

- We introduce a novel semi-supervised approach for table detection. This approach eliminates the need for object proposals and post-processing techniques like Non-maximal Suppression (NMS).

- To the best of our knowledge, this is the first network that optimizes the matching process between object queries and corresponding target features in a semi-supervised setting.
- We conduct comprehensive evaluations on four diverse datasets: PubLayNet, ICDAR-19, TableBank, and Pubtables. Our approach achieves results comparable to CNN-based and transformer-based semi-supervised methods without requiring object proposal processes and Non-maximal Suppression (NMS) in post-processing.

## 2    Related Work

Analyzing document images involves the integral table detection task. This segment summarizes techniques for detecting tables, especially those involving complex structures. Initial methods relied on rules or metadata [40,41,42,43]. Meanwhile, more recent advances employ statistical and deep learning techniques [13,44,45,46], improving system adaptability and generalizability.

### 2.1    Table Detection Approaches

**Rule-based Approaches** Itonori et al. [40] laid the groundwork for table detection. The central focus was identifying tables as distinct text blocks using predefined rules. Building upon this, methods like [42] improved the approach by integrating various techniques, including table detection based on layout [47] or extracting tables from HTML-formatted documents [48]. Although effective for specific document types, these rule-based methods[5,6,7,49,50] lacked the flexibility to be universally applicable.

**Learning-based Approaches** Cesarini et al. [51] deviates from rule-based approaches by pioneering a supervised learning system for identifying table objects in document images. Their approach transforms a document image into an MXY tree model by classifying the blocks surrounded by vertical and horizontal lines as table objects. They further employed Hidden Markov Models [52,53] and an SVM classifier, along with conventional heuristics [54] for table detection. These techniques still needed additional data like ruling lines. In contrast, Deep Learning-based methods, further categorized as object detection, semantic segmentation, and bottom-up approaches, have demonstrated superior accuracy and efficiency over traditional techniques.

**Approaches Based on Semantic Segmentation.** Approaching table detection as a segmentation problem, methods like [55,56,57,58] generate pixel-level segmentation masks and then aggregate the masks to achieve final table detection. These methods utilize existing semantic segmentation networks and outperform traditional methods on various benchmark datasets [59,60,61,62,63,64,65]. Yang et al.'s [55] approach introduced a fully convolutional network (FCN) [66]. They used additional linguistic and visual features to enhance the segmentation results of page objects. He et al. [56] developed a multi-scale FCN that generates segmentation masks and their contours for table/text areas. They isolate the final table blocks after further refining the masks.

**Bottom-Up Methods.** These methods treat table detection as a graph-labeling task with graph nodes as page elements and edges as connections between them. Li et al. [67] used a conventional layout analysis to identify line areas. They then utilized two CNN-CRF networks to categorize these lines into four classes: text, figure, formula, and table. Later, they predicted a cluster for each pair of line areas. Holecek et al. [68] and Riba et al. [69] constructed a graph to establish the document layout and viewed text areas as nodes. They then used graph-neural networks for classifying nodes and edges. These methods require certain assumptions, like the necessity of text line boxes as additional input.

**Object Detection-Focused Techniques** Table detection in document images [70,71] is considered an object detection challenge, treating tables as natural objects. Hao et

al. [72] and Yi et al. [73] utilized R-CNN for table detection, but their performance still depended on heuristic rules, similar to earlier methods. Subsequently, more advanced single-stage object detectors like RetinaNet [74] and YOLO [75], as well as two-stage detectors like Fast R-CNN [8], Faster R-CNN [9], Mask R-CNN [76], and Cascade Mask R-CNN [77], were employed for detecting various document elements, including figures and formulas [78,79,80,81,82,83,13,84]. Additional enhancement techniques, such as image transformations involving coloration and dilation, were applied by [79,82,85]. Siddiqui et al. [86] integrate deformable convolution and RoI-Pooling [87] into Faster R-CNN for improved handling of geometrical changes. Agarwal et al. [83] combined a composite network [88] with deformable convolution to enhance the efficiency of the two-stage Cascade R-CNN. These CNN-based object detectors include heuristic stages like proposal generation and post-processing steps like non-maximal suppression (NMS). Our semi-supervised model treats detection as a set prediction task, eliminating the need for anchor generation and post-processing stages like NMS, resulting in a more streamlined and efficient detection process.

### 2.2   Semi-Supervised Learning in Object Detection

Semi-supervised object detection can be classified into consistency-based methods [89,90] and pseudo-label generation methods [21,22,23,91,92,93,94,95]. Our work focuses on the latter. Earlier works [21,22] employ diverse data augmentation techniques to generate pseudo-labels for unlabeled data. Meanwhile, [23] introduces SelectiveNet for pseudo-label generation by superimposing a bounding box from an unlabeled image onto a labeled image to ensure localization consistency within the labeled dataset. However, this approach involves a complex detection process due to image alteration. STAC [94] proposes to use strong augmentation for pseudo-label creation and weak augmentation for model training. Our method introduces a seamless end-to-end semi-supervised approach for table detection. Similar to other pseudo-label techniques [21,22,23,94,95], it incorporates a multi-level training strategy without the need for anchor generation and post-processing steps like Non-Maximal Suppression (NMS).

## 3   Methodology

First, the paper reviews SAM-DETR, a recent approach for detecting objects using transformers, in Section 3.1. Then, Section 3.2 describes our semi-supervised approach for learning with limited supervision and the generation of pseudo-labels for training.

### 3.1   Revisiting SAM-DETR

DEtection TRansformer (DETR) [96] introduces an encoder-decoder network for object detection. The encoder network extracts features from the image to focus on key details. The decoder then processes these features with object queries, using self-attention and cross-attention mechanisms to identify and locate objects. However, DETR's initial non-selective approach in processing images and object queries can lead to slower detection, especially in semi-supervised learning with limited data. By refining the attention mechanism and enhancing the quality and quantity of object queries, researchers aim to boost DETR's efficiency, accuracy, and training speed [33]. SAM-DETR, as shown in Fig. 1 stands out for its innovative addition of a semantics aligner module and learnable reference boxes within the Transformer decoder part of DETR. Overall, SAM-DETR's enhancements to the original DETR model focus on making the object detection process more efficient in terms of accuracy and speed.
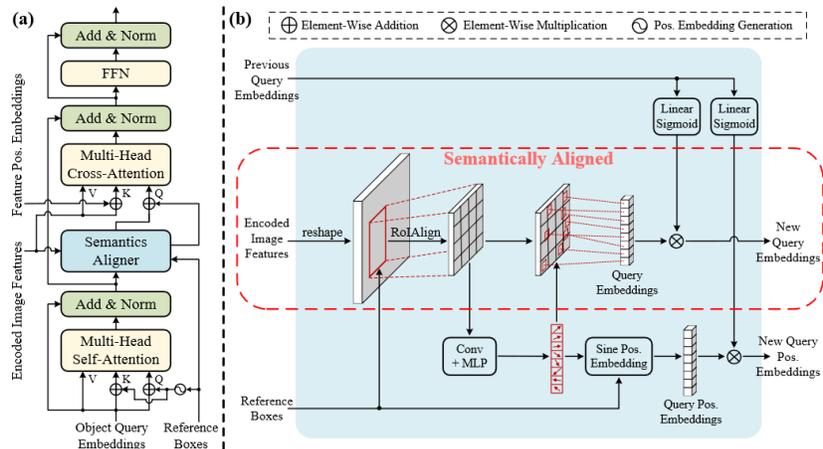
Fig. 1: Overview of SAM-DETR [39]. (a) the architecture of a single decoder layer in SAM-DETR, showing the role of learnable reference boxes in generating position embeddings for each object query. (b) the pipeline of the Semantics Aligner. The process includes the use of reference boxes for feature extraction via RoIAlign, the prediction of salient points in the targeted region, and the generation of new, semantically aligned query embeddings, which are further refined by incorporating attributes from previous queries. Image from [39].

**Semantics Aligner.** Semantic-Aligned Matching focuses on improving the interaction between object queries and encoded image features. Generally, the cross-attention module uses a dot-product method, which is effective in identifying similarities between two vectors. This method typically guides object queries to focus on regions of the image that are more similar. However, the original DETR model does not ensure that object queries and encoded image features are in the same embedding space, leading to less effective matching and requiring extensive training time. To address this, the Semantic-Aligned Matching approach introduces a mechanism to align object queries with encoded image features semantically. This alignment ensures that both are in the same embedding space, making the dot-product a more meaningful measure of similarity. As a result, object queries are more likely to focus on semantically similar regions, enhancing the efficiency and effectiveness of the object detection process.

**Multi-Head Attention and Salient Points.** In DETR, multi-head attention is crucial for focusing on different image parts, enhancing scene understanding. SAM-DETR builds on this by identifying key points on objects, using ConvNet and MLP to predict these points for better alignment and detection. Features from these points are integrated with multi-head attention, allowing each head to concentrate on specific, significant object features, improving accuracy and localization.

**Reweighted Queries.** The Semantics Aligner in DETR aligns object queries with encoded image features but initially misses crucial information from previous embeddings. To address this, it uses a linear projection and sigmoid function to create reweighting coefficients, applied to both new and positional query embeddings. This ensures important features are emphasized and previous data is utilized, significantly enhancing detection.

### 3.2   Semi-Supervised SAM-DETR

We propose a semi-supervised learning approach that improves object detection through semantic alignment and utilizes limited labeled data for training, as shown in Fig. 2. The model leverages fully labeled and unlabeled data for object detection tasks in the semi-
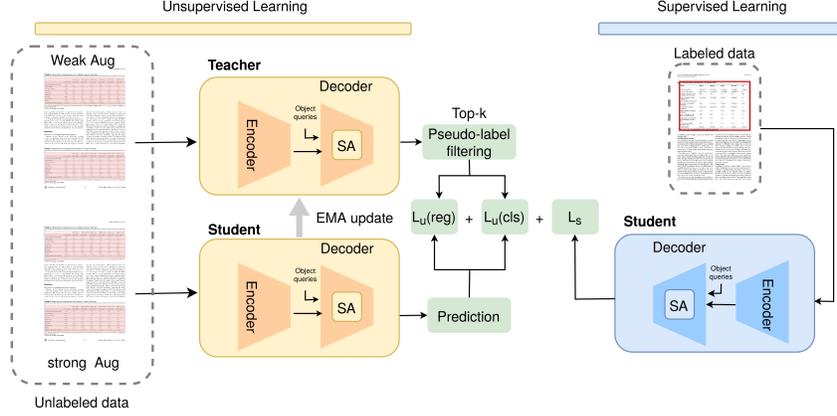
Fig. 2: Illustration of our Semi-Supervised Table Detection Framework. This dual-component system involves a Student module that learns from a mix of labeled data and strongly augmented unlabeled images, and a Teacher module that refines its understanding using weakly augmented unlabeled images. The Student module updates the Teacher module using Exponential Moving-Average (EMA) during training. Within this setup, the Semantics Aligner (SA) is key in the decoder of the student-teacher framework, fine-tuning the relationship between object queries and the image features that have been encoded, ensuring a more effective and accurate detection of tables in various documents.

supervised setting. It consists of two key modules: the student and teacher modules. The student module processes both labeled and unlabeled images. Strong augmentation is applied to unlabeled data, while strong and weak augmentations are applied to labeled data. The teacher module operates on unlabeled images with weak augmentations. It plays a crucial role in generating pseudo-labels for unlabeled data. These pseudo-labels are then employed for supervised training by the student module. Weak augmentation is applied to the unlabeled data for the teacher module to produce more accurate pseudo-labels. In contrast, the student module, designed for more challenging learning, utilizes strong augmentation for unlabeled data. At the start of training, the teacher and student models are randomly initialized. As training progresses, the teacher model is continuously updated by the student model using an exponential moving average (EMA) strategy. For the student module, the student's queries $Q_s$ and features $F_s$ are fed into the decoder. Similarly, in the teacher module, the teacher's queries $Q_t$ and features $F_t$ go through a similar process with the teacher's decoder as follows:

$$\hat{o}_s = \text{Decoder}_s\left(Q_s, F_s\right) \tag{1}$$

$$\hat{o}_t = \text{Decoder}_t\left(Q_t, F_t\right) \tag{2}$$

In the decoder, the Semantics Aligner processes the encoded image features for students $F_s$ and teachers $F_t$, both initially in 1D sequences of dimensions $HW \times d$. The Aligner converts these features into 2D maps with dimensions $H \times W \times d$, using the reference boxes of object queries, denoted as $R_s^{box}$ for the student and $R_t^{box}$ for the teacher. After this transformation, the aligner employs RoIAlign to extract region-level features, represented as $F_s^R$ for the student and $F_t^R$ for the teacher, from the encoded image features. The final step involves generating new object queries, $Q_{\text{new}}$ and their position embeddings $Q_{\text{new pos}}$, through resampling based on $F_s^R$ and $F_t^R$ as follows.

$$F_s^R = \text{RoIAlign}(F_s, R_s^{\text{box}}), \quad F_t^R = \text{RoIAlign}(F_t, R_t^{\text{box}}) \tag{3}$$

$$Q_s^{\text{new}}, Q_{\text{s,pos}}^{\text{new}} = \text{Resample}(F_s^R, R_s^{\text{box}}, Q_s), \tag{4}$$

$$Q_t^{\text{new}}, Q_{\text{t,pos}}^{\text{new}} = \text{Resample}(F_t^R, R_t^{\text{box}}, Q_t) \tag{5}$$

Next, we extract features via a ConvNet and MLP to identify salient points within these regions. These points are then used to create new object query embeddings $Q_s^{\text{new}}$ and $Q_t^{\text{new}}$, ensuring they stay within reference boxes for accuracy. Finally, position embeddings $Q_{\text{s,pos}}^{\text{new}}$ and $Q_{\text{t,pos}}^{\text{new}}$ derived from these points are concatenated, feeding into a multi-head cross-attention module for further processing.

$$R_s^{sp} = MLP(ConvNet(F_s^R)) \tag{6}$$

$$Q_s^{\text{new}} = \text{Concat}\left(\{F_s^R[\ldots, x, y, \ldots] \text{ for } x, y \in R_s^{\text{sp}}\}\right) \tag{7}$$

$$Q_{\text{s,pos}}^{\text{new}} = \text{Concat}(\text{Sin}(R_s^{\text{box}}, R_s^{\text{sp}})) \tag{8}$$

$$R_t^{sp} = MLP(ConvNet(F_t^R)) \tag{9}$$

$$Q_t^{\text{new}} = \text{Concat}\left(\{F_t^R[\ldots, x, y, \ldots] \text{ for } x, y \in R_t^{\text{sp}}\}\right) \tag{10}$$

$$Q_{\text{t,pos}}^{\text{new}} = \text{Concat}(\text{Sin}(R_t^{\text{box}}, R_t^{\text{sp}})) \tag{11}$$

The semantics aligner generates new object queries aligned with image features and incorporates previous query embeddings by generating reweighting coefficients. These coefficients, created through linear projection and sigmoid functions, are applied to new and old query embeddings to emphasize key features. This approach ensures that the valuable information from previous queries is effectively utilized.

$$Q_s^{\text{new}} = Q_s^{\text{new}} \otimes \sigma(Q_s W_s^{\text{RWs1}}), \qquad Q_t^{\text{new}} = Q_t^{\text{new}} \otimes \sigma(Q_t W_t^{\text{RWt1}}) \tag{12}$$

$$Q_{\text{s,pos}}^{\text{new}} = Q_{\text{s,pos}}^{\text{new}} \otimes \sigma(Q_s W_s^{\text{RWs2}}), \qquad Q_{\text{t,pos}}^{\text{new}} = Q_{\text{t,pos}}^{\text{new}} \otimes \sigma(Q_t W_t^{\text{RWt2}}) \tag{13}$$

Here, $W_{\text{RWt1}}$ and $W_{\text{RWt2}}$ are used to denote linear projection functions. The symbol $\sigma(\cdot)$ refers to the sigmoid function, while $\otimes$ represents the operation of element-wise multiplication. The subscripts t and s refer to the teacher and student module, respectively. Combining the semantic alignment capabilities with the semi-supervised approach allows the model to effectively utilize labeled and unlabeled data, leading to improved object detection performance. This approach is particularly useful when labeled data is limited, as it maximizes the information extracted from available resources.

## 4 Pseudo-Label Filtering Framework

In our semi-supervised learning framework, we employ the Top-K pseudo-label filtering technique to augment the training process of our machine learning models, especially when the labeled data is limited. This approach is instrumental in making the most of the unlabeled data. Here, the key strategy is pseudo-labeling, where our model generates labels for the unlabeled data based on its current level of understanding. However, diverging from the traditional method of relying on the single most confident prediction, our top-k approach considers each data point's top 'k' predictions. For instance, if 'k' is set at 3, the model evaluates and includes the three highest probable labels for each piece of unlabeled data in the training process. The benefits of our top-k strategy are significant. Firstly, it broadens the model's exposure to more challenging 'hard samples' data points that are typically difficult to classify and might be overlooked by standard top-1 pseudo-labeling methods. Including a wider range of examples substantially improves the model's learning. Secondly, our approach is effective in cases involving objects or data points with similar features. By

acknowledging and incorporating ambiguity through multiple potential labels, the model is better equipped to handle complex classification scenarios where clear-cut distinctions between categories are not always evident. Implementing the top-k pseudo-label filtering in our semi-supervised learning setting is a pivotal step towards enhancing the model's accuracy and robustness, ensuring a more comprehensive and enhanced learning process. The teacher model generates pseudo boxes for unlabeled images, and the student model is trained on labeled images with ground-truth annotations and unlabeled images with pseudo boxes treated as ground-truth. Therefore, the overall loss is defined as the weighted sum of supervised and unsupervised losses:

$$L = L_s + \alpha L_u, \tag{14}$$

Where $L_s$ represents the supervised loss for labeled images, $L_u$ represents the unsupervised loss for unlabeled images, and $\alpha$ with value 0.25 controls the contribution of the unsupervised loss. Both losses are normalized by the respective number of images in the training data batch:

$$L_s = \frac{1}{N_l} \sum_{i=1}^{N_l} (L_{cls}(I_i^l) + L_{reg}(I_i^l)), \tag{15}$$

$$L_u = \frac{1}{N_u} \sum_{i=1}^{N_u} (L_{cls}(I_i^u) + L_{reg}(I_i^u)), \tag{16}$$

Where $I_i^l$ indicates the $i$-th labeled image, $I_i^u$ indicates the $i$-th unlabeled image, $L_{cls}$ is the classification loss, $L_{reg}$ is the box regression loss, $N_l$ is the number of labeled images, and $N_u$ is the number of unlabeled images. Overall, our semi-supervised learning setting enhances the model's accuracy and robustness, ensuring a more comprehensive learning process.

## 5    Experimental Setup

### 5.1    Datasets

**TableBank:** TableBank [64], a prominent dataset in the field of document analysis, ranks as the second-largest collection for table recognition tasks. This dataset comprises 417,000 document images, annotated via a process of crawling the arXiv database. It categorizes tables into three types: LaTeX images (253,817), Word images (163,417), and a combined set (417,234). Furthermore, TableBank provides data for table structure recognition. In our study, we utilizeonly the table detection component of the TableBank dataset.

**PubLayNet:** PubLayNet [60], a sizable dataset in the public domain, encompasses 335,703 images for training, 11,240 for validation, and 11,405 for testing. It features annotations like polygonal segmentation and bounding boxes for figures, lists, titles, tables, and texts in images sourced from research papers and articles. The dataset's evaluation employed the COCO analytics method [97]. We selectively used 102,514 images from the 86,460 table annotations in PubLayNet for our experiments.

**PubTables:** PubTables-1M [65], specifically tailored for table detection in scientific documents, is an extensive dataset featuring nearly one million tables. It stands out for its comprehensive annotations, including precise location information, crucial for accurately detecting tables within diverse documents. Its large scale and meticulous annotations make it a significant resource for developing and refining table detection algorithms.

**ICDAR-19:** The ICDAR 2019 competition for Table Detection and Recognition (cTDaR) [59] introduced two novel datasets (modern and historical) for the table detection task (TRACK A). To facilitate direct comparisons with previous methods [82], we provide results at an Intersection over Union (IoU) threshold of 0.8 and 0.9.

## 5.2 Evaluation Criteria

We assess the effectiveness of our semi-supervised table detection method through specific evaluation metrics: Precision, Recall, and F1-score. Precision [98] is the ratio of correctly predicted positive observations (True Positives) to the total predicted positive observations (True Positives + False Positives). Recall [98] measures the proportion of actual positives correctly identified (True Positives) out of the total actual positives (True Positives + False Negatives). The F1-score [98] is the harmonic mean of Precision and Recall. Moreover, We evaluate our approach using AP@50 and AP@75, which assess precision at 50% and 75% IoU thresholds, reflecting moderate and high localization accuracy respectively, alongside average recall, measuring our model's capacity to detect all relevant instances

## 5.3 Implementation Details

We use the ResNet-50 backbone on 8 Nvidia RTXA6000 GPUs, initially trained on the ImageNet dataset, to evaluate the effectiveness of our semi-supervised method. We train on a diverse range of datasets, including PubLayNet, ICDAR-19, PubTables, and all subsets of the TableBank dataset, taking randomly 10%, 30%, and 50% labeled data with the remaining as unlabeled. We conduct pseudo-labeling with a 0.7 threshold and optimize using AdamW. Our training spans 120 epochs, reducing the learning rate by 10% after the 110th epoch, and we typically set our batch size to 16. We adopt DETR's data augmentation strategy, which involves horizontal flipping, random cropping, and resizing. Additionally, we apply strong augmentation techniques such as horizontal flips, resizing, patch removal, cropping, conversion to grayscale, and Gaussian blur. For weak augmentation, we focus mainly on horizontal flipping. Setting the number of queries (N) in the decoder to 30 gives the best results. Our resizing approach ensures the image's longest side is at most 1333 pixels and the shortest side is at least 480 pixels. These strategic adjustments and augmentations boost the model's performance and efficiency.

Table 1: Performance of our semi-supervised transformer-based approach on different splits of TableBank dataset with varying percentage label data.

| Dataset | Labels | mAP | $AP^{50}$ | $AP^{75}$ | $AR_L$ |
|---|---|---|---|---|---|
| TableBank-word | 10% | 92.9 | 95.3 | 93.9 | 97.4 |
| | 30% | 94.1 | 95.8 | 94.5 | 98.2 |
| | 50% | 94.3 | 95.8 | 94.8 | 98.3 |
| TableBank-latex | 10% | 91.2 | 97.6 | 96.4 | 95.3 |
| | 30% | 93.7 | 97.3 | 96.3 | 97.7 |
| | 50% | 94.8 | 97.9 | 97.0 | 98.1 |
| TableBank-both | 10% | 92.7 | 95.8 | 94.6 | 93.6 |
| | 30% | 93.8 | 95.2 | 95.2 | 93.6 |
| | 50% | 94.2 | 96.1 | 95.8 | 95.8 |

Table 2: Recall results comparison of our semi-supervised approach with previous semi-supervised table detection approach. Here Def-semi refers to [99].

| Dataset | Labels | Def-semi | Our |
|---|---|---|---|
| TableBank-word | 10% | 87.1 | 97.4 |
| | 30% | 92.1 | 98.2 |
| | 50% | 94.5 | 98.3 |
| TableBank-latex | 10% | 74.3 | 95.3 |
| | 30% | 89.0 | 97.7 |
| | 50% | 91.4 | 98.1 |
| TableBank-both | 10% | 90.1 | 93.6 |
| | 30% | 91.5 | 93.6 |
| | 50% | 95.3 | 95.8 |

# 6 Results and Discussion

## 6.1 TableBank

In our study, we evaluate our approach using the TableBank dataset, examining performance across various splits with different proportions of labeled data: 10%, 30%, and 50%.

Table 1 shows we achieve mAP of 92.9%, 91.2%, and 92.7% by using 10% labels of Table-Bank word, latex, and both splits, respectively. Unlike previous semi-supervised table detection method [99], which employs deformable DETR [30] with a focus on improving the attention mechanism to improve the performance. Our semi-supervised approach optimizes the matching process between object queries and image features. As a result, our semi-supervised strategy achieves significantly higher recall rates than earlier semi-supervised methods, as shown in Tables 2. This improvement shows the effectiveness of semi-supervised table detection, particularly when dealing with limited labeled data. Table 3 presents a

Table 3: Comparative analysis of our semi-supervised approach with previous supervised and semi-supervised methods on the TableBank-Both dataset using 10%, 30%, and 50% labeled data. Here, the results are reported on mAP.

| Method | Approach | Detector | 10% | 30% | 50% |
|---|---|---|---|---|---|
| Ren et al. [9] | supervised | Faster R-CNN | 80.1 | 80.6 | 83.3 |
| Zhu et al. [30] | supervised | Deformable DETR | 80.8 | 82.6 | 86.9 |
| STAC [94] | semi-supervised | Faster R-CNN | 82.4 | 83.8 | 87.1 |
| Unbiased Teacher [100] | semi-supervised | Faster R-CNN | 83.9 | 86.4 | 88.5 |
| Humble Teacher [101] | semi-supervised | Faster R-CNN | 83.4 | 86.2 | 87.9 |
| Soft Teacher [28] | semi-supervised | Faster R-CNN | 83.6 | 86.8 | 89.6 |
| Shehzadi et al. [99] | semi-supervised | Deformable DETR | 84.2 | 86.8 | 91.8 |
| Our | semi-supervised | Sam-DETR | **92.7** | **93.8** | **94.2** |

comparative analysis of our semi-supervised approach against prior supervised and semi-supervised methods using the TableBank-both dataset, which includes splits with 10%, 30%, and 50% labeled data. The outcomes demonstrate that our approach outperforms the earlier methods across these varying levels of labeled data. This is a significant finding, highlighting the effectiveness of our semi-supervised strategy in scenarios with limited labeled data availability.

## 6.2   PubLayNet

We also evaluate the performance of our transformer-based semi-supervised learning model on the PubLayNet dataset, experimenting with different ratios of labeled to unlabeled data (10%, 30%, and 50%). This study aims at understanding the model's performance in scenarios with limited labeled data, a common challenge in real-world applications. Table 4 shows we achieve mAP of 89.9%, 90.9%, and 93.2% by using 10%, 30%, and 50% labels of PubLayNet dataset. We shows the visual analysis of our semi-supervised approach in Fig. 3. Our semi-supervised approach also provides higher recall than the previous semi-supervised approach, as observed in Table 5.

Table 4: Performance of our semi-supervised transformer-based approach on PubLayNet dataset with varying percentage label data.

| Dataset | Label-percent | mAP | $AP^{50}$ | $AP^{75}$ | $AR_L$ |
|---|---|---|---|---|---|
| PubLayNet | 10% | 89.9 | 97.1 | 94.3 | 96.6 |
| | 30% | 90.9 | 97.4 | 94.9 | 96.9 |
| | 50% | 93.2 | 97.7 | 95.0 | 97.3 |

Table 5: Recall results comparison of our approach with previous semi-supervised table detection approach.

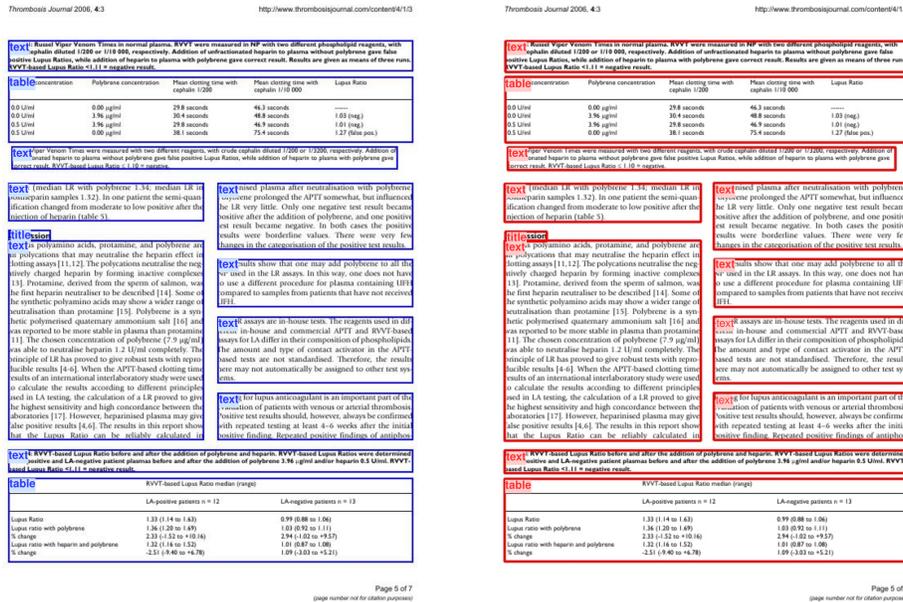| Method | 10% | 30% | 50% |
|---|---|---|---|
| Shehzadi et al. [99] | 91.0 | 93.2 | 96.0 |
| Our | **96.6** | **96.9** | **97.3** |

Fig. 3: Visual Analysis of our semi-supervised approach. Here, blue represents ground truth and red denotes our predictions results using 10% labels on PubLayNet datatset.

We also compare our approach against traditional deep learning methods, both supervised and semi-supervised, to highlight advancements. A key focus is the model's performance with only 10% labeled data, where we observe that our approach achieves the highest mAP score of 89.9, as detailed in Table 6. This shows the effectiveness of our method in leveraging minimal labeled data, demonstrating the significant potential of our approach for practical applications in table detection and recognition.

Table 6: Comparative analysis of our semi-supervised approach with previous supervised and semi-supervised methods on PubLayNet table class dataset using 10%, 30%, and 50% labeled data. Here, the results are reported on mAP.

| Method | Approach | Detector | 10% | 30% | 50% |
|---|---|---|---|---|---|
| Ren et al. [9] | supervised | Faster R-CNN | 83.4 | 86.6 | 87.9 |
| Zhu et al. [30] | supervised | Deformable DETR | 83.9 | 86.8 | 88.1 |
| Soft Teacher [28] | semi-supervised | Faster R-CNN | 88.3 | 89.5 | 92.5 |
| Shehzadi et al. [99] | semi-supervised | Deformable DETR | 88.4 | 90.3 | 92.8 |
| Our | semi-supervised | SAM-DETR | **89.9** | **90.9** | **93.2** |

### 6.3 PubTables

In this subsection, we detail our experimental results for the PubTables dataset in a semi-supervised setting using different percentages of labeled data. Our analysis includes a comparison between our transformer-based semi-supervised method and earlier CNN-based and transformer-based supervised approaches. As shown in Table 7, our semi-supervised

approach achieves a 92.3 mAP score even with only 10% of the data labeled, which highlights the effectiveness of our method in utilizing a smaller amount of labeled data to attain high accuracy.

Table 7: Performance of our semi-supervised transformer-based approach on the PubTables dataset with varying levels of labeled data (10%, 30%, 50%). Results show high accuracy with even a minimal amount of labeled data.

| Dataset | Label | mAP | $AP^{50}$ | $AP^{75}$ | $AR_L$ |
|---------|-------|-----|-----------|-----------|--------|
| PubTables | 10% | 92.3 | 93.7 | 93.8 | 87.8 |
| | 30% | 93.5 | 94.8 | 93.7 | 88.1 |
| | 50% | 93.8 | 94.8 | 94.8 | 88.3 |

Table 8 presents a comparison between our semi-supervised approach and previous supervised methods. While a direct comparison isn't feasible due to different percentages of label data for training, our results are notably comparable. For instance, a Faster R-CNN model trained on fully labeled data achieved an mAP of 82.5, whereas our semi-supervised approach reached an mAP of 92.3 using only 10% labeled data.

Table 8: Comparative Analysis of Semi-Supervised and Supervised Methods. It clearly shows that our semi-supervised model achieves comparable results even with limited data.

| Method | Approach | Detector | mAP | $AP^{50}$ | $AP^{75}$ |
|--------|----------|----------|-----|-----------|-----------|
| Smock et al. [65] | supervised | Faster R-CNN | 82.5 | 98.5 | 92.7 |
| Smock et al. [65] | supervised | DETR | 96.6 | 995 | 98.8 |
| Our | semi-supervised (10%) | SAM-DETR | 92.3 | 93.7 | 93.8 |

**Comparisons with Previous Table Detection Approaches.** In Table 9, we present a comprehensive comparison of our semi-supervised table detection approach against existing supervised and semi-supervised methods. Our approach facilitates learning with signif-

Table 9: Comparative analysis of our semi-supervised approach with previous supervised and semi-supervised methods. Here, the results are reported on mAP.

| Method | Approach | Labels | TableBank | PubLayNet | PubTables |
|--------|----------|--------|-----------|-----------|-----------|
| CDeC-Net [83] | supervised | 100% | 96.5 | 97.8 | - |
| CasTabDetectoRS [45] | supervised | 100% | 95.3 | - | - |
| Faster R-CNN [60] | supervised | 100% | - | 90 | |
| VSR [102] | supervised | 100% | - | 95.69 | |
| Smock et al. [65] | supervised | 100% | - | - | 96.6 |
| Shehzadi et al. [99] | semi-supervised | 10% | 84.2 | 88.4 | - |
| Our | semi-supervised | 10% | 92.7 | 89.9 | 92.3 |

icantly fewer labeled instances. Our semi-supervised method performs well despite limited labeled data, achieving high mAP scores on datasets and outperforming previous semi-supervised models. It shows improved performance in scenarios with scarce labeled data, offering comparable results to fully supervised methods while using only 10% of their labeled data.

## 6.4   ICDAR-19

In our analysis, we additionally conduct an evaluation of the ICDAR-19 TrackA table detection dataset across different Intersection over Union (IoU) thresholds using 50% labeled data. Furthermore, we compare our semi-supervised approach with earlier supervised and semi-supervised strategies, as depicted in Table 10. The results, utilizing 50% labeled data, show that our transformer-based semi-supervised framework surpasses prior semi-supervised methods, demonstrating superior accuracy.

Table 10: Performance comparison between the proposed semi-supervised approach and previous state-of-the-art results on the dataset of ICDAR 19 Track A (Modern).

| Method | Approach | IoU=0.8 | | | IoU=0.9 | | |
|---|---|---|---|---|---|---|---|
| | | Recall | Precision | F1-Score | Recall | Precision | F1-Score |
| TableRadar [59] | supervised | 94.0 | 95.0 | 94.5 | 89.0 | 90.0 | 89.5 |
| NLPR-PAL [59] | supervised | 93.0 | 93.0 | 93.0 | 86.0 | 86.0 | 86.0 |
| Lenovo Ocean [59] | supervised | 86.0 | 88.0 | 87.0 | 81.0 | 82.0 | 81.5 |
| CDeC-Net [83] | supervised | 93.4 | 95.3 | 94.4 | 90.4 | 92.2 | 91.3 |
| HybridTabNet [46] | supervised | 93.3 | 92.0 | 92.8 | 90.5 | 89.5 | 90.2 |
| Shehzadi et al. [99] | semi-supervised (50%) | 71.1 | 82.3 | 76.3 | 66.3 | 76.8 | 71.2 |
| Our | semi-supervised (50%) | 73.5 | 83.8 | 77.2 | 68.4 | 77.8 | 72.1 |

## 7   Ablation Study

In the ablation study, we evaluate the model's performance using only 30% of the labeled data from the PubTables dataset. The study observes the effect of varying the pseudo-labeling confidence threshold, the number of filtered pseudo-labels, and the number of learnable queries, offering insights into their roles in enhancing model performance in document analysis tasks.

**Pseudo-Labeling confidence threshold** The choice of a confidence threshold in pseudo-labeling influences the performance of our semi-supervised approach, as observed in Table 11. A low threshold leads to the filtering of a large number of pseudo-labels. However, these include incorrect pseudo-labels, introducing noise into the training process, and potentially degrading the model's performance. On the other hand, a high threshold ensures the generation of high-quality pseudo-labels, reducing the risk of noise. However, this results in fewer pseudo-labels fed into the student network, thus not fully leveraging the advantages of semi-supervised learning. The balance between generating enough pseudo-labels and ensuring that these pseud-labels are accurate enough to be useful is crucial in optimizing model performance.

**Influence of Learnable queries Quantity** We examine the effect of both increasing and decreasing the number of input queries on the performance of our semi-supervised approach, as highlighted in Table 12. While increasing the queries can improve the model's ability to

Table 11: Performance comparison using different Pseudo-labeling confidence threshold values. The best threshold values are shown in bold.

| Threshold | AP | $AP^{50}$ | $AP^{75}$ |
|---|---|---|---|
| 0.5 | 89.8 | 91.3 | 90.4 |
| 0.6 | 90.4 | 92.1 | 91.5 |
| **0.7** | **93.5** | **94.8** | **93.7** |
| 0.8 | 90.2 | 91.7 | 90.2 |
| 0.9 | 88.6 | 89.3 | 89.1 |

Table 12: Performance comparison using different numbers of learnable queries to the decoder input. Here, the best performance results are shown in bold.

| Queries | AP | $AP^{50}$ | $AP^{75}$ |
|---|---|---|---|
| 10 | 88.5 | 87.8 | 86.8 |
| **30** | **93.5** | **94.8** | **93.7** |
| 60 | 91.8 | 92.8 | 91.5 |
| 100 | 88.6 | 90.2 | 87.3 |
| 300 | 82.1 | 85.3 | 84.1 |

detect and focus on a wide range of features, enhancing accuracy in complex detection tasks, it also leads to more overlapping predictions, necessitating the use of Non-Maximum Suppression (NMS). Conversely, decreasing the number of queries reduces computational complexity but limits the model's detection capabilities. We find that our model achieves the best performance with 30 queries. Deviating from this optimal count, whether by increasing or decreasing the number of queries, significantly impacts the model's accuracy and efficiency.

Table 13: Performance evaluation using top-k pseudo-labels. The best results are in bold.

| Top-k | AP | $AP^{50}$ | $AP^{75}$ |
|---|---|---|---|
| 1 | 90.5 | 93.8 | 91.2 |
| 2 | 91.7 | 94.4 | 91.9 |
| **3** | **93.5** | **94.8** | **93.7** |
| 4 | 92.8 | 94.2 | 92.5 |

**Influence of quantity of Pseudo-label Filtering** In Table 13, we observe the impact of varying quantities of filtered pseudo-labels generated by the teacher network on model performance. While including more pseudo-labels enhances model performance, it is also vital to consider their quality. Selecting more pseudo-labels, such as the top-4, inherently introduces some lower-quality labels into the training process. Including less reliable pseudo-labels can adversely affect the model's performance, highlighting the need for a balanced approach in pseudo-label selection that optimizes quantity and quality to achieve the best model performance.

## 8 Conclusion

Our research addresses the challenge of accurately and efficiently detecting document objects, such as tables and text, in semi-supervised settings. This approach utilizes minimal labeled data and employs student-teacher networks that mutually update during training. Previous transformer-based research focused on improving attention or increasing the number of object queries, which impacts training time and performance. We eliminate the need for NMS and focus on matching between object queries and image features. Our novel approach using SAM-DETR in a semi-supervised setting helps align object queries with target features, significantly reducing false positives and improving the detection of document objects in complex layouts. In short, our semi-supervised method enhances the accuracy of document analysis, particularly in scenarios with limited labeled data.

# References

1. T. M. Breuel and K. Tombre, *Document Analysis Systems: Theory and Practice*. World Scientific Publishing, 2017.
2. R. Kasturi, L. O'Gorman, and V. Govindaraju, "Document image analysis: A primer," *Sadhana - Academy Proceedings in Engineering Sciences*, vol. 27, pp. 3–22, 02 2002.
3. Z. Zhao, M. Jiang, S. Guo, Z. Wang, F. Chao, and K. C. Tan, "Improving deep learning based optical character recognition via neural architecture search," in *2020 IEEE Congress on Evolutionary Computation (CEC)*, 2020, pp. 1–7.
4. D. Van Strien, K. Beelen, M. C. Ardanuy, K. Hosseini, B. McGillivray, and G. Colavizza, "Assessing the impact of ocr quality on downstream nlp tasks," 2020.
5. B. Coüasnon and A. Lemaitre, "Recognition of tables and forms," in *Handbook of Document Image Processing and Recognition*, 2014.
6. R. Zanibbi, D. Blostein, and J. R. Cordy, "A survey of table recognition," *Document Analysis and Recognition*, vol. 7, no. 1, pp. 1–16, 2004.
7. A. M. Jorge, L. Torgo *et al.*, "Design of an end-to-end method to extract information from tables," *International Journal of Document Analysis and Recognition (IJDAR)*, vol. 8, no. 2, pp. 144–171, 2006.
8. R. B. Girshick, "Fast R-CNN," *CoRR*, vol. abs/1504.08083, 2015. [Online]. Available: http://arxiv.org/abs/1504.08083
9. S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," *CoRR*, vol. abs/1506.01497, 2015. [Online]. Available: http://arxiv.org/abs/1506.01497
10. J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," *CoRR*, vol. abs/1612.08242, 2016. [Online]. Available: http://arxiv.org/abs/1612.08242
11. K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2980–2988.
12. T. Orosz, R. Vági, G. M. Csányi, D. Nagy, I. Üveges, J. P. Vadász, and A. Megyeri, "Evaluating human versus machine learning performance in a legaltech problem," *Applied Sciences*, vol. 12, no. 1, 2022. [Online]. Available: https://www.mdpi.com/2076-3417/12/1/297
13. S. Schreiber, S. Agne, I. Wolf, A. Dengel, and 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. 01, 2017, pp. 1162–1167.
14. M. Minouei, K. A. Hashmi, M. R. Soheili, M. Z. Afzal, and D. Stricker, "Continual learning for table detection in document images," *Applied Sciences*, vol. 12, no. 18, 2022. [Online]. Available: https://www.mdpi.com/2076-3417/12/18/8969
15. K. A. Hashmi, D. Stricker, M. Liwicki, M. N. Afzal, and M. Z. Afzal, "Guided table structure recognition through anchor optimization," *CoRR*, vol. abs/2104.10538, 2021. [Online]. Available: https://arxiv.org/abs/2104.10538
16. K. A. Hashmi, A. Pagani, M. Liwicki, D. Stricker, and M. Z. Afzal, "Cascade network with deformable composite backbone for formula detection in scanned document images," *Applied Sciences*, vol. 11, no. 16, 2021. [Online]. Available: https://www.mdpi.com/2076-3417/11/16/7610
17. S. Sinha, K. A. Hashmi, A. Pagani, M. Liwicki, D. Stricker, and M. Z. Afzal, "Rethinking learnable proposals for graphical object detection in scanned document images," *Applied Sciences*, vol. 12, no. 20, 2022. [Online]. Available: https://www.mdpi.com/2076-3417/12/20/10578
18. S. Naik, K. A. Hashmi, A. Pagani, M. Liwicki, D. Stricker, and M. Z. Afzal, "Investigating attention mechanism for page object detection in document images," *Applied Sciences*, vol. 12, no. 15, 2022. [Online]. Available: https://www.mdpi.com/2076-3417/12/15/7486
19. T. Fredriksson, D. Issa Mattos, J. Bosch, and H. Olsson, *Data Labeling: An Empirical Investigation into Industrial Challenges and Mitigation Strategies*, 11 2020, pp. 202–216.
20. J. E. Van Engelen and H. H. Hoos, "A survey on semi-supervised learning," *Machine learning*, vol. 109, no. 2, pp. 373–440, 2020.

21. I. Radosavovic, P. Dollár, R. B. Girshick, G. Gkioxari, and K. He, "Data distillation: Towards omni-supervised learning," *CoRR*, vol. abs/1712.04440, 2017. [Online]. Available: http://arxiv.org/abs/1712.04440

22. B. Zoph, G. Ghiasi, T.-Y. Lin, Y. Cui, H. Liu, E. D. Cubuk, and Q. Le, "Rethinking pre-training and self-training," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 3833–3845. [Online]. Available: https://proceedings.neurips.cc/paper/2020/file/27e9661e033a73a6ad8cefcde965c54d-Paper.pdf

23. Y. Li, D. Huang, D. Qin, L. Wang, and B. Gong, "Improving object detection with selective self-supervised self-training," *CoRR*, vol. abs/2007.09162, 2020. [Online]. Available: https://arxiv.org/abs/2007.09162

24. K. Wang, X. Yan, D. Zhang, L. Zhang, and L. Lin, "Towards human-machine cooperation: Self-supervised sample mining for object detection," *CoRR*, vol. abs/1803.09867, 2018. [Online]. Available: http://arxiv.org/abs/1803.09867

25. P. Tang, C. Ramaiah, R. Xu, and C. Xiong, "Proposal learning for semi-supervised object detection," *CoRR*, vol. abs/2001.05086, 2020. [Online]. Available: https://arxiv.org/abs/2001.05086

26. P. K. Rhee, E. Erdenee, S. D. Kyun, M. U. Ahmed, and S. Jin, "Active and semi-supervised learning for object detection with imperfect data," *Cognitive Systems Research*, vol. 45, pp. 109–123, 2017. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1389041716301127

27. Q. Xie, Z. Dai, E. H. Hovy, M. Luong, and Q. V. Le, "Unsupervised data augmentation," *CoRR*, vol. abs/1904.12848, 2019. [Online]. Available: http://arxiv.org/abs/1904.12848

28. M. Xu, Z. Zhang, H. Hu, J. Wang, L. Wang, F. Wei, X. Bai, and Z. Liu, "End-to-end semi-supervised object detection with soft teacher," *CoRR*, vol. abs/2106.09018, 2021. [Online]. Available: https://arxiv.org/abs/2106.09018

29. H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. M. Ni, and H.-Y. Shum, "Dino: Detr with improved denoising anchor boxes for end-to-end object detection," 2022. [Online]. Available: https://arxiv.org/abs/2203.03605

30. X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable DETR: deformable transformers for end-to-end object detection," *CoRR*, vol. abs/2010.04159, 2020. [Online]. Available: https://arxiv.org/abs/2010.04159

31. Z. Gao, L. Wang, B. Han, and S. Guo, "Adamixer: A fast-converging query-based object detector," 2022. [Online]. Available: https://arxiv.org/abs/2203.16507

32. Z. Zong, G. Song, and Y. Liu, "Detrs with collaborative hybrid assignments training," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 6748–6758.

33. T. Shehzadi, K. A. Hashmi, D. Stricker, and M. Z. Afzal, "Object detection with transformers: A review," 2023.

34. Z. Chen, J. Zhang, and D. Tao, "Recurrent glimpse-based decoder for detection with transformer," *CoRR*, vol. abs/2112.04632, 2021. [Online]. Available: https://arxiv.org/abs/2112.04632

35. F. Li, H. Zhang, S. Liu, J. Guo, L. M. Ni, and L. Zhang, "Dn-detr: Accelerate detr training by introducing query denoising," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 619–13 627.

36. S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, and L. Zhang, "DAB-DETR: dynamic anchor boxes are better queries for DETR," *CoRR*, vol. abs/2201.12329, 2022. [Online]. Available: https://arxiv.org/abs/2201.12329

37. D. Jia, Y. Yuan, H. He, X. Wu, H. Yu, W. Lin, L. Sun, C. Zhang, and H. Hu, "Detrs with hybrid matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19 702–19 712.

38. Y. Zhao, Y. Cai, W. Wu, and W. Wang, "Explore faster localization learning for scene text detection," in *2023 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2023, pp. 156–161.

39. G. Zhang, Z. Luo, Y. Yu, K. Cui, and S. Lu, "Accelerating detr convergence via semantic-aligned matching," 2022.

40. K. Itonori, "Table structure recognition based on textblock arrangement and ruled line position," in *Proceedings of 2nd International Conference on Document Analysis and Recognition (ICDAR '93)*, 1993, pp. 765–768.

41. S. Tupaj, Z. Shi, C. H. Chang, and H. Alam, "Extracting tabular information from text files," *EECS Department, Tufts University, Medford, USA*, vol. 1, 1996.

42. S. Chandran and R. Kasturi, "Structural recognition of tabulated data," in *Proceedings of 2nd International Conference on Document Analysis and Recognition (ICDAR '93)*, 1993, pp. 516–519.

43. Y. Hirayama, "A method for table structure analysis using dp matching," in *Proceedings of 3rd International Conference on Document Analysis and Recognition*, vol. 2, 1995, pp. 583–586 vol.2.

44. S. A. Siddiqui, M. I. Malik, S. Agne, A. Dengel, and S. Ahmed, "Decnt: Deep deformable cnn for table detection," *IEEE Access*, vol. 6, pp. 74 151–74 161, 2018.

45. K. A. Hashmi, A. Pagani, M. Liwicki, D. Stricker, and M. Z. Afzal, "Castabdetectors: Cascade network for table detection in document images with recursive feature pyramid and switchable atrous convolution," *Journal of Imaging*, vol. 7, 2021.

46. D. Nazir, K. A. Hashmi, A. Pagani, M. Liwicki, D. Stricker, and M. Z. Afzal, "Hybridtabnet: Towards better table detection in scanned document images," *Applied Sciences*, vol. 11, no. 18, 2021. [Online]. Available: https://www.mdpi.com/2076-3417/11/18/8396

47. P. Pyreddy and W. B. Croft, "Tintin: a system for retrieval in text tables," in *Digital library*, 1997.

48. A. Pivk, P. Cimiano, Y. Sure, M. Gams, V. Rajkovič, and R. Studer, "Transforming arbitrary tables into logical form with tartar," *Data & Knowledge Engineering*, vol. 60, no. 3, pp. 567–595, 2007. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0169023X06000620

49. S. Khusro, A. Latif, and I. Ullah, "On methods and tools of table detection, extraction and annotation in pdf documents," *Journal of Information Science*, vol. 41, no. 1, pp. 41–57, 2015.

50. D. W. Embley, M. Hurst, D. Lopresti, and G. Nagy, "Table-processing paradigms: a research survey," *International Journal of Document Analysis and Recognition (IJDAR)*, vol. 8, no. 2, pp. 66–86, 2006.

51. F. Cesarini, S. Marinai, L. Sarti, and G. Soda, "Trainable table location in document images," in *2002 International Conference on Pattern Recognition*, vol. 3, 2002, pp. 236–240 vol.3.

52. A. C. e. Silva, "Learning rich hidden markov models in document analysis: Table location," in *2009 10th International Conference on Document Analysis and Recognition*, 2009, pp. 843–847.

53. A. Silva, "Parts that add up to a whole: a framework for the analysis of tables," *Edinburgh University, UK*, 2010.

54. T. Kasar, P. Barlas, S. Adam, C. Chatelain, and T. Paquet, "Learning to detect tables in scanned document images using line information," in *2013 12th International Conference on Document Analysis and Recognition*. IEEE, 2013, pp. 1185–1189.

55. X. Yang, M. E. Yümer, P. Asente, M. Kraley, D. Kifer, and C. L. Giles, "Learning to extract semantic structure from documents using multimodal fully convolutional neural network," *CoRR*, vol. abs/1706.02337, 2017. [Online]. Available: http://arxiv.org/abs/1706.02337

56. D. He, S. Cohen, B. Price, D. Kifer, and C. L. Giles, "Multi-scale multi-task fcn for semantic page segmentation and table detection," in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 01, 2017, pp. 254–261.

57. I. Kavasidis, S. Palazzo, C. Spampinato, C. Pino, D. Giordano, D. Giuffrida, and P. Messina, "A saliency-based convolutional neural network for table and chart detection in digitized documents," *CoRR*, vol. abs/1804.06236, 2018. [Online]. Available: http://arxiv.org/abs/1804.06236

58. S. Paliwal, V. D, R. Rahul, M. Sharma, and L. Vig, "Tablenet: Deep learning model for end-to-end table detection and tabular data extraction from scanned document images," *CoRR*, vol. abs/2001.01469, 2020. [Online]. Available: http://arxiv.org/abs/2001.01469

59. L. Gao, Y. Huang, H. Déjean, J.-L. Meunier, Q. Yan, Y. Fang, F. Kleber, and 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.

60. X. Zhong, J. Tang, and A. J. Yepes, "Publaynet: largest dataset ever for document layout analysis," in *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, Sep. 2019, pp. 1015–1022.

61. A. Mondal, P. Lipps, and C. V. Jawahar, "IIIT-AR-13K: A new dataset for graphical object detection in documents," *CoRR*, vol. abs/2008.02569, 2020. [Online]. Available: https://arxiv.org/abs/2008.02569

62. M. C. Göbel, T. Hassan, E. Oro, and G. Orsi, "Icdar 2013 table competition," *2013 12th International Conference on Document Analysis and Recognition*, pp. 1449–1453, 2013.

63. L. Gao, X. Yi, Z. Jiang, L. Hao, and Z. Tang, "Icdar2017 competition on page object detection," in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 01, 2017, pp. 1417–1422.

64. M. Li, L. Cui, S. Huang, F. Wei, M. Zhou, and Z. Li, "Tablebank: A benchmark dataset for table detection and recognition," 2019.

65. B. Smock, R. Pesala, and R. Abraham, "PubTables-1M: Towards comprehensive table extraction from unstructured documents," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 4634–4642.

66. J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," *CoRR*, vol. abs/1411.4038, 2014. [Online]. Available: http://arxiv.org/abs/1411.4038

67. X.-H. Li, F. Yin, and C.-L. Liu, "Page object detection from pdf document images by deep structured prediction and supervised clustering," in *2018 24th International Conference on Pattern Recognition (ICPR)*, 2018, pp. 3627–3632.

68. M. Holecek, A. Hoskovec, P. Baudis, and P. Klinger, "Line-items and table understanding in structured documents," *CoRR*, vol. abs/1904.12577, 2019. [Online]. Available: http://arxiv.org/abs/1904.12577

69. P. Riba, L. Goldmann, O. R. Terrades, D. Rusticus, A. Fornés, and J. Lladós, "Table detection in business document images by message passing networks," *Pattern Recognition*, vol. 127, p. 108641, 2022. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0031320322001224

70. M. Minouei, K. A. Hashmi, M. R. Soheili, M. Z. Afzal, and D. Stricker, "Continual learning for table detection in document images," *Applied Sciences*, vol. 12, no. 18, 2022. [Online]. Available: https://www.mdpi.com/2076-3417/12/18/8969

71. A. Kölsch, M. Z. Afzal, M. Ebbecke, and M. Liwicki, "Real-time document image classification using deep cnn and extreme learning machines," in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 01, 2017, pp. 1318–1323.

72. L. Hao, L. Gao, X. Yi, and Z. Tang, "A table detection method for pdf documents based on convolutional neural networks," *2016 12th IAPR Workshop on Document Analysis Systems (DAS)*, pp. 287–292, 2016.

73. X. Yi, L. Gao, Y. Liao, X. Zhang, R. Liu, and Z. Jiang, "Cnn based page object detection in document images," in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 01, 2017, pp. 230–235.

74. T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *CoRR*, vol. abs/1708.02002, 2017. [Online]. Available: http://arxiv.org/abs/1708.02002

75. Y. Fang, B. Liao, X. Wang, J. Fang, J. Qi, R. Wu, J. Niu, and W. Liu, "You only look at one sequence: Rethinking transformer in vision through object detection," *CoRR*, vol. abs/2106.00666, 2021. [Online]. Available: https://arxiv.org/abs/2106.00666

76. K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask R-CNN," *CoRR*, vol. abs/1703.06870, 2017. [Online]. Available: http://arxiv.org/abs/1703.06870

77. Z. Cai and N. Vasconcelos, "Cascade R-CNN: delving into high quality object detection," *CoRR*, vol. abs/1712.00726, 2017. [Online]. Available: http://arxiv.org/abs/1712.00726

78. N. D. Vo, K. Nguyen, T. V. Nguyen, and K. Nguyen, "Ensemble of deep object detectors for page object detection," in *Proceedings of the 12th International Conference on Ubiquitous Information Management and Communication*, ser. IMCOM '18. New York, NY, USA: Association for Computing Machinery, 2018. [Online]. Available: https://doi.org/10.1145/3164541.3164644

79. A. Gilani, S. R. Qasim, I. Malik, and F. Shafait, "Table detection using deep learning," in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 01, 2017, pp. 771–776.

80. Y. Huang, Q. Yan, Y. Li, Y. Chen, X. Wang, L. Gao, and Z. Tang, "A yolo-based table detection method," in *2019 International Conference on Document Analysis and Recognition (ICDAR)*, 2019, pp. 813–818.

81. X. Zheng, D. Burdick, L. Popa, and N. X. R. Wang, "Global table extractor (GTE): A framework for joint table identification and cell structure recognition using visual context," *CoRR*, vol. abs/2005.00589, 2020. [Online]. Available: https://arxiv.org/abs/2005.00589

82. D. Prasad, A. Gadpal, K. Kapadni, M. Visave, and K. Sultanpure, "Cascadetabnet: An approach for end to end table detection and structure recognition from image-based documents," *CoRR*, vol. abs/2004.12629, 2020. [Online]. Available: https://arxiv.org/abs/2004.12629

83. M. Agarwal, A. Mondal, and C. V. Jawahar, "Cdec-net: Composite deformable cascade network for table detection in document images," in *2020 25th International Conference on Pattern Recognition (ICPR)*, 2021, pp. 9491–9498.

84. T. Shehzadi, K. A. Hashmi, D. Stricker, M. Liwicki, and M. Z. Afzal, "Bridging the performance gap between detr and r-cnn for graphical object detection in document images," *arXiv preprint arXiv:2306.13526*, 2023.

85. S. Arif and F. Shafait, "Table detection in document images using foreground and background features," in *2018 Digital Image Computing: Techniques and Applications (DICTA)*, 2018, pp. 1–8.

86. S. A. Siddiqui, M. I. Malik, S. Agne, A. Dengel, and S. Ahmed, "Decnt: Deep deformable cnn for table detection," *IEEE Access*, vol. 6, pp. 74 151–74 161, 2018.

87. J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, "Deformable convolutional networks," *CoRR*, vol. abs/1703.06211, 2017. [Online]. Available: http://arxiv.org/abs/1703.06211

88. Y. Liu, Y. Wang, S. Wang, T. Liang, Q. Zhao, Z. Tang, and H. Ling, "Cbnet: A novel composite backbone network architecture for object detection," *CoRR*, vol. abs/1909.03625, 2019. [Online]. Available: http://arxiv.org/abs/1909.03625

89. J. Jeong, S. Lee, J. Kim, and N. Kwak, "Consistency-based semi-supervised learning for object detection," in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019. [Online]. Available: https://proceedings.neurips.cc/paper/2019/file/d0f4dae80c3d0277922f8371d5827292-Paper.pdf

90. P. Tang, C. Ramaiah, R. Xu, and C. Xiong, "Proposal learning for semi-supervised object detection," *CoRR*, vol. abs/2001.05086, 2020. [Online]. Available: https://arxiv.org/abs/2001.05086

91. T. Shehzadi, K. A. Hashmi, A. Pagani, M. Liwicki, D. Stricker, and M. Z. Afzal, "Mask-aware semi-supervised object detection in floor plans," *Applied Sciences*, vol. 12, no. 19, 2022. [Online]. Available: https://www.mdpi.com/2076-3417/12/19/9398

92. G. Kallempudi, K. A. Hashmi, A. Pagani, M. Liwicki, D. Stricker, and M. Z. Afzal, "Toward semi-supervised graphical object detection in document images," *Future Internet*, vol. 14, no. 6, 2022. [Online]. Available: https://www.mdpi.com/1999-5903/14/6/176

93. T. Shehzadi, K. A. Hashmi, D. Stricker, and M. Z. Afzal, "Sparse semi-detr: Sparse learnable queries for semi-supervised object detection," *arXiv preprint arXiv:2404.01819*, 2024.

94. K. Sohn, Z. Zhang, C. Li, H. Zhang, C. Lee, and T. Pfister, "A simple semi-supervised learning framework for object detection," *CoRR*, vol. abs/2005.04757, 2020. [Online]. Available: https://arxiv.org/abs/2005.04757

95. K. Wang, X. Yan, D. Zhang, L. Zhang, and L. Lin, "Towards human-machine cooperation: Self-supervised sample mining for object detection," *CoRR*, vol. abs/1803.09867, 2018. [Online]. Available: http://arxiv.org/abs/1803.09867

96. N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*. Springer, 2020, pp. 213–229.

97. T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," *CoRR*, vol. abs/1405.0312, 2014. [Online]. Available: http://arxiv.org/abs/1405.0312

98. D. M. W. Powers, "Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation," *CoRR*, vol. abs/2010.16061, 2020. [Online]. Available: https://arxiv.org/abs/2010.16061

99. T. Shehzadi, K. Azeem Hashmi, D. Stricker, M. Liwicki, and M. Zeshan Afzal, "Towards end-to-end semi-supervised table detection with deformable transformer," in *Document Analysis and Recognition - ICDAR 2023*, G. A. Fink, R. Jain, K. Kise, and R. Zanibbi, Eds.   Cham: Springer Nature Switzerland, 2023, pp. 51–76.

100. Y. Liu, C. Ma, Z. He, C. Kuo, K. Chen, P. Zhang, B. Wu, Z. Kira, and P. Vajda, "Unbiased teacher for semi-supervised object detection," *CoRR*, vol. abs/2102.09480, 2021. [Online]. Available: https://arxiv.org/abs/2102.09480

101. Y. Tang, W. Chen, Y. Luo, and Y. Zhang, "Humble teachers teach better students for semi-supervised object detection," *CoRR*, vol. abs/2106.10456, 2021. [Online]. Available: https://arxiv.org/abs/2106.10456

102. P. Zhang, C. Li, L. Qiao, Z. Cheng, S. Pu, Y. Niu, and F. Wu, "VSR: A unified framework for document layout analysis combining vision, semantics and relations," *CoRR*, vol. abs/2105.06220, 2021. [Online]. Available: https://arxiv.org/abs/2105.06220