# Bridging the Performance Gap between DETR and R-CNN for Graphical Object Detection in Document Images.

Tahira Shehzadi[1,2,3*], Khurram Azeem Hashmi[1,2,3], Didier Stricker[1,2,3], Marcus Liwicki[4] and Muhammad Zeshan Afzal[1,2,3]

[1]Department of Computer Science, Technical University of Kaiserslautern, 67663, Germany.
[2]Mindgarage Lab, Kaiserslautern, 67663, Germany.
[3]Augmented Vision, German Research Institute for Artificial Intelligence (DFKI), Kaiserslautern, 67663, Germany.
[4]Department of Computer Science, Luleå University of Technology, Luleå, 97187, Sweden.

*Corresponding author(s). E-mail(s): tahira.shehzadi@dfki.de;

## Abstract

This paper takes an important step in bridging the performance gap between DETR and R-CNN for graphical object detection. Existing graphical object detection approaches have enjoyed recent enhancements in CNN-based object detection methods, achieving remarkable progress. Recently, Transformer-based detectors have considerably boosted the generic object detection performance, eliminating the need for hand-crafted features or post-processing steps such as Non-Maximum Suppression (NMS) using object queries. However, the effectiveness of such enhanced transformer-based detection algorithms has yet to be verified for the problem of graphical object detection. Essentially, inspired by the latest advancements in the DETR, we employ the existing detection transformer with few modifications for graphical object detection. We modify object queries in different ways, using points, anchor boxes and adding positive and negative noise to the anchors to boost performance. These modifications allow for better handling of objects with varying sizes and aspect ratios, more robustness to small variations in object positions and sizes, and improved image discrimination

between objects and non-objects. We evaluate our approach on the four graphical datasets: PubTables, TableBank, NTable and PubLaynet. Upon integrating query modifications in the DETR, we outperform prior works and achieve new state-of-the-art results with the mAP of 96.9%, 95.7% and 99.3% on TableBank, PubLaynet, PubTables, respectively. The results from extensive ablations show that transformer-based methods are more effective for document analysis analogous to other applications. We hope this study draws more attention to the research of using detection transformers in document image analysis.

# 1 Introduction

Over the past few decades, the rapid growth in digital transformation in different forms like web pages, scientific publications, invoices, and financial statements has increased the storage and production of digital documents [1]. Digitization and information extraction from such a large collection of documents is impossible for humans. So, many researchers are attracted to automatic information extraction from document images. These documents contain text and graphical page objects such as formulas, figures and tables. Even though modern Optical Character Recognition (OCR) networks [2–4] can extract text information, they fail to interpret graphical page objects. The graphical object detection task in document analysis is challenging as the documents may have high variability in layout, complex backgrounds, small object size, similarity with text, and limited training data. Thus, it is essential to have an accurate graphical object detection system for document analysis.

For the graphical object detection task, previous methods used rule-based approaches [5–8] and, then, CNN-based object detectors such as R-CNN [9, 10], Faster R- CNN [11], and Cascade Mask R-CNN [12] have been presented. Thus, the improvement in object detection networks is directly reflected in state-of-the-art graphical object detection systems. Earlier, people used the classical detectors in the object detection domain and transformer-based networks in the sequence prediction domain. Recently, Carion et al. [13] proposed a transformer-based architecture for object detection that performs better than all CNN-based object detectors. For the first time, Smock et al. [14] used a simple DEtection TRansformer (DETR) framework for table detection and structure recognition tasks on the PubTables [14] dataset and observed excellent results. However, rapid progress in transformers still need to be employed in document analysis.

Transformer-based object detectors [15] remove the need for hand-designed components like non-maximum suppression (NMS) and anchor design used

in CNN-based object detectors by introducing the concept of object queries. Object queries are learned vectors that are used to attend to all the objects in the image simultaneously and predict their class and location. In contrast to CNN-based object detectors, which generate a fixed number of proposals or anchors that are then classified and refined, transformer-based detectors use an attention mechanism to dynamically attend to relevant parts of the image and predict the objects. This allows them to detect objects of different sizes and aspect ratios without the need for hand-designed anchor boxes. Additionally, because transformer-based detectors do not rely on a fixed set of proposals, they can avoid the need for post-processing steps such as NMS. This simplifies the training and inference pipelines and can lead to faster and more accurate object detection.

Modifying object queries with points, anchor boxes, positive noise, and negative noise can improve the performance of graphical page object detection by providing additional localization cues, handling non-standard aspect ratios more effectively, and better differentiating between foreground and background graphical objects. These modifications can also make training more efficient and improve the generalization capabilities of the detector. By incorporating these techniques, graphical object detection models can achieve higher accuracy and efficiency, while reducing false positives and improving the overall detection performance. This paper helps in bridging the performance gap between DETR and R-CNN for graphical object detection in document images. We use the potential of the detection transformer with improved object queries. Additionally, two different pre-processing approaches are applied to the training process, further boosting the table detection performance and helping it learn more effectively. We evaluate transformer-based detectors on four popular graphical object detection datasets: TableBank [16], Publaynet [17], NTables [18] and PubTables [14] and compare their performance with CNN-based detectors [19–24].

- We present an end-to-end trainable graphical object detection framework that operates on a DEtection TRansformer (DETR) equipped with improved object queries.
- To the best of our knowledge, for the first time, extensive experiments are conducted on transformers for the graphical object detection domain and leverage the potential of detection transformers for this task.
- We use different pre-processing techniques that boost the graphical object detection performance.
- We accomplish state-of-the-art performance of transformer-based detectors on four publicly available graphical object detection datasets in scanned document images and camera-captured images.

The remaining part of this paper is arranged as follows. Section 2 discusses the previous work on graphical object detection using popular computer vision

approaches. Section 3 explains transformer-based detector and its submodules. Section 4 discusses the implementation details and performance analysis. Section 5 concludes the paper and discusses future directions.

# 2 Related work

## 2.1 Traditional Approaches

Previous methods used rule-based approaches for graphical objects such as table detection problems in the document analysis domain. These approaches locate tables in predefined scenarios using different rules like vertical and horizontal lines [5, 6], text layouts [25], keywords [7], or formal patterns [8]. These approaches need manual work for tuning hyper-parameters and designing rules. Many machine learning approaches [26, 27] are proposed to minimize these heuristics dependencies. These machine learning-based approaches are explained briefly in [28]. Even though these approaches boost table detection performance, they use handcrafted features, Thus these methods cannot be generalized.

## 2.2 Deep Learning Based Approaches

With the increasing progress in deep learning, many CNN-based table detection approaches have been presented that significantly improve performance. The researchers classified these approaches into semantic segmentation, bottom-up and object detection approaches.

**Semantic Segmentation Based Approaches.** Graphical object detection such as table detection is also considered a segmentation problem, and applied the segmentation mask using current semantic segmentation networks like Fully-Convolutional Networks (FCN) [29] to detect tables at the pixel level in [30–32]. In [30], the author designed a multimodal-based FCN for graphical objects in document analysis using linguistic and visual features and improved the segmentation performance. He et al. [31] presented a multi-task and multi-scale-based FCN for predicting masks for the table, figure, text and their related contours. Then they filtered these segmentation masks by the Conditional Random field (CRF) network to get table areas in document images. In [32], the author gave a saliency-based FCN for multi-scales reasoning with CRF to detect charts and tables in digital documents. However, detection transformers can handle objects of different sizes and aspect ratios more effectively than semantic segmentation approaches, as they do not rely on fixed-size segmentation masks. This makes detection transformers more suitable for detecting objects with irregular shapes or objects that are densely packed together.

**Bottom-up Approaches.** The bottom-up approaches analyze the text documents as graphs where objects (e.g., text, table) are considered nodes, and

edges show the relationship between these objects to solve the table detection as a graph labelling problem. In [33], the author applied popular layout analysis approaches to form clusters and classify page objects (table, figure, text, formula) using CNN-based CRF networks. The overlapping regions of the same cluster and class merged to detect page objects. In [34, 35], the authors represented text areas (text lines, words) as nodes and page layouts as graphs. They then used graph neural networks to classify edges and nodes. Finally, they extracted the table class from subgraphs where detected tables are present in the nodes. Li et al. [36] presented a document analysis problem as a sequence-labelling problem. They classified each word in the sequence into predefined object categories, including tables, using pre-trained language models. All these approaches also need correct text/word bounding boxes as an extra input. However, detection transformers can be more efficient than bottom-up approaches for graphical page object detection, as they can process the entire image in a single forward pass through the neural network, without the need for additional post-processing steps.

**Object Detection-based Approaches.** In [9, 10], the authors used R-CNN [37] for detecting tables. However, the performance of these methods depends on handcrafted features and heuristic rules, as in earlier approaches. The researchers [12, 38–40] also used more advanced single-stage [41, 42] and two-stage [19–21, 43] object detection approaches to detect tables, formulas and figures in the document analysis domain. Further, augmentation approaches are used to boost the performance of these networks for detecting tables. As Prasad et al. [40], Arif et al. [44] and Gilani et al. [45] used image transformation methods of dilation, coloration and distance transformation to increase the training data to get more information as features from input table images. In [46], the author integrated deformable convolution with deformable (Region of Interest) RoI pooling network in Faster R-CNN [19] to make the network more efficient for the geometric transformation such as translation, rotation, reflection and dilation. Sun et al. [11] enhance the localization performance by using Faster R-CNN [19] for corner detection and adjusting the table bounding boxes by corners using a post-processing network. Anyhow, these corner boxes are handcrafted, and their size has no clear meaning that increases the false detection rate of corner boxes [11]. Agarwal et al. [12] used a composite backbone network along with deformable-convolution filters in Cascade Mask R-CNN [43] to boost detection results in the table analysis domain. Even though these methods improve the results on many benchmark datasets [16, 17, 47–51], they use large memory and have high computational complexity. Compared to CNN, transformer networks make predictions without needing non-maximum suppression, making the network more simple and efficient. Recently, Smock et al. [14] used a simple detection transformer (DETR) framework for table detection and structure recognition tasks on the PubTables dataset and observed excellent results.

## 2.3 Object Detection with Transformers

Previously, people used classical detectors [19–24, 52] in the object detection domain and transformer-based networks in the sequence prediction domain [53]. Recently, Carion et al. [13] proposed a transformer-based architecture for object detection. This Detection transformer (DETR) [54] uses image features as learning queries and predicts the bounding boxes using bipartite graph matching. This architecture removes the need for hand-designed anchors [55] and non-maximum suppression (NMS) to provide a simpler and optimal object detection framework [56]. Compared with anchor-based object detectors [19–24, 37, 57], transformers consider object detection a prediction problem and extract features from an image using learning queries that eliminate the need for non-maximum suppression. However, transformer-based detector like DETR has slower convergence during training than previous detectors like Faster R-CNN [19]. For example, DETR needs 500 epochs, while Faster R-CNN needs 12 epochs to get the same performance on the coco detection dataset.

Several works [58–63] have tried to find the main reason for the slow convergence of transformer-based detection networks. Some of them tried to improve the encoder-decoder architecture. Sun et al. [60] proposed only encoder-based DETR by considering the low performance of the attention module in the decoder network as the main reason for slow training convergence. Dai et al. [63] proposed the regions-of-interest-based dynamic decoder that uses Region of Interest (ROI). Recent works [58, 59, 61, 62] focus on spatial positions than multiple positions for DETR queries to extract features from the image. Conditional DETR [58] breaks down every query into positional and content segments to have clear similarities with the spatial region in the image. Deformable-DETR [62] apply deformable attention mechanism in the decoder module to improve training convergence. Efficient-DETR [64] proposes a dense prediction network to take top-K object queries. But all this work uses only 2D points as anchor regions and ignores object scaling. However, DAB-DETR [59] takes 4D coordinates as learnable queries and continuously improves them in every layer. Recently, DN-DETR [65] indicated that the bipartite matching approach used in Hungarian loss is one of the reasons for transformer slow training convergence and presented a denoising training-based method to overcome the slow training convergence issue. Following DAB-DETR and DN-DETR, DINO [66] proposed the Contrastive DeNoising(CDN) module that adds extra Denoising (DN) loss. All these modified versions of DETR make the transformer network faster and boost performance. We combine all these modifications for graphical object detection task and observe a remarkable performance boost.

# 3 Method

This section first describes an overview of the transformer-based detector's main modules and then employs different mechanisms to improve the object queries. The whole network is shown in Figure 1.
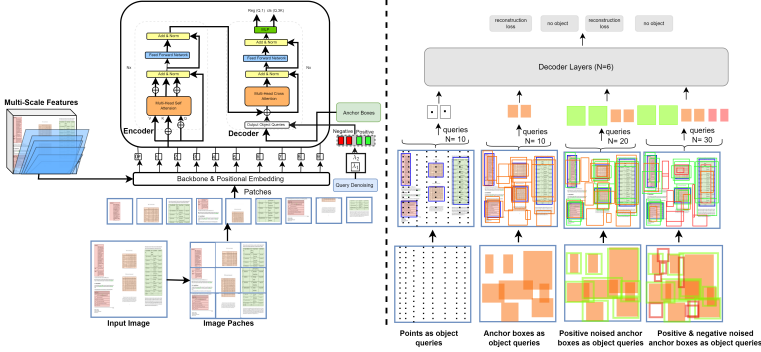


**Fig. 1**: Presented Transformer-based graphical object detection framework. We divide the input image into small equal-size patches, add position embeddings, and embed the resulting patches along with input multi-scale features to the transformer encoder. We use different forms of object queries such as points, anchor boxes, the addition of positive and negative noises to the anchor boxes in the decoder and observe network performance. Here, blue rectangles represent ground truth (GT), black dots represent points as object queries, brown rectangles denote anchor boxes as object queries, green rectangles indicate positive noised anchor boxes for foreground class, and red rectangles show negative noised anchor boxes for the background class. These object queries are taken as decoder input to provide final graphical object labels and locations.

## 3.1 Revisiting DETR

**Positional Encoding and Backbone Network** The transformer encoder takes a 1D sequence as input token embeddings. The backbone network ResNet-50 [67] extracts the features from the input, reduces the channel dimension, and converts the spatial dimension into one dimension as the transformer network takes input as one vector. We transform the input image $I \in R^{(H_i \times W_i \times C_i)}$ into equal size 2D patches $I_p \in R^{(N_i \times (P_i^2 . C_i))}$. Where $(H_i \times W_i)$ is the original image spatial resolution, $C_i$ is the channels/bands number, $(P_i.P_i)$ is the resolution of each image patch, $N_i = \frac{H_i.W_i}{(P_i.P_i)}$ is the total number of input patches, also considered as the transformer's input sequence length. The standard transformer model also considers any spatial relationships between the input features. So the encoder network takes two inputs: first is the feature vector of the input image as small image patches, and second is

the positional encoding of the input feature vector. The resulting sequence $Z_0$ is fed as input to the transformer encoder as follows:

$$Z_0 = I_p + M_{PE} \tag{1}$$

Here, $I_p \in R^{(N_i \times (P_i^2 . C_i))}$ is the input image patches and $M_{PE} \in R^{(N_i+1) \times D}$ represents the positional embeddings. The positional encoding contains information of the absolute or relative position of the patches to make use of the spatial information.

**Encoder-Decoder Network** The extracted features from the backbone network as one single vector and their position within the input vector are fed to the encoder network. Here, the self-attention layer provides key, query and value matrices which feed to the multi-head attention to find the attention probabilities of the input vector. The DETR decoder takes N number of object queries in parallel with the encoder output. During training, the model is trained to generate the same set N of predicted bounding boxes regardless of the order of the object queries. DETR searches for a permutation of N elements $\sigma \in N$ that results in the lowest cost for the matching. The cost is defined as the sum of the negative log-likelihood of the predicted class probabilities and the box loss for each matched object as follows:

$$\hat{\sigma} = \arg\min_{\sigma \in N} \sum_{k}^{N} \mathcal{L}_{match}(y_k, \hat{y}_{\sigma(k)}), \tag{2}$$

Here, $y_k$ is a set of ground truth objects, $\hat{y}_{\sigma(k)}$ is the predicted objects. By training the model using a permutation-invariant loss, DETR is able to learn to detect objects in an image regardless of their order, which makes it more robust to variations in the input data. The term $\mathcal{L}_{match}(y_k, \hat{y}\sigma(k))$ is the one-to-one matching cost for direct prediction without duplicates between predicted objects and ground truth as shown in the following equation:

$$\mathcal{L}_{match}(y_k, \hat{y}_{\sigma(k)}) = -\mathbb{1}_{\{c_k \neq \phi\}}\hat{p}_{\sigma(k)}(c_k) + \mathbb{1}_{\{c_k \neq \phi\}}\mathcal{L}_{bbox}(b_k, \hat{b}_{\sigma(k)}) \tag{3}$$

Where $\hat{p}_{\sigma(k)}$ and $c_k$ are the predicted class labels and target labels, respectively, $b_k$ and $\hat{b}_{\hat{\sigma}}(k)$ are ground truth and predicted bounding boxes, respectively.

To match the object queries N with the ground-truth objects in the image, DETR uses the Hungarian algorithm to find a bipartite matching between the set of object queries and the set of ground-truth objects. To do this, the set of ground-truth objects is padded with empty objects ($\phi$) to make it the same size as the set of object queries. The next step is to compute the Hungarian loss $\mathcal{L}_H$ in Equation (4) by determining the optimal matching between ground truth (GT) and detected boxes regarding bounding box location and class. It is defined as a linear combination of two terms: a negative log-likelihood term for the class predictions, and a box loss term for the predicted bounding boxes as follows.

$$\mathcal{L}_H(y, \hat{y}) = \sum_{i=1}^{N} [-log\hat{p}_{\hat{\sigma}(k)}(c_k) + \mathbb{1}_{\{c_k \neq \phi\}}\mathcal{L}_{bbox}(b_k, \hat{b}_{\hat{\sigma}}(k))] \tag{4}$$

Where $\hat{\sigma}$ is the optimal-assignment factor from Equation (2). The negative log-likelihood term measures the difference between the predicted class probabilities and the ground-truth class labels. The box loss term measures the difference between the predicted bounding boxes and the ground-truth bounding boxes. This network doesn't need NMS to remove redundant predictions as it uses Hungarian loss that learns to make non-redundant predictions. However, the DETR network has several challenges, such as optimizing the network because of its slow training convergence and performance drops for small objects.

## 3.2 Proposed Advancements

**Object queries.**     Object queries can be viewed as a replacement for anchor boxes in object detection models, such as the popular Faster R-CNN [19]. During training, the model learns a fixed number of object queries that represent specific object classes. In inference, detection transformer directly outputs a set of object detection using a fixed number of object queries, eliminating the need for anchor boxes and non-maximum suppression. The use of object queries enables direct, end-to-end optimization for object detection and improves interpretability of the model's outputs. We modify the object queries to observe the effect of quantity and quality of object queries on detection transformer performance for the graphical object detection task.

**Object queries as points.**     Graphical object detection can be challenging due to the presence of multiple objects at the same location. For example, text, tables, and other graphical elements can overlap or appear in close proximity to each other. This issue can be resolved by using points as objects queries. We used two type of points as grid and learned points. Grid points are fixed points placed at regular intervals in the image, while learned points are initialized randomly and updated during training to better match the objects in the image. During training, the model learns a fixed number of points, which are used to define regions in the image where objects are likely to be present. The model then predicts the presence and location of objects near each point. This allows the model to predict multiple objects at one position, which is particularly useful for graphical object detection.

**Object queries as anchor boxes.**     In traditional object detection tasks [19, 68], anchor boxes are used as references during training and prediction to generate the final bounding box predictions. For graphical object detection such as tables, the object shapes and sizes are much more regular and predictable, since tables generally consist of rows and columns of consistent dimensions. The dynamic anchor boxes are better suited to the regular and predictable nature of table structures, and result in more accurate and reliable table boundary predictions.

In this case, decoder network takes positional queries as anchor boxes and needs keys $k_i$, queries $q_i$ and values $v_i$ having a cross-attention module to find features probing. Given a bounding box coordinates $a_i = (a'_i, y'_i, w'_i, h'_i)$ with content query $C_q$ to detect the contents of a table cell in a document image,

its positional query $P_i$ is formed as follows:

$$P_i = MLP(PE(a_i)) \tag{5}$$

$$PE(a_i) = Conc(PE(x'_i), PE(y'_i), PE(w'_i), PE(h'_i)) \tag{6}$$

Here, $PE(a_i)$ is the positional encoding, it is the overall positional encoding of the bounding box coordinates by finding the positional encoding of its component and then concatenating them by the *Conc* function. The positional query $P_i$ effectively captures the complex visual patterns present in document images, such as the varying thicknesses and styles of table borders, and the presence of different types of content within table cells.

In the self-attention module, the queries, keys, and values all have the same



**Fig. 2**: Illustrating design of attention module. The attention network only takes a finite number of samples near the reference point, irrespective of feature map size. Here, Sampling is performed at multiple scales using a deformation field computed by a deformable convolutional neural network. For clear visualization, we show sampling points with attention weights as a circle where the circle color represents its attention weight: blue indicates low intensity while red indicates high intensity. The rectangle represents the predicted bounding box in the decoder.

content information, while key and query has extra positional information $PE(a_i)$. In cross attention, the queries are derived from the dynamic anchor boxes $a_i$, while the keys and values are derived from the feature map. The multi-head self-attention consider deformable attention to allow the model to attend to different parts of the image at different resolutions and scales as shown in Figure 2. Here, attention network only takes a finite number of samples near the reference point, irrespective of feature map size. Considering only

a small number of keys for each query converges the network faster.

**Object queries as noised anchor boxes.** Anchor boxes can be noisy and imprecise, leading to inaccurate detections. Adding positive noise to better fit for foreground objects and negative noise to better fit for background objects modifies the anchor boxes used in the detection process to better fit the rectangular/square shape of graphical objects. For this, positive noise is added to the anchor boxes to expand their size and account for small variations in object size and shape. Negative noise, on the other hand, is used for detecting background objects as shown in Figure 1. By combining positive and negative noise, we can more accurately capture the range of object sizes and aspect ratios present in the data while filtering out noise and irrelevant regions. $\lambda_1$ controls the amount of positive noise added to the anchor boxes, while $\lambda_2$ controls the amount of negative noise. These values are learned during training and are used to adjust the size of the anchor boxes, allowing them to better match the objects in the image. The initial anchor boxes at the decoder input is represented as $a_i$, where $a_i = (x'_i, y'_i, w'_i, h'_i)$, while we have N number of GT boxes represented as $b_i$ where $b_i = (x_i, y_i, w_i, h_i)$. We remove those anchors that are farther away from the ground truth anchor as follows:

$$AMD(k) = \frac{1}{k}\Sigma\{Max_K(\{\parallel b_0 - a_0 \parallel_1, ..., \parallel b_{N-1} - a_{N-1} \parallel_1\}, k)\} \quad (7)$$

Where $(b_i - a_i)$ is the distance between bounding box b and anchor and $Max_K$ is the module that selects the top K anchors. Usually, the value of $\lambda_2$ is small to improve model performance as hard negative examples are closer ground truth boxes. If an image has N number of objects to be detected, we will have a total of 3N queries generating 2N positive and negative anchors for N ground truths.

## 3.3 Pre-processing Methods

We use some transformation techniques that help model learning more accurately during training. Document images contain blank spaces and content or text regions. As the network detects graphical objects in document data, we apply different transformation techniques to form the text or table regions thicker and reduce the blank space regions. For that, we use two types of approaches as smudge transform and dilation transform.

**Dilation Transform.** In this transformation, black pixel regions are thicker by transforming the original table image. First, we convert the original image into a binary image and then apply the dilation transform by a kernel filter (2x2 size) for one iteration. We select this kernel size because it gives the best results. The original image is represented as a), and the image after dilation is represented as b) in Figure 3.

**Smudge Transform.** In this transformation, black pixel regions spread around the end of the black regions. First, we convert the original image into

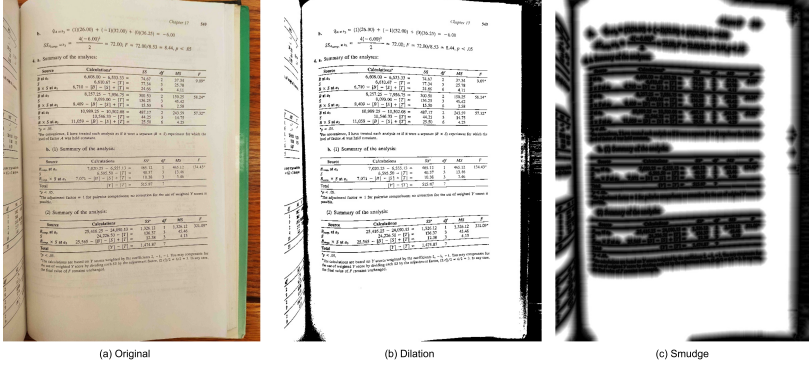(a) Original          (b) Dilation          (c) Smudge

**Fig. 3**: Pre-processing techniques to form the text or table regions thicker and reduce the blank space regions. Here, (a) shows the original data sample, (b) contains a dilation transformation-based augmentation sample, and (c) provides a smudge transformation-based augmentation sample.

a binary image and then apply the smudge transformation by various distance transforms. Gilani et al. [45] explained the original algorithm that applies Linear Distance Transform, Max Distance Transform and Euclidean Distance Transform to the input image. Furthermore, We also normalize and tune the parameters to boost the results. The original image is represented as a), and the image after the smudge transformation is represented as c) in Figure 3.

# 4  Experiments

## 4.1  Dataset and Evaluation Setup

The performance of the transformer-based table detection model is evaluated using precision, recall, F-Score and mean average precision (mAP) in the context of MS COCO [69] evaluation on the four most extensive graphical object detection datasets: TableBank [16], PubLayNet [17], PubTables [14] and NTable [18] dataset. Please refer to the supplementary material for a detailed explanation of the dataset and evaluation setup.

## 4.2  Results

### 4.2.1  TableBank

We validate the performance of detection transformer by modifying the object queries on the raw (without pre-processing), dilation and smudge transformation of the TableBank dataset in Table 1. Here, the term "both" represents combination of latex+word split. For the data group of TableBank$_{both}$ with dilation and smudge transformation, we achieve an mAP of 93.4% for grid points used as object queries, 95.1% for anchor boxes as object queries, 94.9% for positive noised anchor boxes as object queries and 96.9% for positive and

**Table 1**: Comparison between transformer-based detector results on raw (without pre-processing), dilation and smudge transformation of the Table-Bank dataset. Here, term $Q_b$ represents object queries as anchor boxes, $Q_p$ denotes object queries with positive noise and $Q_n$ indicates object queries with negative noise. The IoU thresholds are set to 0.5 and 0.75 for average precision calculation and also calculate average recall for large objects. AR represents Average Recall for a large area. The best results are highlighted.

| Methods | Preprocessing | mAP | $AP^{50}$ | $AP^{75}$ | AR |
|---|---|---|---|---|---|
| DETR | raw(latex) | 86.6 | 96.9 | 93.7 | 92.6 |
| | Dilation+Smudge | 87.2±0.21 | 97.5 | 94.2 | 93.5 |
| DETR + $Q_b$ | raw (latex) | 87.7 | 97.3 | 94.8 | 93.5 |
| | Dilation+Smudge | 88.2±0.53 | 98.0 | 95.4 | 94.3 |
| DETR + $Q_b$ + $Q_p$ | raw (latex) | 88.9 | 97.2 | 94.7 | 94.6 |
| | Dilation+Smudge | 89.4±0.61 | 97.6 | 94.9 | 95.4 |
| **DETR + $Q_b$ + $Q_p$ + $Q_n$** | **raw (latex)** | **91.5** | **97.4** | **94.8** | **97.8** |
| | **Dilation+Smudge** | **92.6±1.23** | **98.2** | **95.0** | **98.4** |
| DETR | raw(word) | 93.4 | 96.6 | 95.1 | 96.7 |
| | Dilation+Smudge | 93.9±0.31 | 97.3 | 95.9 | 97.2 |
| DETR + $Q_b$ | raw (word) | 92.3 | 97.4 | 95.9 | 95.0 |
| | Dilation+Smudge | 93.1±0.40 | 97.9 | 96.4 | 95.8 |
| DETR + $Q_b$ + $Q_p$ | raw (word) | 95.0 | 98.1 | 96.5 | 96.6 |
| | Dilation+Smudge | 95.5±1.5 | 98.7 | 96.9 | 96.9 |
| **DETR + $Q_b$ + $Q_p$ + $Q_n$** | **raw (word)** | **96.3** | **98.3** | **96.5** | **99.1** |
| | **Dilation+Smudge** | **96.8±1.24** | **98.8** | **96.9** | **99.7** |
| DETR | raw (both) | 92.5 | 97.2 | 95.1 | 96.8 |
| | Dilation+Smudge | 93.4±1.1 | 97.5 | 95.8 | 97.3 |
| DETR + $Q_b$ | raw (both) | 91.7 | 97.6 | 95.4 | 96.4 |
| | Dilation+Smudge | 95.1±1.23 | 97.8 | 96.9 | 97.4 |
| DETR + $Q_b$ + $Q_p$ | raw (both) | 94.3 | 98.5 | 97.1 | 97.7 |
| | Dilation+Smudge | 94.9±0.5 | 98.8 | 97.9 | 98.8 |
| **DETR + $Q_b$ + $Q_p$ + $Q_n$** | **raw (both)** | **95.8** | **98.9** | **97.2** | **98.8** |
| | **Dilation+Smudge** | **96.9±0.41** | **99.4** | **97.6** | **99.1** |

negative noised anchor boxes as object queries.

Figure 4 exhibits the performance analysis of detection transformers using different types of object queries as input using Average Precision (AP) with IoU threshold values ranging from 0.5 to 1 on all splits TableBank dataset. We can observe that detection transformer that uses points as object queries shows the lowest performance, while noised anchor boxes as object queries has the highest result on all threshold values. The qualitative analysis of table detection for the TableBank$_{both}$ dataset is illustrated in Figure 5. Analysis of incorrect results reveals that the network fails to localize accurate tabular areas or gives false positives.

**Comparison with State-of-the-art Methods** We also compare the results of the detection transformer using different type of object queries with earlier
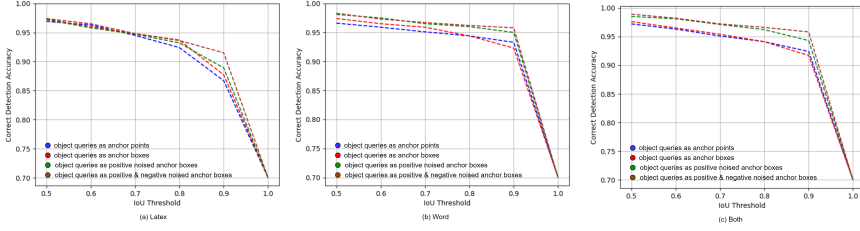
**Fig. 4**: Performance analysis of DETR and modifications in object queries in terms of AP over the IoU thresholds range from 0.5 to 1.0 on word, latex and both splits of the TableBank dataset. Here, the blue color shows results with simple DETR, the red highlights results with anchor boxes as object queries, the green represents results with positive noised anchor boxes, and the brown denotes results with positive and negative noised anchor boxes.



**Fig. 5**: Transformer-based detector with improved object queries (by adding positive and negative noised anchors) results on TableBank$_{both}$ dataset. Green color denotes true positives, blue exhibits false negatives and red exhibits false positives. Here, (a) exhibits true positive detection samples, (b) contains true positive and false positive detection samples, and (c) provides false negative detection.

CNN-based approaches on the TableBank dataset. Table 2 shows that transformer with noised anchor boxes has outperformed the previous state-of-the-art methods.

### 4.2.2 PubLayNet

We validate the performance of detection transformer by modifying the object queries on the raw, dilation and smudge transformation of the PubLayNet

**Table 2**: Comparison between the transformer-based detectors and previous state-of-the-art results on the TableBank dataset without pre-processing (raw data). Here, term $Q_b$ represents object queries as anchor boxes, $Q_p$ denotes object queries with positive noise and $Q_n$ indicates object queries with negative noise. The IoU threshold value is set to 0.5. The best results are highlighted.

| Model | Split | $AP^{50}$ | AR |
|---|---|---|---|
| CascadeTabNet [40] | TableBank$_{latex}$ | 95.9 | 97.2 |
| HybridTabNet [70] | TableBank$_{latex}$ | 97.7 | 98.3 |
| CasTabDetectoRS [71] | TableBank$_{latex}$ | 98.3 | 98.4 |
| **DETR + $Q_b$ + $Q_p$ + $Q_n$** | **TableBank$_{latex}$** | **97.4** | **97.8** |
| CascadeTabNet [40] | TableBank$_{word}$ | 94.3 | 95.5 |
| HybridTabNet [70] | TableBank$_{word}$ | 95.5 | 98.5 |
| CasTabDetectoRS [71] | TableBank$_{word}$ | 96.7 | 98.5 |
| **DETR + $Q_b$ + $Q_p$ + $Q_n$** | **TableBank$_{word}$** | **98.3** | **99.1** |
| CascadeTabNet [40] | TableBank$_{both}$ | 94.4 | 95.7 |
| HybridTabNet [70] | TableBank$_{both}$ | 96.3 | 98.6 |
| CasTabDetectoRS [71] | TableBank$_{both}$ | 97.4 | 98.2 |
| **DETR + $Q_b$ + $Q_p$ + $Q_n$** | **TableBank$_{both}$** | **0.989** | **0.988** |

dataset in Table 3. we consider all PubLayNet classes (table, title, figure, list and text). We get the best results having mAP of 95.7% with noised anchor boxes taken as object queries. Anchor boxes are typically used in object detection models to help the model localize and classify objects in an image. However, in the context of graphical page object detection, anchor boxes used are too small or too large, they may not capture the full extent of an object on the page, leading to incorrect predictions. Noisy anchor boxes can help to capture more diverse object shapes and sizes, especially when the objects in the images have a wide variety of aspect ratios and sizes. By introducing random noise into the anchor boxes, the model is forced to learn more robust features and adapt to a wider range of object sizes and shapes, which can improve its overall performance.

The qualitative analysis for the PubLayNet dataset is illustrated in (a) part of Figure 6. Here, the performance (AP) using noised anchor boxes as object queries is highest on all IoU threshold values represented with a brown dotted line.

**Comparison with state-of-the-art methods.** We also compare the results of transformer with modifications in object queries with earlier approaches on the PubLayNet dataset. Table 4 shows that transformer with noised anchor boxes as object queries has outperformed the previous state-of-the-art methods for detecting graphical objects in document images.

**Table 3**: Comparison between DETR and its submodules on raw (without pre-processing), dilation and smudge transformation of the PubLayNet dataset. Here, term $Q_b$ represents object queries as anchor boxes, $Q_p$ denotes object queries with positive noise and $Q_n$ indicates object queries with negative noise. The IoU thresholds are set to 0.5 and 0.75. AR represents average precision for a large area. The best results are highlighted.

| Methods | Preprocessing | mAP | AP$^{50}$ | AP$^{75}$ | AR |
|---|---|---|---|---|---|
| DETR | raw | 93.6 | 95.3 | 94.1 | 94.5 |
| | Dilation+Smudge | 93.9±0.12 | 95.6 | 94.4 | 94.7 |
| DETR + $Q_b$ | raw | 94.2 | 95.8 | 94.6 | 94.9 |
| | Dilation+Smudge | 94.7±0.31 | 96.0 | 95.0 | 95.1 |
| DETR + $Q_b$ + $Q_p$ | raw | 94.8 | 96.3 | 95.5 | 95.3 |
| | Dilation+Smudge | 95.0±0.10 | 96.5 | 95.9 | 95.5 |
| DETR + $Q_b$ + $Q_p$ + $Q_n$ | **raw** | **95.2** | **96.9** | **96.2** | **95.7** |
| | **Dilation+Smudge** | 95.7±0.11 | 97.5 | 96.7 | 96.9 |

**Table 4**: Comparison between the transformer-based detectors and previous state-of-the-art results on PubLayNet validation set without pre-processing (raw data). The term $Q_b$ represents object queries as anchor boxes, $Q_p$ denotes object queries with positive noise and $Q_n$ indicates object queries with negative noise. Here, the mAP is for all these graphical objects. The best results are exhibited.

| Model | Framework | Backbone | Table | Text | Title | List | Figure | mAP |
|---|---|---|---|---|---|---|---|---|
| PubLayNet [17] | Mask R-CNN | ResNet-101 | 96.0 | 91.6 | 84.0 | 88.6 | 94.9 | 91.0 |
| UDoc [72] | Faster R-CNN | ResNet-50 | 97.3 | 93.9 | 88.5 | 93.7 | 96.4 | 93.9 |
| DiT$_B$ [73] | Cascade R-CNN | Transformer | 97.6 | 94.4 | 88.9 | 94.8 | 96.9 | 94.5 |
| LayoutLMv3$_B$ [74] | Cascade R-CNN | Transformer | 97.9 | 94.5 | 90.6 | 95.5 | 97.0 | 95.1 |
| **Our** | **DETR+$Q_b$+ $Q_p$+$Q_n$** | **ResNet-50** | 98.1 | 94.7 | 91.8 | 96.4 | 97.5 | 95.7 |

### 4.2.3  PubTables

PubTables is the largest table dataset on which we evaluate the capabilities of detection transformer and modifications in object queries. Table 5 shows the results of all modules on the raw (without pre-processing), dilation and smudge transformation of the PubTables dataset. It shows that the nosied anchor boxes as object queries shows the best results.

The qualitative analysis on the PubTables dataset is illustrated in (b) part of Figure 6. It provides the Average precision (AP) on all IoU thresholds. The improved query network, represented with a brown dotted line, gives the highest mAP on all IoU threshold values.

**Comparison with State-of-the-art Methods** Recently, Smock et al.[14] show an 82.5% mAP, 98.5% AP at 0.5 and 92.7% AP at 0.75 IoU threshold on the PuTables dataset on the Faster R-CNN detector as shown in Table 6.
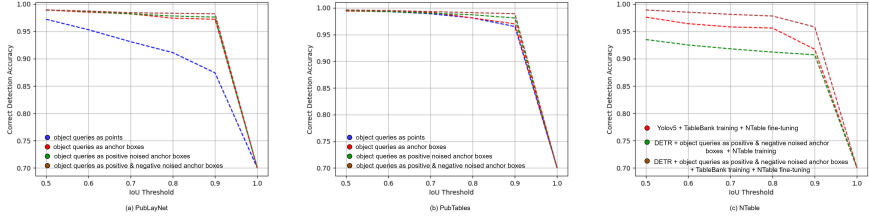
**Fig. 6**: Performance analysis of DETR with modifications in object queries in terms of AP over the IoU threshold values range from 0.5 to 1.0 on the a) PubLayNet, b) PubTables and c) NTable datasets. In part a) and b) of this figures, the blue color shows results with simple DETR, the red highlights results with anchor boxes as object queries, the green represents results with positive noised anchor boxes, and the brown denotes results with positive and negative noised anchor boxes. In part c), we compare results with previous YOLOv5 framework on NTable dataset for camera-based table images.

**Table 5**: Comparison between transformer-based detector results on raw (without pre-processing), dilation and smudge transformation of the PubTables dataset. Here, term $Q_b$ represents object queries as anchor boxes, $Q_p$ denotes object queries with positive noise and $Q_n$ indicates object queries with negative noise. The IoU thresholds are set to 0.5 and 0.75 for average precision and also calculate average recall for large objects. The best results are highlighted.

| Methods | Preprocessing | mAP | $AP^{50}$ | $AP^{75}$ | AR |
|---|---|---|---|---|---|
| DETR | raw | 97.6 | 99.5 | 98.9 | 98.5 |
| | Dilation + Smudge | 97.8±0.11 | 99.8 | 99.3 | 98.9 |
| DETR + $Q_b$ | raw | 98.0 | 99.4 | 99.2 | 99.2 |
| | Dilation + Smudge | 98.6±0.21 | 99.8 | 99.5 | 99.5 |
| DETR + $Q_b$ + $Q_p$ | raw | 98.1 | 99.5 | 99.0 | 99.2 |
| | Dilation + Smudge | 98.6±0.02 | 99.8 | 99.5 | 99.5 |
| **DETR + $Q_b$ + $Q_p$ + $Q_n$** | **raw** | **98.9** | **99.6** | **99.3** | **99.4** |
| | **Dilation + Smudge** | 99.3±0.42 | 99.9 | 99.8 | 99.6 |

As the PubTables dataset is newly released, no previous work on this dataset is available.

### 4.2.4 NTable

We evaluate modifications in object queries on camera-based table detection task using NTable dataset. In Table 7, we have AP value of 92.5% at IoU threshold of 0.6 trained on NTable-cam and NTable-gen data split.

We train the detection transformer network with noised anchors as object queries on TableBank$_{both}$ data split and then fine tune it on NTable-cam and NTable-gen dataset. It achieves AP value of 98.5% at IoU threshold of 0.6.

**Table 6**: Comparison between the transformer-based detectors and previous state-of-the-art results on PubTables dataset without pre-processing (raw data). Here, term $Q_b$ represents object queries as anchor boxes, $Q_p$ denotes object queries with positive noise and $Q_n$ indicates object queries with negative noise. The best results are exhibited.

| Model | Framework | mAP | $AP_{50}$ | $AP_{75}$ | AR |
|-------|-----------|-----|-----------|-----------|-----|
| Smock et al.[14] | Faster R-CNN | 82.5 | 98.5 | 92.7 | 86.6 |
| Smock et al.[14] | DETR | 96.6 | 99.5 | 98.8 | 98.1 |
| Our | DETR+ $Q_b + Q_p + Q_n$ | 98.9 | 99.6 | 99.3 | 99.4 |

**Table 7**: Comparison between the transformer-based detectors and previous state-of-the-art results on NTable dataset without pre-processing (raw data). Here, term $Q_b$ represents object queries as anchor boxes, $Q_p$ denotes object queries with positive noise and $Q_n$ indicates object queries with negative noise. The best results are exhibited.

| Model | Framework | Train | Fine-tuning | IoU | AP | AR | F-Score |
|-------|-----------|-------|-------------|-----|-----|-----|---------|
| NTable [18] | YOLOv5 | TableBank | Ntable-cam + Ntable-gen | 0.6 | 96.4 | 99.7 | 98.0 |
| | | | | 0.8 | 95.6 | 98.8 | 97.2 |
| Our | DETR | TableBank | Ntable-cam + Ntable-gen | 0.6 | 95.9 | 99.3 | 97.6 |
| | | | | 0.8 | 94.9 | 98.4 | 96.6 |
| Our | DETR + $Q_b$+ $Q_p + Q_n$ | NTable-cam+ NTable-gen | - | 0.6 | 92.5 | 98.2 | 95.3 |
| | | | | 0.8 | 91.2 | 97.3 | 94.2 |
| Our | DETR+$Q_b$+ $Q_p + Q_n$ | TableBank | NTable-cam + NTable-gen | 0.6 | 98.5 | 99.8 | 99.1 |
| | | | | 0.8 | 97.8 | 99.2 | 98.5 |

We also compare these results with earlier approach on NTable dataset.The qualitative analysis for the NTable dataset is illustrated in (c) part of Figure 6. Here, traditional anchor boxes use a fixed set of scales and aspect ratios chosen based on prior knowledge or heuristics. However, these fixed anchor boxes may not be optimal for datasets with a wide range of object sizes and shapes, such as the NTable dataset. Dynamic noised anchor boxes address this limitation by adjusting the anchor box scales and aspect ratios during training based on the dataset's distribution of object sizes and shapes.

## 4.3  Ablation Studies

In this section, we perform a series of ablation studies as pre-processing, quality and quantity of object queries on network performance. This ablation study is performed on PubTables dataset.

**Influence of pre-processing modules.** We study the effect of transformation approaches used in our method on performance. In Table 8, the best performance is achieved using the dilation transform and smudge transform together on the table data. The results are on PubTables dataset with positive and negative noised anchor boxes as object queries. These approaches make

the table regions more prominent, which increases the performance.

**Table 8**: Effectiveness of Pre-processing as smudge and dilation transformation on network performance.

| Dilation | Smudge | mAP | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|---|
| ✗ | ✗ | 98.9 | 99.6 | 99.3 |
| ✔ | ✗ | 99.0 | 99.7 | 99.5 |
| ✗ | ✔ | 99.2 | 99.8 | 99.7 |
| ✔ | ✔ | **99.3** | **99.9** | **99.8** |

**Table 9**: Effectiveness of the number of object queries on network performance.

| N | mAP | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|
| 5 | 97.0 | 98.7 | 97.0 |
| **10** | **98.9** | **99.6** | **99.3** |
| 50 | 98.7 | 99.5 | 99.2 |
| 100 | 98.6 | 99.4 | 99.0 |

**Influence of learnable object queries quantity** We study the effect of the quantity of learnable queries on detection transformer performmance. Table 9 shows different numbers of object queries and their effect on the performance. These results are on PubTables dataset with positive and negative noised anchor boxes as object queries without pre-processing on raw data. The best performance is achieved by setting the value of N to 10, lower or higher values will cause a significant performance drop. By setting a lower value of N, the model may not provide boxes to all objects and reduce performance by classifying some objects as false negatives. Similarly, by setting a large value of N, the model will overfit and classify no object regions as false positives that cause performance drop.

**Influence of object queries quality.** We analyze the effect object queries quality on the detection transformer performance. These object queries can be either points, anchor boxes, positive noised anchor boxes or positive and negative noised anchor boxes. In Table 10, We can observe that noised anchor boxes as object queries increase the performance as it gives better spatial prior for the attention module. Adding positive and negative noise to anchor boxes improves models' performance by introducing random perturbations during training to make the model more robust to variations in object size, aspect ratio, and position. Positive noise improves mAP to 98.1%. By com-

**Table 10**: Influence of modifying object queries as points, anchor boxes, positive noised anchor boxes and positive & negative noised anchor boxes as object queries on network performance. These are the results on Pubtables dataset on raw data without pre-processing.

| Points | Boxes | Positive noise | Negative noise | mAP | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|---|---|---|
| ✔ | ✗ | ✗ | ✗ | 97.0 | 99.5 | 98.9 |
| ✗ | ✔ | ✗ | ✗ | 98.0 | 99.4 | 99.2 |
| ✗ | ✔ | ✔ | ✗ | 98.1 | 99.5 | 99.0 |
| ✗ | ✔ | ✔ | ✔ | 98.9 | 99.6 | 99.3 |

bining positive and negative noise, the model becomes more robust to object appearance and position variations and can better differentiate between foreground and background regions, improving overall accuracy and robustness. Positive and negative noised anchors improve the mAP to 98.9%. This training produces more anchors by adding noise and selects the best ones closer to ground truth that improves performance.

# 5  Conclusion and Future Work

This paper bridges the performance gap between state-of-the-art CNN-based graphical object detection algorithms and detection transformers. We perform different experiments on detection transformers and observe their effectiveness and efficiency for the graphical object detection task. The transformer-based approach can achieve higher table detection accuracy without relying on handcrafted components like non-maximum suppression (NMS) and anchor design used in CNN-based object detectors. Furthermore, experimental results show that performance of detection tarnsformer depends on quality and quantity of object queries fed as input to transformer decoder. Consequently, these transformer-based detectors have achieved state-of-the-art performance on the four large public datasets, including TableBank, PubLayNet, PubTables and NTable dataset. In Future, we plan to extend the transformer-based model for table structure recognition and content extraction.
**Supplementary information** The supplementary material contains a detailed explanation of the dataset, evaluation setup and Implementation Details.

# Declarations

**Conflict of interest** The authors declare that they have no conflict of interest.

# References

[1] Reis, J., Amorim, M., Melão, N., Matos, P.: Digital transformation: A literature review and guidelines for future research. In: Rocha, Á., Adeli, H., Reis, L.P., Costanzo, S. (eds.) Trends and Advances in Information Systems and Technologies, pp. 411–421. Springer, Cham (2018)

[2] Zhao, Z., Jiang, M., Guo, S., Wang, Z., Chao, F., Tan, K.C.: Improving deep learning based optical character recognition via neural architecture search. In: 2020 IEEE Congress on Evolutionary Computation (CEC), pp. 1–7 (2020). https://doi.org/10.1109/CEC48606.2020.9185798

[3] Hashmi, K.A., Bymana Ponnappa, R., Bukhari, S.S., Jenckel, M., Dengel, A.: Feedback learning: Automating the process of correcting and completing the extracted information. In: 2019 International Conference on Document Analysis and Recognition Workshops (ICDARW), vol. 5, pp. 116–121 (2019). https://doi.org/10.1109/ICDARW.2019.40091

[4] van Strien, D.A., Beelen, K., Ardanuy, M.C., Hosseini, K., McGillivray, B., Colavizza, G.: Assessing the impact of ocr quality on downstream nlp tasks. In: ICAART (2020)

[5] Anh, T.T., In-Seop, N., Soo-Hyung, K.: A hybrid method for table detection from document image. In: 2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR), pp. 131–135 (2015). https://doi.org/10.1109/ACPR.2015.7486480

[6] Gatos, B., Danatsas, D., Pratikakis, I., Perantonis, S.J.: Automatic table detection in document images. In: ICAPR (2005)

[7] Harit, G., Bansal, A.: Table detection in document images using header and trailer patterns. In: Proceedings of the Eighth Indian Conference on Computer Vision, Graphics and Image Processing. ICVGIP '12. Association for Computing Machinery, New York, NY, USA (2012). https://doi.org/10.1145/2425333.2425395. https://doi.org/10.1145/2425333.2425395

[8] Wangt, Y., Phillipst, I.T., Haralick, R.: Automatic table ground truth generation and a background-analysis-based table structure extraction method. In: Proceedings of Sixth International Conference on Document Analysis and Recognition, pp. 528–532 (2001). https://doi.org/10.1109/ICDAR.2001.953845

[9] Yi, X., Gao, L., Liao, Y., Zhang, X., Liu, R., Jiang, Z.: Cnn based page object detection in document images. In: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR), vol. 01, pp. 230–235 (2017). https://doi.org/10.1109/ICDAR.2017.46

[10] Borges Oliveira, D.A., Viana, M.P.: Fast cnn-based document layout analysis. In: 2017 IEEE International Conference on Computer Vision Workshops (ICCVW), pp. 1173–1180 (2017). https://doi.org/10.1109/ICCVW.2017.142

[11] Sun, N., Zhu, Y., Hu, X.: Faster r-cnn based table detection combining corner locating. 2019 International Conference on Document Analysis and Recognition (ICDAR), 1314–1319 (2019)

[12] Agarwal, M., Mondal, A., Jawahar, C.V.: Cdec-net: Composite deformable cascade network for table detection in document images. In: 2020 25th International Conference on Pattern Recognition (ICPR), pp.

9491–9498 (2021). https://doi.org/10.1109/ICPR48806.2021.9411922

[13] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems, vol. 30. Curran Associates, Inc., ??? (2017). https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

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

[15] Shehzadi, T., Hashmi, K.A., Stricker, D., Afzal, M.Z.: 2D Object Detection with Transformers: A Review (2023)

[16] Li, M., Cui, L., Huang, S., Wei, F., Zhou, M., Li, Z.: TableBank: A Benchmark Dataset for Table Detection and Recognition (2019)

[17] Zhong, X., Tang, J., Yepes, A.J.: Publaynet: largest dataset ever for document layout analysis. In: 2019 International Conference on Document Analysis and Recognition (ICDAR), pp. 1015–1022 (2019). https://doi.org/10.1109/ICDAR.2019.00166. IEEE

[18] Zhu, Z., Gao, L., Li, Y., Huang, Y., Du, L., Lu, N., Wang, X.: Ntable: A dataset for camera-based table detection. In: Document Analysis and Recognition – ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part II, pp. 117–129. Springer, Berlin, Heidelberg (2021). https://doi.org/10.1007/978-3-030-86331-9_8. https://doi.org/10.1007/978-3-030-86331-9_8

[19] Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: towards real-time object detection with region proposal networks. CoRR **abs/1506.01497** (2015) https://arxiv.org/abs/1506.01497

[20] Girshick, R.B.: Fast R-CNN. CoRR **abs/1504.08083** (2015) https://arxiv.org/abs/1504.08083

[21] He, K., Gkioxari, G., Dollár, P., Girshick, R.B.: Mask R-CNN. CoRR **abs/1703.06870** (2017) https://arxiv.org/abs/1703.06870

[22] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S.E., Fu, C., Berg, A.C.: SSD: single shot multibox detector. CoRR **abs/1512.02325** (2015) https://arxiv.org/abs/1512.02325

[23] Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense

object detection. CoRR **abs/1708.02002** (2017) https://arxiv.org/abs/1708.02002

[24] Ge, Z., Liu, S., Wang, F., Li, Z., Sun, J.: YOLOX: exceeding YOLO series in 2021. CoRR **abs/2107.08430** (2021) https://arxiv.org/abs/2107.08430

[25] Kieninger, T., Dengel, A.: The t-recs table recognition and analysis system. In: Selected Papers from the Third IAPR Workshop on Document Analysis Systems: Theory and Practice. DAS '98, pp. 255–269. Springer, Berlin, Heidelberg (1998)

[26] Cesarini, F., Marinai, S., Sarti, L., Soda, G.: Trainable table location in document images. In: 2002 International Conference on Pattern Recognition, vol. 3, pp. 236–2403 (2002). https://doi.org/10.1109/ICPR.2002.1047838

[27] Silva, A.C.e.: Learning rich hidden markov models in document analysis: Table location. In: 2009 10th International Conference on Document Analysis and Recognition, pp. 843–847 (2009). https://doi.org/10.1109/ICDAR.2009.185

[28] Costa e Silva, A., Jorge, A., Torgo, L.: Design of an end-to-end method to extract information from tables. Document Analysis and Recognition **8**, 144–171 (2006). https://doi.org/10.1007/s10032-005-0001-x

[29] Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3431–3440 (2015). https://doi.org/10.1109/CVPR.2015.7298965

[30] Yang, X., Yümer, M.E., Asente, P., Kraley, M., Kifer, D., Giles, C.L.: Learning to extract semantic structure from documents using multimodal fully convolutional neural network. CoRR **abs/1706.02337** (2017) https://arxiv.org/abs/1706.02337

[31] He, D., Cohen, S., Price, B., Kifer, D., Giles, C.L.: 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, pp. 254–261 (2017). https://doi.org/10.1109/ICDAR.2017.50

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

[33] Li, X.-H., Yin, F., Liu, C.-L.: Page object detection from pdf document images by deep structured prediction and supervised clustering. In: 2018 24th International Conference on Pattern Recognition (ICPR), pp. 3627–3632 (2018). https://doi.org/10.1109/ICPR.2018.8546073

[34] Riba, P., Dutta, A., Goldmann, L., Fornés, A., Ramos, O., Lladós, J.: Table detection in invoice documents by graph neural networks. In: 2019 International Conference on Document Analysis and Recognition (ICDAR), pp. 122–127 (2019). https://doi.org/10.1109/ICDAR.2019.00028

[35] Holecek, M., Hoskovec, A., Baudis, P., Klinger, P.: Line-items and table understanding in structured documents. CoRR **abs/1904.12577** (2019) https://arxiv.org/abs/1904.12577

[36] Li, M., Xu, Y., Cui, L., Huang, S., Wei, F., Li, Z., Zhou, M.: Docbank: A benchmark dataset for document layout analysis. CoRR **abs/2006.01038** (2020) https://arxiv.org/abs/2006.01038

[37] Girshick, R.B., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. CoRR **abs/1311.2524** (2013) https://arxiv.org/abs/1311.2524

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

[39] Saha, R., Mondal, A., Jawahar, C.V.: Graphical object detection in document images. CoRR **abs/2008.10843** (2020) https://arxiv.org/abs/2008.10843

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

[41] Redmon, J., Divvala, S.K., Girshick, R.B., Farhadi, A.: You only look once: Unified, real-time object detection. CoRR **abs/1506.02640** (2015) https://arxiv.org/abs/1506.02640

[42] Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense object detection. CoRR **abs/1708.02002** (2017) https://arxiv.org/abs/1708.02002

[43] Cai, Z., Vasconcelos, N.: Cascade R-CNN: high quality object detection and instance segmentation. CoRR **abs/1906.09756** (2019) https://arxiv.

org/abs/1906.09756

[44] Arif, S., Shafait, F.: Table detection in document images using foreground and background features. In: 2018 Digital Image Computing: Techniques and Applications (DICTA), pp. 1–8 (2018). https://doi.org/10.1109/DICTA.2018.8615795

[45] Gilani, A., Qasim, S.R., Malik, I., Shafait, F.: Table detection using deep learning. In: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR), vol. 01, pp. 771–776 (2017). https://doi.org/10.1109/ICDAR.2017.131

[46] Siddiqui, S.A., Malik, M.I., Agne, S., Dengel, A., Ahmed, S.: Decnt: Deep deformable cnn for table detection. IEEE Access **6**, 74151–74161 (2018). https://doi.org/10.1109/ACCESS.2018.2880211

[47] Gao, L., Huang, Y., Déjean, H., Meunier, J.-L., Yan, Q., Fang, Y., Kleber, F., Lang, E.: Icdar 2019 competition on table detection and recognition (ctdar). In: 2019 International Conference on Document Analysis and Recognition (ICDAR), pp. 1510–1515 (2019). IEEE

[48] Mondal, A., Lipps, P., Jawahar, C.V.: IIIT-AR-13K: A new dataset for graphical object detection in documents. CoRR **abs/2008.02569** (2020) https://arxiv.org/abs/2008.02569

[49] Göbel, M.C., Hassan, T., Oro, E., Orsi, G.: Icdar 2013 table competition. 2013 12th International Conference on Document Analysis and Recognition, 1449–1453 (2013)

[50] Gao, L., Yi, X., Jiang, Z., Hao, L., Tang, Z.: Icdar2017 competition on page object detection. In: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR), vol. 01, pp. 1417–1422 (2017). https://doi.org/10.1109/ICDAR.2017.231

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

[52] Shehzadi, T., Hashmi, K.A., Pagani, A., Liwicki, M., Stricker, D., Afzal, M.Z.: Mask-aware semi-supervised object detection in floor plans. Applied Sciences **12**(19) (2022). https://doi.org/10.3390/app12199398

[53] Li, Y., Li, Q., Meng, S., Hou, J.: Transformer-based rating-aware sequential recommendation. In: Algorithms and Architectures for Parallel Processing: 21st International Conference, ICA3PP 2021, Virtual Event, December 3–5, 2021, Proceedings, Part I, pp. 759–774. Springer, Berlin,

Heidelberg (2021). https://doi.org/10.1007/978-3-030-95384-3_47. https://doi.org/10.1007/978-3-030-95384-3_47

[54] Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: European Conference on Computer Vision, pp. 213–229 (2020). Springer

[55] Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense object detection. CoRR **abs/1708.02002** (2017) https://arxiv.org/abs/1708.02002

[56] Shehzadi, T., Hashmi, K.A., Stricker, D., Liwicki, M., Afzal, M.Z.: Towards End-to-End Semi-Supervised Table Detection with Deformable Transformer (2023)

[57] He, K., Zhang, X., Ren, S., Sun, J.: Spatial pyramid pooling in deep convolutional networks for visual recognition. CoRR **abs/1406.4729** (2014) https://arxiv.org/abs/1406.4729

[58] Meng, D., Chen, X., Fan, Z., Zeng, G., Li, H., Yuan, Y., Sun, L., Wang, J.: Conditional DETR for fast training convergence. CoRR **abs/2108.06152** (2021) https://arxiv.org/abs/2108.06152

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

[60] Sun, Z., Cao, S., Yang, Y., Kitani, K.: Rethinking transformer-based set prediction for object detection. CoRR **abs/2011.10881** (2020) https://arxiv.org/abs/2011.10881

[61] Wang, Y., Zhang, X., Yang, T., Sun, J.: Anchor : query design for transformer-based detector. CoRR **abs/2109.07107** (2021) https://arxiv.org/abs/2109.07107

[62] Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable DETR: deformable transformers for end-to-end object detection. CoRR **abs/2010.04159** (2020) https://arxiv.org/abs/2010.04159

[63] Dai, X., Chen, Y., Yang, J., Zhang, P., Yuan, L., Zhang, L.: Dynamic detr: End-to-end object detection with dynamic attention. 2021 IEEE/CVF International Conference on Computer Vision (ICCV), 2968–2977 (2021)

[64] Yao, Z., Ai, J., Li, B., Zhang, C.: Efficient DETR: improving end-to-end object detector with dense prior. CoRR **abs/2104.01318** (2021) https://arxiv.org/abs/2104.01318

[65] Li, F., Zhang, H., Liu, S., Guo, J., Ni, L.M., Zhang, L.: Dn-detr: Accelerate detr training by introducing query denoising. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13619–13627 (2022)

[66] Zhang, H., Li, F., Liu, S., Zhang, L., Su, H., Zhu, J., Ni, L.M., Shum, H.-Y.: DINO: DETR with Improved DeNoising Anchor Boxes for End-to-End Object Detection. arXiv (2022). https://doi.org/10.48550/ARXIV.2203.03605. https://arxiv.org/abs/2203.03605

[67] Szegedy, C., Ioffe, S., Vanhoucke, V.: Inception-v4, inception-resnet and the impact of residual connections on learning. CoRR **abs/1602.07261** (2016) https://arxiv.org/abs/1602.07261

[68] Girshick, R.B.: Fast R-CNN. CoRR **abs/1504.08083** (2015) https://arxiv.org/abs/1504.08083

[69] Lin, T., Maire, M., Belongie, S.J., Bourdev, L.D., Girshick, R.B., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: common objects in context. CoRR **abs/1405.0312** (2014) https://arxiv.org/abs/1405.0312

[70] Nazir, D., Hashmi, K.A., Pagani, A., Liwicki, M., Stricker, D., Afzal, M.Z.: Hybridtabnet: Towards better table detection in scanned document images. Applied Sciences **11**(18) (2021). https://doi.org/10.3390/app11188396

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

[72] Gu, J., Kuen, J., Morariu, V.I., Zhao, H., Barmpalios, N., Jain, R., Nenkova, A., Sun, T.: Unified Pretraining Framework for Document Understanding. arXiv (2022). https://doi.org/10.48550/ARXIV.2204.10939. https://arxiv.org/abs/2204.10939

[73] Li, J., Xu, Y., Lv, T., Cui, L., Zhang, C., Wei, F.: DiT: Self-supervised Pre-training for Document Image Transformer. arXiv (2022). https://doi.org/10.48550/ARXIV.2203.02378. https://arxiv.org/abs/2203.02378

[74] Huang, Y., Lv, T., Cui, L., Lu, Y., Wei, F.: LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking. arXiv (2022). https://doi.org/10.48550/ARXIV.2204.08387. https://arxiv.org/abs/2204.08387