

GraphKD: Exploring Knowledge Distillation Towards Document Object Detection with Structured Graph Creation

Ayan Banerjee¹[0000-0002-0269-2202], Sanket Biswas¹[0000-0001-6648-8270], Josep Lladós¹[0000-0002-4533-4739], and Umapada Pal²[0000-0002-5426-2618]

¹ Computer Vision Center & Computer Science Department
Universitat Autònoma de Barcelona, Spain
{abanerjee, sbiswas, josep}@cvc.uab.es
² CVPR Unit, Indian Statistical Institute, Kolkata, India
umapada@isical.ac.in

Abstract. Object detection in documents is a key step to automate the structural elements identification process in a digital or scanned document through understanding the hierarchical structure and relationships between different elements. Large and complex models, while achieving high accuracy, can be computationally expensive and memory-intensive, making them impractical for deployment on resource constrained devices. Knowledge distillation allows us to create small and more efficient models that retain much of the performance of their larger counterparts. Here we present a graph-based knowledge distillation framework to correctly identify and localize the document objects in a document image. Here, we design a structured graph with nodes containing proposal-level features and edges representing the relationship between the different proposal regions. Also, to reduce text bias an adaptive node sampling strategy is designed to prune the weight distribution and put more weightage on non-text nodes. We encode the complete graph as a knowledge representation and transfer it from the teacher to the student through the proposed distillation loss by effectively capturing both local and global information concurrently. Extensive experimentation on competitive benchmarks demonstrates that the proposed framework outperforms the current state-of-the-art approaches. The code will be available at: github.com/ayanban011/GraphKD

Keywords: Knowledge Distillation · Document Object Detection · Graph Neural Network.

1 Introduction

Document Object Detection (DOD) has indeed become an essential task in Document Understanding (DU) because any task related to document understanding entails the need to obtain a structured representation that helps to distinguish between text, images, tables, headings, footers, and other design elements [46,51].

The main goal of DOD is to understand the structure and content of a document, typically as a precursor to more detailed processing or analysis. For instance, in the field of Optical Character Recognition (OCR) [27], DOD can be used to determine which parts of a page contain text that needs to be recognized, versus which parts might contain pictures or other non-textual elements. Similarly, in Document Retrieval [20], Key Information Extraction [59], and Visual Question Answering [65], DOD helps to localize the region where the key information or the answer has been located. For this reason, throughout the last decade, remarkable progress has been observed, starting from the convolution based algorithms (e.g. Faster-RCNN [60], Mask R-CNN [76], RetinaNet [42], and so on) to the large-scale multi-modal transformers [32,61]. With the increasing complexity of document layouts, the model complexity and the number of model parameters have also increased. Although the conventional approaches [7,37] demonstrate significant performance, we cannot use them in our edge devices (we are penalized by the computation cost). On the other hand, tiny networks can be used for edge devices but we are penalized by the object detection performance (i.e. efficiency). In order to solve this trade-off between memory and efficiency we proposed a graph-based knowledge distillation technique where we trained large networks in the GPUs and encoded their learned features in a graph data structure to transfer it to the tiny networks by optimizing a distillation loss, so that we can use the same in the edge devices. According to our best knowledge, this is the first attempt to perform knowledge distillation (KD) in order to solve the Document Object Detection (DOD) task.

However, KD in object detection is not easy because we need to consider the spatial location of the multiple objects along with their scale variation and imbalanced distributions. Conventional KD techniques [29,28,43] are often fail to deal with the feature imbalance problem and also unable to identify the missing instance-level relations as they perform the distillation through one of the following techniques. 1) *Logit-based*: here we distill only the logits to the final softmax layer. This technique always loses the fine-grained information as it doesn't consider the feature maps of the intermediate layers. 2) *Feature-based*: It performs layer-to-layer feature distillation but suffers with feature alignment problems. So, only homogeneous distillation (e.g. ResNet101-ResNet50) is possible. 3) *Hybrid*: It performs a layer-to-layer distillation of the feature maps as well as logit-to-logit distillation and the final softmax layer but it diminishes the transferability. In order to tackle these issues, we construct a structured instance graph where we collect the regional features from the RoI instances of the RPN and store them in the graph nodes. On the other hand, the edges represent the relations between nodes and are measured by their feature similarity. So, the nodes help to overcome the feature imbalance problem and edges excavate the missing instance relation. Not only that graphs also preserved the structural information during the embedding transfer and adapted easily to the topological structure between teacher and student even with different widths and depths enabling us to perform heterogeneous distillation (EfficientNet to ResNet, MobileNet to ResNet, and so on).

Even the distillation through the graph for the documents creates some additional issues. Initially, a notable challenge arises due to the large number of text nodes providing text bias which is propagated as a noise for non-text nodes (i.e. Table, Figure, etc.) distillation. So, we cannot use simple L2 distance [17] in order to compute the similarity between nodes and edges during distillation. Hence we used cosine similarity to compute the similarity between the teacher and student nodes as it considers the orthogonality of the representation and Mahalanobis distance for the edges which is sensitive to the outliers. Also, we propose a sample mining approach for the "Text" nodes to strategically remove less important text-associated edges. This not only cuts down the biased edges but also improves the regularization of false negatives in distillation framework.

The overall contributions of this work can be summarized as follows:

- A new task knowledge distillation in document object detection has been proposed and solved through a structure instance graph creation whose node contains the RoI instances of the RPN and the edges represent the relationship between the nodes through feature similarity. To the best of our knowledge, this is the first time where KD for DOD has been explored.
- A new sample mining technique has been introduced to get rid of the text bias and improve the regularization of false negatives during distillation.
- Last but not the least, we utilize cosine similarity for the node-to-node distillation to tackle the orthogonality and Mahalanobis distance for edge-to-edge distillation to tackle the outliers which allows us to perform heterogeneous distillation which is one of the most critical problems in distillation till date.

The rest of the paper is organized in the following way: In Section 2 we review state-of-the-art approaches (SOTA) for DOD and KD. We describe the *GraphKD* in Section 3. We introduce our experimental evaluation as well as ablation studies in Section 4 and discuss the extensive experimentation to consolidate our claims. Finally, Section 5 concludes and guides the future research directions.

2 Related Work

Mainstream DOD algorithms were initially dominated by the classical heuristic rule-based algorithms before the rise of deep learning. Then, Convolution frameworks took the lead until transformers-based architectures demonstrated remarkable performance. This section provides a brief overview of SOTA methodologies for DOD till date. As there is no work on KD for DOD, we are trying to introduce a generalized SOTA on KD for object detection.

Heuristic Rule-based DOD: This refers to a method of identifying and extracting structured information from documents using a set of predefined rules and heuristics. Heuristics methods can be further classified into top-down, bottom-up, and hybrid approaches based on the parsing directions. Top-down strategies [33,35] perform iteratively partitioning a document image into regions

until a distinct region is identified. This strategy offers quicker implementation but sacrifices generalization and is effective primarily on specific document types. On the other hand, Bottom-up approaches [4,54] perform pixel grouping, merging, and other set operations to create homogeneous regions around similar objects, separating them from dissimilar ones. Although bottom-up approaches can tackle complex layouts, they are computationally expensive, especially for large or high-resolution documents. Moreover, to take advantage of both, hybrid strategies [13,26] leverage both bottom-up and top-down cues for fast and efficient results. While effective for table detection in the pre-deep learning era, these methods fall short for other complex categories.

Convolution-based DOD: the use of convolutional neural networks (CNNs) leverages the power of convolutional operations to learn hierarchical features and spatial relationships in the document images, enabling the network to effectively distinguish between different types of document components. In 2015, Faster-RCNN [60] provides a strong baseline for table detection and further extended to solve page segmentation [39]. With an extra segmentation loss and additional branch on top of Faster-RCNN for mask prediction, Mask-RCNN [39] provides the instance segmentation benchmark for newspaper elements. Similarly, RetinaNet [42] provides a complex convolution benchmark for keyword detection in document images, specifically focusing on text region detection. Meanwhile, DeepDeSRT [56] introduces a novel image transformation strategy to set a new benchmark for table detection and structure recognition. It identifies visual features of table structures and feeds them into a fully convolutional network with skip pooling. A similar FCNN-based framework [50,53] has been utilized for historical documents which surpasses the previous convolutional auto-encoder benchmark with transfer learning paradigms [15,14,55] on ICDAR2017 Page Object Detection (POD) dataset. A new cross-domain DOD benchmark was established in [38] to apply domain adaptation strategies to solve the domain shift problem with DeepLabv3+ [45] and YOLO [24]. However, the problem isn't fully solved until the arrival of transformers due to the lack of cross-attention mechanism. Lastly, a new benchmark [67] for vision-based layout detection utilized a recurrent convolutional neural network with VoVNet-v2 backbone, generating synthetic PDF documents from ICDAR-2013 and GROTOAP datasets, achieving a new benchmark for scientific document segmentation.

Transformer based DOD: Nowadays, modern transformers excel in the performance of DOD by using positional embedding and self-attention mechanisms [41]. Here, DiT [37] set a new baseline for DOD, using self-supervised pretraining on large-scale unlabeled document images, but its applicability to small magazine datasets like PRIMA is limited. To enhance performance, TILT [52] introduces a mechanism that concurrently learns textual semantics, visual features, and layout information using an encoder-decoder Transformer. A similar auto-encoder mechanism [66] establishes a new baseline for the PubLayNet dataset (AP: 95.95) by extracting text information through OCR. Recently, LayoutLmv3 [32] achieved state-of-the-art results in visual document understanding tasks

through joint learning of text, layout, and visual features. While excelling for large-scale datasets, it also falls short for small-scale datasets like DiT. Other recent approaches [21,76,61,25,72] leverage joint pretraining for various VDU tasks, including document visual question answering. While beneficial for downstream tasks through unified pretraining, they exhibit a pretraining bias, hindering performance in domain shifts and struggling to learn class information with low instance numbers due to a lack of weight prioritization. SwinDocSegmenter [7] tries to solve this problem through contrastive denoising training and hybrid bipartite matching, however, it is computationally expensive (223M model parameters to train). SemiDocSeg [6] tries to improve the training strategy by training only the categories that occur rarely in the document and letting the other categories learn through the co-occurrence matrix and support set. It improves the training time of the system and reduces the annotation time but it still has 223M parameters to train. This highlights the necessity of a KD technique for optimizing the number of trainable parameters.

Knowledge Distillation: According to our best knowledge, there is no such work on knowledge distillation for document object detection. So, we are striving to achieve the closest state-of-the-art results available for knowledge distillation in document object detection. KD strategies can be categorized into three main categories: *response-based* KD [1,5,31,47,69,73] seeks to match the final layer predictions of the teacher model; *feature-based* KD [2,11,12,16,30,36] aims to mimic features extracted from intermediate hidden layers of the deep network and *relation-based* KD [9,34,70,22] which exploits the relations between different layers or sampled data points. However, the latter approach is more geared toward pixel-based semantic segmentation tasks. While feature-based KD is more versatile, it is more expensive and harder to implement than soft teacher predictions. While offline methods [23,40] consider an existing frozen teacher model, online methods [71,74] update both student and teacher networks jointly. Self-distillation [63,64] represents a special case of online KD, which employs the same network as both the teacher and student, progressively outperforming the network’s performance, albeit disregarding the aim of efficiency.

However, all the methods described above, perform layerwise distillation so the two backbones should follow the same architecture (i.e. CNN to CNN, ViT to ViT). There is no chance of heterogeneous distillation (i.e. ViT to CNN or vice-versa). This problem motivates us to propose graph-based knowledge distillation (GraphKD) which gives us the freedom about the choice of backbone. In fact, we can distill an entire teacher model to a simple student backbone, as it performs node-to-node as well as edge-to-edge distillation via a distillation loss.

3 Method

We introduce a distillation framework in Fig. 1 where we are trying to generate a structured instance graph constructed using regional objects within both the teacher and student backbone. This graph effectively leverages the profound insights embedded within detection networks, essentially representing a knowledge

framework integrated into the detection system. Distilling this structured graph not only facilitates thorough knowledge transfer but also preserves the entirety of the topological structure corresponding to its embedding space. It provides

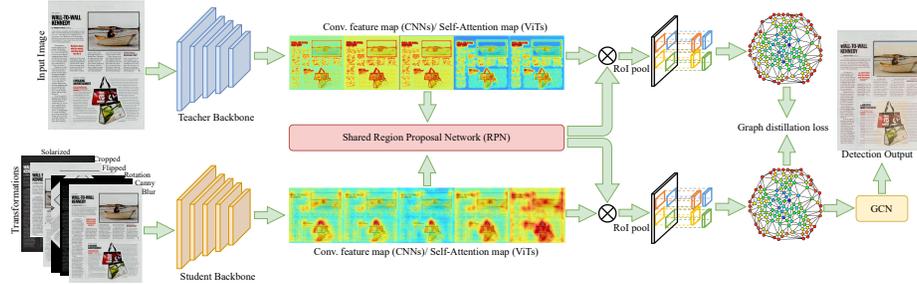


Fig. 1: **GraphKD** creates a graph from RoI pooled features of both teacher and student networks and utilize graph distillation loss from knowledge transfer. Finally, a Graph Convolution Network has been used to predict the object classes.

us the flexibility to perform knowledge distillation from any large network to a smaller network (e.g. ResNet18) without depending on its architectural similarity as we are performing a node-to-node distillation between two structured graph instances which is a main drawback of the SOTA distillation strategy.

3.1 Structured Instance Graph Construction

Within the structured graph, each node is a representation of an individual instance within an image, and this representation is in the form of a vectorized feature associated with that specific instance. Also, the connection between two instances is established as the edge connecting their respective nodes, and it is determined by measuring their similarity within the embedding space. Indeed, it's essential to highlight that the definition and meaning of an edge in this context are fundamentally distinct from the concept of pixel similarity as outlined in [17,44]. It is notable that, unlike the other approaches that handle the entire backbone feature map, our focus is directed toward constructing graphs using Region of Interest (RoI) pooled features (See 2nd and 3rd diagram of Fig. 2). These RoI features are derived from the RPN proposals and are then transmitted to the subsequent detection head for further processing. Our nodes are obtained through pooling and extraction using adaptable semantic proposals of different scales and sizes, thereby fostering robust semantic connections among them while those [17,44] re-sampled and uniformly distributed pixel blocks with the same sizes within an image. The robust connections established between instances that are transferred from the teacher to the student play a pivotal role in enabling interpretable distillation within our approach. Along with that, this framework shared the teacher and student ROIs to make a perfect alignment of the anchor boxes so the same regional features get highlighted, and extracted by two different backbones (teacher and student).

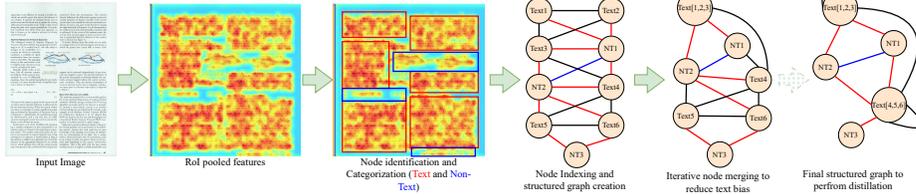


Fig. 2: **Structured graph creation:** Here first we extracted the RoI pooled features and classified them into "Text" and "Non-text" based on their covariance. Then we initialize the node in the identified RoI regions and define the adjacency edges. Lastly, we iteratively merge the text node with an adaptive sample mining strategy to reduce text bias.

3.2 Node Definition

We create nodes directly from RoI pooled features and categorize them into "non-text" or "text" based on the covariance of the RoI pooled features. As depicted in the Fig. 2 (2nd diagram), the "text" region contains high feature covariance compared to the "non-text" nodes. However, as depicted in the Fig. 2, there are 6 "text" nodes and 3 "non-text" nodes, which leads to a large bias towards the text region in the final object detection performance (i.e. lots of regions will be misclassified as text region). In order to reduce this text bias, we concatenate the adjacency text nodes whose edge distance is below a certain threshold (i.e. node merging and edge reduction). In the late stage, when we need to perform the node separation to detect individual text regions we utilize an adaptive text weight loss, which performs a weight base separation until maximum IoU with the ground truth is achieved.

3.3 Edge Definition

The edges of the structured graph \mathcal{G}_s as $\xi_s = [e_{ij}]_{n \times n}$ where n is the size of the node feature set. e_{pq} is the edge between p-th and q-th nodes, denoting the cosine similarity of the corresponding instances in the graph embedding space as defined in Eq. (1).

$$e_{pq} = \frac{v_p \cdot v_q}{\|v_p\| \cdot \|v_q\|} \quad (1)$$

Where, v_i is the node of the i^{th} instance. This definition of the edges is invariable to the length of the feature belonging to the corresponding node. This helps the framework to get rid of backbone similarity. The graph we are creating is always a complete, symmetric, and undirected graph (i.e. $e_{pq} = e_{qp}$) with elements all being 1 in the principal diagonal. We utilize this cosine similarity to create the weighted edges because it helps to put small weights between the edges of the similar nodes (i.e. text-text, non-text(NT)-non-text(NT)) and higher weights between the edges of the different nodes (text-NT). This helps a lot in thresholding-based node merging and edge classification in the later stage.

3.4 Node Indexing

We’ve found that using dense edges generated by the entire set of nodes can hinder the training process (See 4th diagram of Fig. 2). This is because a significant number of ”Text” nodes contribute to an excessive loss during the distillation of object-related edges. To counteract this issue, it’s beneficial to create a more focused set of edges using only the non-text nodes.

However, eliminating all edges related to ”Text” nodes can lead to a significant loss of information early in the training process. This is because some of these edges represent hard negative samples that provide valuable insights during training. To address this, we’ve introduced a technique called Object Samples Mining [48]. This method identifies and selects pertinent non-text nodes, which are then combined with all the merged ”Text” nodes to form edges.

Algorithm 1 Node Indexing

Require: $V_T^{text}, V_S^{text}, T$ $\triangleright T$: Teacher model; S : Student Model;
Ensure: V_T^{mine}, V_S^{mine} $\triangleright mine$: Nodes after mining
Initialize threshold t .
 $\mathcal{L}_{RoI_{T}^{text}} \leftarrow T.RoIclsLoss(V_T^{text})$
 $V_{ind} \leftarrow V_T^{text} : \forall V_T^{text} \models L_{RoI_{T}^{text}} \gg t$
 $V_T^{mine} \leftarrow SelectIndex(V_T^{bg}, V_{text})$
 $V_S^{mine} \leftarrow SelectIndex(V_S^{bg}, V_{text})$
Return V_T^{mine}, V_S^{mine}

Algorithm 1 presents a method that merges specific ”Text” samples based on a criterion: their classification losses in the teacher model exceed a threshold t (empirically decided). This suggests that these ”Text” samples are more likely to be misclassified (i.e. leads to text bias). Consequently, they can be appropriately included in the set of edges that only contain non-text nodes. This inclusion ensures that a dense graph is maintained.

In this Algorithm 1, samples with high confidence to be classified to the ”Text” are not directly added to the set. Here, we perform node indexing through the *SelectIndex* function as defined in [58]. After this sample mining, we applied graph distillation loss to perform KD from the teacher-to-student model.

3.5 Graph Distillation Loss

The graph distillation loss \mathcal{L}_g is defined as the variance between structured graphs of teacher $\mathcal{G}_{s,t}$ and student $\mathcal{G}_{s,s}$, consisting of graph node and edge loss (\mathcal{L}_v , and \mathcal{L}_ξ) respectively. We simply compute the Mahalanobis distance function to evaluate this loss as depicted in Eq. (2).

$$\mathcal{L}_g = \lambda_1 \cdot \mathcal{L}_v^{text} + \lambda_2 \cdot \mathcal{L}_v^{nt} + \lambda_3 \cdot \mathcal{L}_\xi = \frac{\lambda_1}{N_{nt}} \sum_{i=1}^{N_{nt}} \left\| \frac{v_i^{t,nt} - v_i^{s,nt}}{\sigma_{nt}} \right\| + \frac{\lambda_2}{N_{text}} \sum_{i=1}^{N_{text}} \left\| \frac{v_i^{t,text} - v_i^{s,text}}{\sigma_{text}} \right\| + \frac{\lambda_3}{N^2} \sum_{i=1}^N \sum_{j=1}^N \left\| \frac{e_{ij}^t - e_{ij}^s}{\sigma_\xi} \right\| \quad (2)$$

where, λ_i represents the penalty coefficient in order to balance the text, non-text, and edge components of graph distillation loss. We set λ_1 and λ_3 to 0.3 based on the Optuna [3] search, and $\lambda_2 = \alpha \cdot \frac{N_{nt}}{N_{text}}$ as an adaptive loss weight for text nodes to mitigate the imbalance problem, where α is empirically determined to ensure that the loss is on a similar scale as other distillation losses.

The graph node loss \mathcal{L}_v is the imitation loss [49] between node set which makes an accurate alignment between teacher and student backbone extracted features through instance matching (See Algorithm 2). Conventionally, in KD, there’s a preference for matching the feature map directly between two networks. But when it comes to detection models, not every pixel in the feature maps contributes to the classification and box regression loss calculation. Instead of using the entire feature map, we use selected features from both text and non-text regions to compute the graph node loss. This method helps the student model to pay attention to the RoIs and assimilate relevant object knowledge.

Algorithm 2 Graph Distillation Loss

Require: Image x , transformations x_t , teacher T , student S

Ensure: \mathcal{L}_g

```

 $\mathcal{F}_T \leftarrow T.\text{backbone}(x)$ 
 $\mathcal{F}_s \leftarrow S.\text{backbone}(x_t)$ 
proposals  $\leftarrow S.\text{rpn}(x, x_t, \mathcal{F}_s)$ 
 $V_T^{nt}, V_T^{text} \leftarrow \text{RoIPool}(\mathcal{F}_T, \text{proposals})$ 
 $V_S^{nt}, V_S^{text} \leftarrow \text{RoIPool}(\mathcal{F}_S, \text{proposals})$ 
 $V_T^{mine}, V_S^{mine} \leftarrow \text{Node Indexing}(V_T^{text}, V_S^{text}, T)$ 
 $V_T \leftarrow V_T^{nt} \oplus V_T^{mine}$ 
 $V_S \leftarrow V_S^{nt} \oplus V_S^{mine}$ 
 $\xi_T \leftarrow \text{CosineSimilarity}(V_T)$ 
 $\xi_S \leftarrow \text{CosineSimilarity}(V_S)$ 
 $\mathcal{L}_v^{text} \leftarrow \text{MahalanobisDistance}(V_T^{text}, V_S^{text})$ 
 $\mathcal{L}_v^{nt} \leftarrow \text{MahalanobisDistance}(V_T^{nt}, V_S^{nt})$ 
 $\mathcal{L}_\xi \leftarrow \text{MahalanobisDistance}(\xi_T, \xi_S)$ 
Return  $\mathcal{L}_v^{text} + \mathcal{L}_v^{nt} + \mathcal{L}_\xi$ 

```

As depicted in Algorithm 2 \mathcal{L}_ξ is also an imitation loss between the entire edge set. This helps in aligning the edges of the student with those of the teacher. In our tests, imitating features doesn’t fully tap into the depth of available knowledge. If the node loss doesn’t capture the intricate semantic relationships effectively, then the edge loss steps in to drive the learning of these pairwise interactions. Hence, to synchronize the knowledge structure between the student and teacher, it’s vital to craft the edge loss in a way that encapsulates the overarching structured data within detectors. The only thing left is to match the output logits. We utilize a KL Divergence loss in order to incorporate these logits matching as defined in [68].

After this stage, the student model can correctly identify object regions (i.e. the bounding boxes) in a document. However, it doesn’t have any about object categories except "Text" and "non-text" labels. In order to, further categorize

the "non-text node" and predict the correct object labels we utilize a GNN [62] with a cross-entropy loss for node classification.

4 Experimental Evaluation

For validation purposes, we have considered four important benchmark datasets (PubLayNet [75], PRIMA [19], Historical Japanese [57], and DocLayNet [51]) which cover most of the existing document object categories (see Table 1). Our experimentation shows that the proposed GraphKD provides similar results to the large-scale supervised models (SwinDocSegmenter [7], LayoutLMv3 [32], DocSegTr [8], and so on) with a lesser no. of parameters.

Table 1: Experimental dataset description (instance level)

PubLayNet		PRIMA			Historical Japanese			DocLayNet			
Object	Train	Eval	Object	Train	Eval	Object	Train	Eval	Object	Train	Eval
Text	2,343,356	88,625	Text	6401	1531	Body	1443	308	Caption	20280	1543
Title	627,125	18,801	Image	761	163	Row	7742	1538	Footnote	5964	387
Lists	80,759	4239	Table	37	10	Title	33,637	7271	Formula	22367	1966
Figures	109,292	4327	Math	35	7	Bio	38,034	8207	List-item	170889	10522
Tables	102,514	4769	Separator	748	155	Name	66,515	7257	Page-footer	64717	3994
-	-	-	other	86	25	Position	33,576	7256	Page-header	50700	3366
-	-	-	-	-	-	Other	103	29	Picture	39621	3534
-	-	-	-	-	-	-	-	-	Section-header	18003	8550
-	-	-	-	-	-	-	-	-	Table	30115	2394
-	-	-	-	-	-	-	-	-	Text	431222	29940
-	-	-	-	-	-	-	-	-	Title	4423	335
Total	3,263,046	120,761	Total	8068	1891	Total	181,097	31,866	Total	994123	66531

It has been observed in Table 1, that PRIMA [19] has only 8068 no. of training instances however "Table", "Math", and "Other" region have only 37, 35, and 86 no. of training instances which is not sufficient to define an object class makes the dataset more challenging to solve. On the other hand, PubLayNet [75] and HJ [57] have only 5 and 7 object classes respectively, and all the classes have sufficient no. of training instances except the "Others" in HJ. However, it is very difficult to define the "Others" class properly. Last but not the least, DocLayNet [51] has 11 classes and provides a large variability bias towards some classes (i.e. the training samples are not equally distributed. For example, List-item has $\approx 170K$ training instances whereas, Title has only $\approx 4K$ training instances.).

Now, in order to evaluate the results on the aforementioned dataset, we utilize the IoU score which assesses instance segmentation accuracy, with COCO benchmarks using mean AP across IoU thresholds (0.5 to 0.95 with a step size of 0.05) to calculate mAP. This metric has been used throughout the rest of the paper to evaluate the performance of the proposed settings.

4.1 Ablation Studies

In this context, We have performed two sets of the ablation. In the first set (see Table 2), we have used the various combinations of components to emphasize the distinct contribution of each element in our graph distillation approach. Our framework incorporates three distinct modules that contribute to the graph distillation loss: 1) edge, 2) text node, and 3) non-text node. It has been observed that maintaining a consistent edge structure between student and teacher models adds $\approx 3\%$ AP points to the distillation outcome. This suggests that beyond direct pixel-to-pixel imitation, simply ensuring alignment in relationships can serve as a critical regularization to maintain topological configuration. This validates the effectiveness of our approach.

The introduction of student features mirroring non-text object nodes leads to a notable increase of $\approx 4\%$ in Average Precision (AP), surpassing the impact achieved by considering edges. This implies that distilling non-text features effectively encourages the student networks to prioritize regions containing non-text instances (i.e. helps to reduce the text bias). This underscores the significance of feature alignment within these foreground-labeled regions as more crucial for the student to emulate than the broader and potentially noisy high-dimensional feature maps. Additionally, when edges collaborate with nodes, it produces even more promising outcomes, affirming the effectiveness of both components of the graph, as they complement each other.

Table 2: Ablation study on the building blocks of GraphKD with DocLayNet using resnet50 to resnet18 setup

Student	Edge	Non-text	Text	AP	AP@50	AP@75	APs	APm	API
✓				33.1	52.5	36.0	19.1	36.6	41.3
✓	✓			36.6	59.8	40.1	21.2	38.1	41.4
✓	✓	✓		40.2	62.1	33.8	24.7	34.2	41.2
✓	✓	✓	✓	42.1	80.5	36.3	29.0	34.7	41.8
Teacher				61.2	87.9	46.3	30.1	40.2	42.7

Lastly, incorporating the imitation of student features within text regions results in an additional increase in AP, specifically a gain of $\approx 2\%$ when compared to the gains achieved with foreground nodes. This indicates that, even when the primary focus is on imitating non-text object nodes, the inclusion of text nodes becomes significant in the process of student distillation, particularly when balanced using our adaptive Mahalanobis distance loss. (NOTE: in Table 2 there is a large performance gap between the performance of teacher and student models as ResNet50 and ResNet18 pose different numbers of layers in each of the ResNet blocks, so some important feature has been lost during the feature compression in node-to-node distillation.)

In the second set, we have performed the ablation of different distance functions we can use to perform the node-to-node and edge-to-edge knowledge distillation (see Table 3). It has been observed that, L1 and L2 distance as a loss

function performs the worst as it only computes the absolute error or least square errors between the nodes or the edges. It neither considers the similarity nor their covariance. So when we are distilling using L1 or L2 losses it just penalizes the error based on the feature values of the node or the weights of the edges.

Table 3: Ablation study on the loss functions combinations of GraphKD with DocLayNet using resnet50 to resnet18 setup

	Node2Node	Edge2Edge	AP	AP@50	AP@75	APs	APm	API
L1		L1	3.0	6.1	2.8	0.4	1.6	2.2
		L2	7.1	8.3	3.5	1.2	1.6	4.3
		Cosine	7.1	10.5	8.2	1.5	2.1	6.7
		Mahalanobis	8.9	11.2	9.0	4.4	3.3	10.2
L2		L1	10.2	11.9	14.6	6.0	7.5	12.7
		L2	10.2	13.4	14.8	6.8	7.8	20.7
		Cosine	10.2	15.1	16.3	7.3	9.9	23.2
		Mahalanobis	12.9	18.8	23.2	10.0	12.0	24.9
Mahalanobis		L1	19.5	24.8	25.6	13.1	15.4	25.5
		L2	23.5	61.7	29.4	14.4	23.0	28.1
		Cosine	24.3	63.6	29.7	20.6	23.6	30.6
		Mahalanobis	28.5	64.5	31.6	21.9	23.6	31.7
Cosine		L1	30.4	66.8	31.7	23.3	26.9	36.9
		L2	38.9	68.5	32.7	25.6	28.3	37.2
		Cosine	38.9	79.5	33.5	27.5	34.1	41.5
		Mahalanobis	42.1	80.5	36.3	29.0	34.7	41.8

On the other hand, when we are using the Cosine distance as a loss function it considers the similarity which is beneficial for node-to-node distillation but not for edge-to-edge distillation. Similarly, Mahalanobis distance considers covariance but not similarity which is also beneficial for edge-to-edge distance. That’s why the best combination has been obtained by using Cosine distance for node-to-node distillation and Mahalanobis distance for edge-to-edge distillation.

4.2 Quantitative Evaluation

In order to establish a robust quantitative evaluation of GraphKD, we obtain Homogeneous (ResNet152-ResNet101 and ResNet101-Resnet50) as well as Heterogeneous (ResNet50-ResNet18, Resnet101-EfficientnetB0, and ResNet50-MobileNetv2) knowledge distillation across all four competitive benchmarks (PublayNet, PRIMA, Historical Japanese, and DoclayNet respectively). In homogeneous distillation, the number of layers in every ResNet block is the same for the teacher and the student network only we reduce the number of ResNet blocks from teacher to student. On the other hand, in heterogeneous distillation, both the number of layers in every Resnet block and the number of ResNet blocks have been reduced from teacher to student which leads to poor performance compared

to homogeneous distillation (which needs only block-wise feature compression) due to double compression of RoI pooled features during knowledge distillation.

Table 4: Graph-based knowledge distillation on PublayNet dataset

	#params (t)	#params (s)	Text	Title	List	Table	Figure	AP	AP@50	AP@75	APs	APm	API
R50-R18	25.6M	11.1	18.9	32.2	32.9	28.9	27.0	28.0	72.8	6.7	13.7	27.0	28.6
R101-R50	44.5M	25.6M	91.0	82.9	85.2	95.0	88.9	88.6	97.0	94.5	38.6	73.8	93.2
R152-R101	60.2M	44.5M	90.9	82.3	85.6	95.3	89.1	88.8	97.0	94.7	38.6	73.9	93.7
R101-EB0	44.5M	5.3M	19.2	34.8	28.1	29.1	26.8	27.6	72.2	7.1	13.0	25.0	28.1
R50-MNv2	25.6M	3.4M	18.5	32.1	32.9	29.3	28.2	28.2	72.8	7.2	13.5	24.2	29.6
DocSegTr [8]	-	168M	91.1	75.6	91.5	97.9	97.1	90.4	97.9	95.8	-	-	-
LayoutLMv3 [32]	-	368M	94.5	90.6	95.5	97.9	97.9	95.1	-	-	-	-	-
SwinDocSegmenter [7]	-	223M	94.5	87.1	93.0	97.9	97.2	93.7	97.9	96.2	-	-	-

Moreover, in Table 4 we have obtained our quantitative evaluation on the PublayNet dataset. For a fair comparison, we compared the performance of our distilled network with state-of-the-art supervised approaches [8,7,32] and it is obvious that we cannot outperform those transformer based large scale networks however, we can reduce the performance gap and increase the efficiency. For example, the performance gap between LayoutLMv3 [32] and distilled ResNet101 Faster-RCNN network is $\approx 7\%$ (Overall AP) however, we can reduce 323.5M(368M-44.5M) number of model parameters. Also, it outperforms DocSegTr (168M) [8] for the "Title" class by $\approx 8\%$ as the Transformers are not the best choice for small object detection.

Table 5: Graph-based knowledge distillation on PRIMA dataset

	#params(t)	#params(s)	Text	Image	Table	Math	Separator	Other	AP	AP@50	AP@75	APs	APm	API
R50-R18	25.6M	11.1M	38.5	47.7	41.6	8.5	17.5	5.4	26.5	52.1	22.3	30.8	30.9	28.3
R101-R50	44.5M	25.6M	76.6	60.8	38.3	4.1	23.4	7.1	35.0	51.0	39.0	39.6	42.1	36.6
R152-R101	60.2M	44.5M	79.9	64.6	44.0	30.2	25.9	6.7	41.9	58.9	42.0	41.4	41.7	43.5
R101-EB0	44.5M	5.3M	20.3	19.0	15.6	4.4	13.1	3.2	12.6	40.7	1.69	16.4	18.7	12.2
R50-MNv2	25.6M	3.4M	19.5	16.6	28.0	7.1	12.2	5.9	14.9	40.2	2.5	16.1	20.1	15.2
DocSegTr [8]	-	168M	75.2	64.3	59.4	48.4	1.8	3.0	42.5	54.2	45.8	-	-	-
LayoutLMv3 [32]	-	368M	70.8	50.1	42.5	46.5	9.6	17.4	40.3	-	-	-	-	-
SwinDocSegmenter [7]	-	223M	87.7	75.9	49.8	78.1	27.5	7.0	54.3	69.3	52.9	-	-	-

On the other hand, the distilled ResNet101 outperforms LayoutLMv3 [32] by $\approx 1\%$ which clearly shows large networks are not quite effective for small datasets like PRIMA (see Table 5) as in the later stage of training instance they only learn noise due to unavailability of data. Not only that, distilled EfficientnetB0 and MobileNetv2 which have only 5.3M and 3.4M parameters respectively, also outperform the LayoutLMv3 [32] and DocSegTr [8] for the "Separator class". Similarly, distilled ResNet101 outperforms SwinDocsegmenter [7] for the "other" classes which shows the real power of the knowledge distillation.

Similarly, for the Historical Japanese dataset (see Table 6) all the supervised methods perform better than the distilled networks as all the classes of this dataset are quite separable from each other. However, an interesting observation has been noticed in the performance of the DoclayNet dataset. Distilled ResNet50 and ResNet101 outperform DocSegTr [8] and LayoutLMv3 [32] by a significant margin for "Caption", "Page-footer", and "Picture" which shows how the bias factor affects the supervised training and the potential of knowledge dis-

Table 6: Graph-based knowledge distillation on Historical Japanese dataset

	#params (t)	#params (s)	Body	Row	Title	Bio	Name	Position	Other	AP	AP@50	AP@75	APs	APm	API
R50-R18	25.6M	11.1	37.6	25.9	50.1	28.0	44.0	35.7	12.4	33.4	67.2	22.7	20.7	23.3	28.9
R101-R50	44.5M	25.6M	88.2	98.3	84.9	94.6	68.6	83.5	30.2	78.3	86.0	84.7	38.2	44.6	71.5
R152-R101	60.2M	44.5M	98.4	98.3	85.1	94.5	68.6	84.2	28.5	79.7	87.4	86.0	37.6	43.6	74.8
R101-EB0	44.5M	5.3M	31.1	22.0	49.9	25.2	20.2	65.3	18.4	33.1	65.9	25.9	19.1	23.5	29.4
R50-MNv2	25.6M	3.4M	30.1	26.5	61.6	23.3	36.0	65.1	20.3	37.5	71.1	29.0	16.7	29.0	30.7
DocSegTr [8]	-	168M	99.0	99.1	93.2	94.7	70.3	87.4	43.7	83.1	90.1	88.1	-	-	-
LayoutLMv3 [32]	-	368M	99.0	99.0	92.9	94.7	67.9	87.8	38.7	82.7	-	-	-	-	-
SwinDocSegmenter [7]	-	223M	99.7	99.0	89.5	86.2	83.8	93.0	40.5	84.5	90.7	88.2	-	-	-

tillation to overcome it by forming a more generalized efficient network which we can use in our edge devices.

Table 7: Graph-based knowledge distillation on DocLayNet dataset

	#params(t)	#params(s)	Caption	Footnote	Formula	List item	Page footer	Page header	Picture	Section header	Table	Text	Title	AP
R50-R18	25.6M	11.1M	53.4	23.3	30.6	39.7	45.7	44.2	58.4	43.0	45.6	41.5	37.3	42.1
R101-R50	44.5M	25.6M	77.3	46.4	48.1	72.3	60.4	63.0	72.9	59.3	73.5	77.5	64.6	65.0
R152-R101	60.2M	44.5M	78.9	58.1	53.7	75.3	59.3	67.3	76.4	61.6	78.0	80.3	69.0	68.9
R101-EB0	44.5M	5.3M	36.35	20.7	27.5	32.1	26.4	34.6	24.8	32.5	25.2	25.0	32.9	28.9
R50-MNv2	25.6M	3.4M	28.8	13.5	20.2	24.4	26.4	22.3	25.4	25.8	24.6	24.3	23.8	23.6
DocSegTr [8]	-	168M	70.1	73.7	63.5	81.0	58.9	72.0	72.0	68.4	82.2	85.4	79.9	73.4
LayoutLMv3 [32]	-	368M	71.5	71.8	63.4	80.8	59.3	70.0	72.7	69.3	82.9	85.8	80.4	73.5
SwinDocSegmenter [7]	-	223M	83.5	64.8	62.3	82.3	65.1	66.3	84.7	66.5	87.4	88.2	63.2	76.8

4.3 Comparative Study

In order to, establish the superiority of the proposed method we have performed a detailed comparative study with a feature-based knowledge distillation method ReviewKD [16], a logit-based knowledge distillation method NKD [18], and with a hybrid approach SimKD [10]. The result of this comparative study has been depicted in Table 8. It has been observed that the performance of the DOD has improved quite a lot through GraphKD than the rest.

It is worth noting that, the performance of the feature-based method (ReviewKD [16]) is better compared to the logit-based method as the logits are like labels. That’s why this NKD [18] will be effective for the classification task not for the object detection task. On the other hand, when we combine the logits and features in SimKD [10] it outperforms the individuals as it gets the class-level information along with the features. Lastly, GraphKD creates nodes from the anchor boxes and also performs the node indexing which gets rid of the two complex task feature matching and logit matching through an easier node matching, which affects performance gain.

5 Conclusion

In this paper, we have presented *GraphKD*, a graph-based knowledge distillation strategy for document object detection to optimize the number of parameters so that, we use them in edge devices. According to our best knowledge, this is the first time, we have explored knowledge distillation for DOD where we used RoI pooled features for structured graph creation and performed an effective

Table 8: Comparative study of GraphKD with the state-of-the-art approaches

Method	Backbone	PublayNet			PRIMA			HJ			DocLayNet		
		AP	AP@50	AP@75									
ReviewKD [16]	R50-R18	20.9	56.6	3.4	17.7	29.1	12.0	19.5	36.2	19.7	18.2	60.2	21.2
	R101-R50	72.7	94.4	90.2	22.7	38.8	26.0	62.8	79.7	74.2	61.1	84.3	74.1
	R152-R101	77.1	95.3	90.8	26.2	44.7	30.7	66.7	77.2	69.8	63.7	84.6	75.0
	R101-EB0	20.1	62.5	6.7	7.1	28.6	0.9	22.1	50.9	22.8	12.7	50.1	7.6
	R50-MNv2	19.8	63.6	6.3	9.2	20.7	1.0	25.4	58.7	15.7	14.2	44.2	5.1
NKD [18]	R50-R18	17.9	29.0	1.7	9.7	10.7	8.1	27.5	61.4	20.7	14.0	58.9	16.0
	R101-R50	68.3	89.6	81.2	12.2	28.7	16.9	52.8	59.9	74.2	37.2	57.0	40.4
	R152-R101	70.7	90.7	82.3	16.8	34.2	20.9	56.9	65.2	69.8	41.5	62.1	45.2
	R101-EB0	17.7	58.0	2.7	7.0	18.7	0.7	11.7	40.8	20.3	10.2	40.1	5.6
	R50-MNv2	13.2	63.0	3.1	8.2	10.7	0.9	15.4	38.2	17.3	12.1	34.2	5.9
SimKD [10]	R50-R18	24.7	65.4	5.1	23.6	32.7	17.2	19.1	58.0	17.4	33.2	77.8	35.7
	R101-R50	80.2	95.2	92.1	29.7	48.8	36.0	72.3	80.7	79.7	62.7	85.2	74.2
	R152-R101	81.1	95.9	92.2	36.2	54.7	40.7	76.7	82.1	84.2	64.6	88.6	77.2
	R101-EB0	24.3	71.4	7.0	9.1	38.6	0.9	27.1	55.9	15.9	22.0	60.2	9.7
	R50-MNv2	26.1	70.4	7.0	10.2	30.7	1.7	27.5	61.4	20.7	19.7	54.2	7.1
GraphKD	R50-R18	28.0	72.8	6.7	26.5	52.1	22.3	33.4	67.2	22.7	42.1	80.5	36.3
	R101-R50	88.6	97.0	94.5	35.0	51.0	39.0	78.3	86.0	84.7	65.0	86.7	74.3
	R152-R101	88.8	97.0	94.7	41.9	58.9	42.0	79.7	87.4	86.0	68.9	89.0	78.0
	R101-EB0	27.6	72.2	7.1	12.6	40.7	1.6	33.1	65.9	25.9	28.9	68.1	13.9
	R50-MNv2	28.2	72.8	7.2	14.9	40.2	2.5	37.5	71.1	29.0	23.6	62.4	7.3

sample mining strategy to reduce the text bias. We also performed node-to-node distillation based on the similarity (using Cosine distance) and edge-to-edge distillation based on their covariance (i.e. Mahalanobis distance). From the results, it can be concluded that we are successful in reducing the number of parameters however, it also affects the performance ($\approx 7\%$ for PublayNet and DocLayNet and $\approx 13\%$ for PRIMA compared to SwinDocSegmenter). Also, we are unable to incorporate Transformers as a backbone in order to perform cross-architecture distillation (Transformers to ResNet) effectively due to their different data handling mechanism. It will be our future perspective of this work.

Acknowledgment

This work has been partially supported by the Spanish project PID2021-126808OB-I00, the Catalan project 2021 SGR 01559 and the PhD Scholarship from AGAUR (2021FIB-10010). The Computer Vision Center is part of the CERCA Program/Generalitat de Catalunya.

References

1. Aditya, S., Saha, R., Yang, Y., Baral, C.: Spatial knowledge distillation to aid visual reasoning. In: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 227–235 (2019)
2. Ahn, S., Hu, S.X., Damianou, A., Lawrence, N.D., Dai, Z.: Variational information distillation for knowledge transfer. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9163–9171 (2019)

3. Akiba, T., Sano, S., Yanase, T., Ohta, T., Koyama, M.: Optuna: A next-generation hyperparameter optimization framework. In: Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining. pp. 2623–2631 (2019)
4. Asi, A., Cohen, R., Kedem, K., El-Sana, J.: Simplifying the reading of historical manuscripts. In: 2015 13th International Conference on Document Analysis and Recognition (ICDAR). pp. 826–830. IEEE (2015)
5. Ba, J., Caruana, R.: Do deep nets really need to be deep? Advances in neural information processing systems (2014)
6. Banerjee, A., Biswas, S., Lladós, J., Pal, U.: Semidocseg: Harnessing semi-supervised learning for document layout analysis (2023)
7. Banerjee, A., Biswas, S., Lladós, J., Pal, U.: Swindocsegmenter: An end-to-end unified domain adaptive transformer for document instance segmentation. arXiv preprint arXiv:2305.04609 (2023)
8. Biswas, S., Banerjee, A., Lladós, J., Pal, U.: Docsegtr: An instance-level end-to-end document image segmentation transformer. arXiv preprint arXiv:2201.11438 (2022)
9. Chawla, A., Yin, H., Molchanov, P., Alvarez, J.: Data-free knowledge distillation for object detection. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 3289–3298 (2021)
10. Chen, D., Mei, J.P., Zhang, H., Wang, C., Feng, Y., Chen, C.: Knowledge distillation with the reused teacher classifier. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (2022)
11. Chen, D., Mei, J.P., Zhang, Y., Wang, C., Wang, Z., Feng, Y., Chen, C.: Cross-layer distillation with semantic calibration. In: Proceedings of the AAAI Conference on Artificial Intelligence (2021)
12. Chen, G., Choi, W., Yu, X., Han, T., Chandraker, M.: Learning efficient object detection models with knowledge distillation. Advances in neural information processing systems **30** (2017)
13. Chen, J., Lopresti, D.: Table detection in noisy off-line handwritten documents. In: 2011 International Conference on Document Analysis and Recognition. pp. 399–403. IEEE (2011)
14. Chen, K., Seuret, M., Hennebert, J., Ingold, R.: Convolutional neural networks for page segmentation of historical document images. In: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR). vol. 1, pp. 965–970. IEEE (2017)
15. Chen, K., Seuret, M., Liwicki, M., Hennebert, J., Ingold, R.: Page segmentation of historical document images with convolutional autoencoders. In: 2015 13th International Conference on Document Analysis and Recognition (ICDAR). pp. 1011–1015. IEEE (2015)
16. Chen, P., Liu, S., Zhao, H., Jia, J.: Distilling knowledge via knowledge review. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2021)
17. Chen, Y., Chen, P., Liu, S., Wang, L., Jia, J.: Deep structured instance graph for distilling object detectors. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4359–4368 (2021)
18. Chi, Z., Zheng, T., Li, H., Yang, Z., Wu, B., Lin, B., Cai, D.: Normkd: Normalized logits for knowledge distillation. arXiv preprint arXiv:2308.00520 (2023)
19. Clausner, C., Antonacopoulos, A., Pletschacher, S.: Icdar2019 competition on recognition of documents with complex layouts-rdcl2019. In: Proceedings of the

- International Conference on Document Analysis and Recognition. pp. 1521–1526 (2019)
20. Coquenot, D., Chatelain, C., Paquet, T.: Dan: a segmentation-free document attention network for handwritten document recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023)
 21. Da, C., Luo, C., Zheng, Q., Yao, C.: Vision grid transformer for document layout analysis. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 19462–19472 (2023)
 22. Dai, X., Jiang, Z., Wu, Z., Bao, Y., Wang, Z., Liu, S., Zhou, E.: General instance distillation for object detection. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 7842–7851 (2021)
 23. De Rijk, P., Schneider, L., Cordts, M., Gavrila, D.: Structural knowledge distillation for object detection. *Advances in Neural Information Processing Systems* **35**, 3858–3870 (2022)
 24. Deng, Q., Ibrayim, M., Hamdulla, A., Zhang, C.: The yolo model that still excels in document layout analysis. *Signal, Image and Video Processing* pp. 1–10 (2023)
 25. Douzon, T., Duffner, S., Garcia, C., Espinas, J.: Long-range transformer architectures for document understanding. In: *International Conference on Document Analysis and Recognition*. pp. 47–64. Springer (2023)
 26. Fang, J., Gao, L., Bai, K., Qiu, R., Tao, X., Tang, Z.: A table detection method for multipage pdf documents via visual separators and tabular structures. In: *2011 International Conference on Document Analysis and Recognition*. pp. 779–783. IEEE (2011)
 27. Fateh, A., Fateh, M., Abolghasemi, V.: Enhancing optical character recognition: Efficient techniques for document layout analysis and text line detection. *Engineering Reports* p. e12832 (2023)
 28. Gong, L., Lin, S., Zhang, B., Shen, Y., Li, K., Qiao, R., Ren, B., Li, M., Yu, Z., Ma, L.: Adaptive hierarchy-branch fusion for online knowledge distillation. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 37, pp. 7731–7739 (2023)
 29. Gou, J., Xiong, X., Yu, B., Du, L., Zhan, Y., Tao, D.: Multi-target knowledge distillation via student self-reflection. *International Journal of Computer Vision* pp. 1–18 (2023)
 30. Heo, B., Lee, M., Yun, S., Choi, J.Y.: Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 33, pp. 3779–3787 (2019)
 31. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015)
 32. Huang, Y., Lv, T., Cui, L., Lu, Y., Wei, F.: Layoutlmv3: Pre-training for document ai with unified text and image masking. In: *Proceedings of the 30th ACM International Conference on Multimedia*. pp. 4083–4091 (2022)
 33. Journet, N., Eglin, V., Ramel, J.Y., Mullot, R.: Text/graphic labelling of ancient printed documents. In: *Eighth International Conference on Document Analysis and Recognition (ICDAR’05)*. pp. 1010–1014. IEEE (2005)
 34. Kang, Z., Zhang, P., Zhang, X., Sun, J., Zheng, N.: Instance-conditional knowledge distillation for object detection. *Advances in Neural Information Processing Systems* **34**, 16468–16480 (2021)
 35. Kise, K., Sato, A., Iwata, M.: Segmentation of page images using the area voronoi diagram. *Computer Vision and Image Understanding* **70**(3), 370–382 (1998)

36. Komodakis, N., Zagoruyko, S.: Paying more attention to attention: improving the performance of convolutional neural networks via attention transfer. In: ICLR (2017)
37. Li, J., Xu, Y., Lv, T., Cui, L., Zhang, C., Wei, F.: Dit: Self-supervised pre-training for document image transformer. In: Proceedings of the 30th ACM International Conference on Multimedia. pp. 3530–3539 (2022)
38. Li, K., Wigington, C., Tensmeyer, C., Zhao, H., Barmpalios, N., Morariu, V.I., Manjunatha, V., Sun, T., Fu, Y.: Cross-domain document object detection: Benchmark suite and method. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12915–12924 (2020)
39. Li, X.H., Yin, F., Liu, C.L.: Page segmentation using convolutional neural network and graphical model. In: Document Analysis Systems: 14th IAPR International Workshop, DAS 2020, Wuhan, China, July 26–29, 2020, Proceedings 14. pp. 231–245. Springer (2020)
40. Li, Z., Xu, P., Chang, X., Yang, L., Zhang, Y., Yao, L., Chen, X.: When object detection meets knowledge distillation: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023)
41. Liao, H., RoyChowdhury, A., Li, W., Bansal, A., Zhang, Y., Tu, Z., Satzoda, R.K., Manmatha, R., Mahadevan, V.: Doctr: Document transformer for structured information extraction in documents. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 19584–19594 (2023)
42. Lin, G.S., Tu, J.C., Lin, J.Y.: Keyword detection based on retinanet and transfer learning for personal information protection in document images. *Applied Sciences* **11**(20), 9528 (2021)
43. Lin, H., Han, G., Ma, J., Huang, S., Lin, X., Chang, S.F.: Supervised masked knowledge distillation for few-shot transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19649–19659 (2023)
44. Liu, Y., Chen, K., Liu, C., Qin, Z., Luo, Z., Wang, J.: Structured knowledge distillation for semantic segmentation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2604–2613 (2019)
45. Markewich, L., Zhang, H., Xing, Y., Lambert-Shirzad, N., Jiang, Z., Lee, R.K.W., Li, Z., Ko, S.B.: Segmentation for document layout analysis: not dead yet. *International Journal on Document Analysis and Recognition (IJDAR)* pp. 1–11 (2022)
46. Mathur, P., Jain, R., Mehra, A., Gu, J., Dernoncourt, F., Tran, Q., Kaynig-Fittkau, V., Nenkova, A., Manocha, D., Morariu, V.I., et al.: Layerdoc: Layer-wise extraction of spatial hierarchical structure in visually-rich documents. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 3610–3620 (2023)
47. Mirzadeh, S.I., Farajtabar, M., Li, A., Levine, N., Matsukawa, A., Ghasemzadeh, H.: Improved knowledge distillation via teacher assistant. In: Proceedings of the AAAI conference on artificial intelligence. vol. 34, pp. 5191–5198 (2020)
48. Mooney, R.J., Bunescu, R.: Mining knowledge from text using information extraction. *ACM SIGKDD explorations newsletter* **7**(1), 3–10 (2005)
49. Negrinho, R., Gormley, M., Gordon, G.J.: Learning beam search policies via imitation learning. *Advances in Neural Information Processing Systems* **31** (2018)
50. Oliveira, S.A., Seguin, B., Kaplan, F.: dhsegment: A generic deep-learning approach for document segmentation. In: 2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR). pp. 7–12. IEEE (2018)
51. Pfitzmann, B., Auer, C., Dolfi, M., Nassar, A.S., Staar, P.: Doclaynet: A large human-annotated dataset for document-layout segmentation. In: Proceedings of

- the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. pp. 3743–3751 (2022)
52. Powalski, R., Borchmann, L., Jurkiewicz, D., Dwojak, T., Pietruszka, M., Palka, G.: Going full-tilt boogie on document understanding with text-image-layout transformer. In: Document Analysis and Recognition–ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part II 16. pp. 732–747. Springer (2021)
 53. Rahal, N., Vöggtlin, L., Ingold, R.: Layout analysis of historical document images using a light fully convolutional network. In: International Conference on Document Analysis and Recognition. pp. 325–341. Springer (2023)
 54. Saabni, R., El-Sana, J.: Language-independent text lines extraction using seam carving. In: 2011 International Conference on Document Analysis and Recognition. pp. 563–568. IEEE (2011)
 55. Saha, R., Mondal, A., Jawahar, C.: Graphical object detection in document images. In: 2019 International Conference on Document Analysis and Recognition (ICDAR). pp. 51–58. IEEE (2019)
 56. Schreiber, S., Agne, S., Wolf, I., Dengel, A., Ahmed, S.: Deepdesrt: Deep learning for detection and structure recognition of tables in document images. In: 2017 14th IAPR international conference on document analysis and recognition (ICDAR). vol. 1, pp. 1162–1167. IEEE (2017)
 57. Shen, Z., Zhang, K., Dell, M.: A large dataset of historical japanese documents with complex layouts. In: Proceedings of the IEEE Conference on CVPRW. pp. 548–549 (2020)
 58. Shrivastava, A., Gupta, A., Girshick, R.: Training region-based object detectors with online hard example mining. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 761–769 (2016)
 59. Stanisławek, T., Graliński, F., Wróblewska, A., Lipiński, D., Kaliska, A., Rosalska, P., Topolski, B., Biecek, P.: Kleister: key information extraction datasets involving long documents with complex layouts. In: International Conference on Document Analysis and Recognition. pp. 564–579. Springer (2021)
 60. Sun, N., Zhu, Y., Hu, X.: Faster r-cnn based table detection combining corner locating. In: 2019 International Conference on Document Analysis and Recognition (ICDAR). pp. 1314–1319. IEEE (2019)
 61. Tang, Z., Yang, Z., Wang, G., Fang, Y., Liu, Y., Zhu, C., Zeng, M., Zhang, C., Bansal, M.: Unifying vision, text, and layout for universal document processing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19254–19264 (2023)
 62. Wang, Y., Weng, X., Kitani, K.: Joint detection and multi-object tracking with graph neural networks. arXiv preprint arXiv:2006.13164 1(2) (2020)
 63. Wu, A., Deng, C.: Single-domain generalized object detection in urban scene via cyclic-disentangled self-distillation. In: Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition. pp. 847–856 (2022)
 64. Wu, D., Chen, P., Yu, X., Li, G., Han, Z., Jiao, J.: Spatial self-distillation for object detection with inaccurate bounding boxes. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6855–6865 (2023)
 65. Wu, X., Zheng, D., Wang, R., Sun, J., Hu, M., Feng, F., Wang, X., Jiang, H., Yang, F.: A region-based document vqa. In: Proceedings of the 30th ACM International Conference on Multimedia. pp. 4909–4920 (2022)
 66. Yang, H., Hsu, W.: Transformer-based approach for document layout understanding. In: 2022 IEEE International Conference on Image Processing (ICIP). pp. 4043–4047. IEEE (2022)

67. Yang, H., Hsu, W.H.: Vision-based layout detection from scientific literature using recurrent convolutional neural networks. In: 2020 25th international conference on pattern recognition (ICPR). pp. 6455–6462. IEEE (2021)
68. Yang, X., Yang, X., Yang, J., Ming, Q., Wang, W., Tian, Q., Yan, J.: Learning high-precision bounding box for rotated object detection via kullback-leibler divergence. *Advances in Neural Information Processing Systems* **34**, 18381–18394 (2021)
69. Yang, Z., Zeng, A., Li, Z., Zhang, T., Yuan, C., Li, Y.: From knowledge distillation to self-knowledge distillation: A unified approach with normalized loss and customized soft labels. *arXiv preprint arXiv:2303.13005* (2023)
70. Zhang, L., Ma, K.: Improve object detection with feature-based knowledge distillation: Towards accurate and efficient detectors. In: *International Conference on Learning Representations* (2020)
71. Zhang, L., Ma, K.: Structured knowledge distillation for accurate and efficient object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023)
72. Zhang, P., Li, C., Qiao, L., Cheng, Z., Pu, S., Niu, Y., Wu, F.: Vsr: a unified framework for document layout analysis combining vision, semantics and relations. In: *Document Analysis and Recognition–ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part I 16*. pp. 115–130. Springer (2021)
73. Zhao, B., Cui, Q., Song, R., Qiu, Y., Liang, J.: Decoupled knowledge distillation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022)
74. Zheng, Z., Ye, R., Wang, P., Ren, D., Zuo, W., Hou, Q., Cheng, M.M.: Localization distillation for dense object detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 9407–9416 (2022)
75. Zhong, X., Tang, J., Yepes, A.J.: Publaynet: largest dataset ever for document layout analysis. In: *Proceedings of the International Conference on Document Analysis and Recognition*. pp. 1015–1022 (2019)
76. Zhong, Z., Wang, J., Sun, H., Hu, K., Zhang, E., Sun, L., Huo, Q.: A hybrid approach to document layout analysis for heterogeneous document images. In: *International Conference on Document Analysis and Recognition*. pp. 189–206. Springer (2023)

GraphKD: Exploring Knowledge Distillation Towards Document Object Detection with Structured Graph Creation

Supplementary

Material

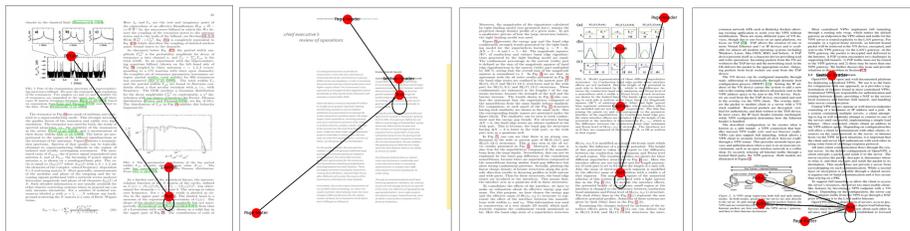
1 Implementation Details

We train *GraphKD* with an Adam optimizer with an initial learning rate of $1e - 5$ with a cosine annealing scheduler of 5000 cycles utilizing weight decay of $1e - 6$. To accumulate the final model we train for 90K iterations with a learning rate reduction by an order of magnitude 10 in the range of 70K to 80K. We prepared all models on NVIDIA A40 GPU of 48G RAM with 1 day of training utilizing stochastic learning. We utilize Pytorch and Detectron2 to build this framework. We need to tune three hyperparameters: namely edge distance for node merging, classification loss threshold during node indexing, and the covariance threshold during node creation. We select those values as 0.1, 0.6, and 0.8 respectively. Also, in the graph distillation loss we need to tune λ_1 , λ_2 , and λ_3 in order to balance all the loss terms via regularization, we selected them as 0.001, 0.008, and, 0.003 respectively. (Note: all these hyperparameter values have been selected via hyperparameter grid search).

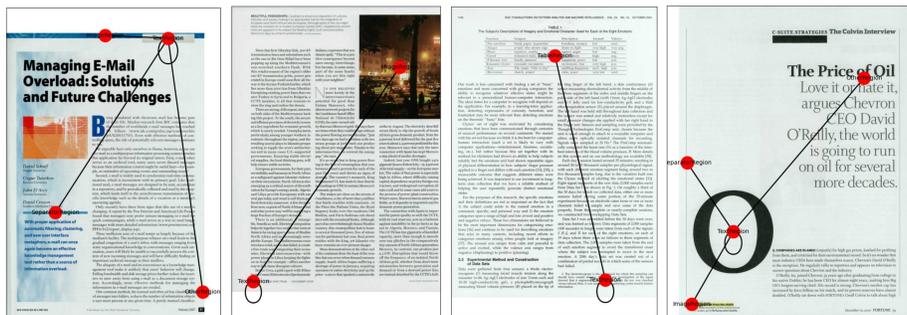
2 A deep dive into the dataset description

From Table 1 of the main paper of *GraphKD*, it has been observed that PubLayNet [75] has the highest number of training/validation instances while PRIMA [19] has the lowest. On the other hand, DocLayNet [51] has a large number of class instances although the data points are not well distributed. Besides that, it is important to understand the local and global relationships between the class instances. In order to identify that, we perform a UMAP visualization over the validation set as depicted in Fig. 1.

From the Fig. 1 it has been observed that PubLayNet [75] has a very dense data distribution among all the 5 classes (which helps to gain the instance segmentation performance) while PRIMA [19] suffers from the scarcity of the data problem. Other the other hand, in the historical Japanese dataset [57], classes are well separated, so it is easier to instantly segment them even with the low number of data points. Similarly, the DocLayNet has a large inter-class variability, with a large bias to the "list-item" and "text".



(a) Graph Creation on DocLayNet



(b) Graph Creation on PRIMA

Fig. 2: **Graph creation without node indexing:** Here one node represents the whole instances of each class as the nodes are developed on the feature embedding space.

provide comparable performance with large-scale transformer-based approaches. However, It is worth examining the performance of each feature extraction backbone for complex document object detection with supervised training. Table 1 depicts the performance of supervised training of all the networks on four competitive benchmarks.

Table 1: Performance of various convolution networks with supervised training

Backbone	PublayNet [75]			PRIMA [19]			HJ [57]			DocLayNet [51]		
	AP	AP@50	AP@75	AP	AP@50	AP@75	AP	AP@50	AP@75	AP	AP@50	AP@75
R18	24.7	70.9	12.3	22.1	42.1	17.3	29.7	60.1	21.8	32.7	70.8	30.1
R50	85.4	96.2	92.7	31.3	46.7	36.2	75.6	82.3	80.4	62.1	87.9	46.3
R101	86.7	96.9	93.1	39.7	57.8	40.5	76.9	87.1	85.7	64.9	88.4	76.2
R152	90.2	98.7	95.1	42.8	60.2	44.9	80.2	89.9	87.2	69.1	90.8	80.1
EBO	23.7	69.2	17.1	11.1	30.8	1.2	30.1	62.9	23.9	26.9	58.1	11.7
MNv2	24.2	70.8	21.1	12.6	35.9	2.0	34.8	70.5	26.1	20.6	56.4	6.3

It has been observed that all the trained in a supervised manner obtained poor performance compared to the distilled one as the teacher model is typically

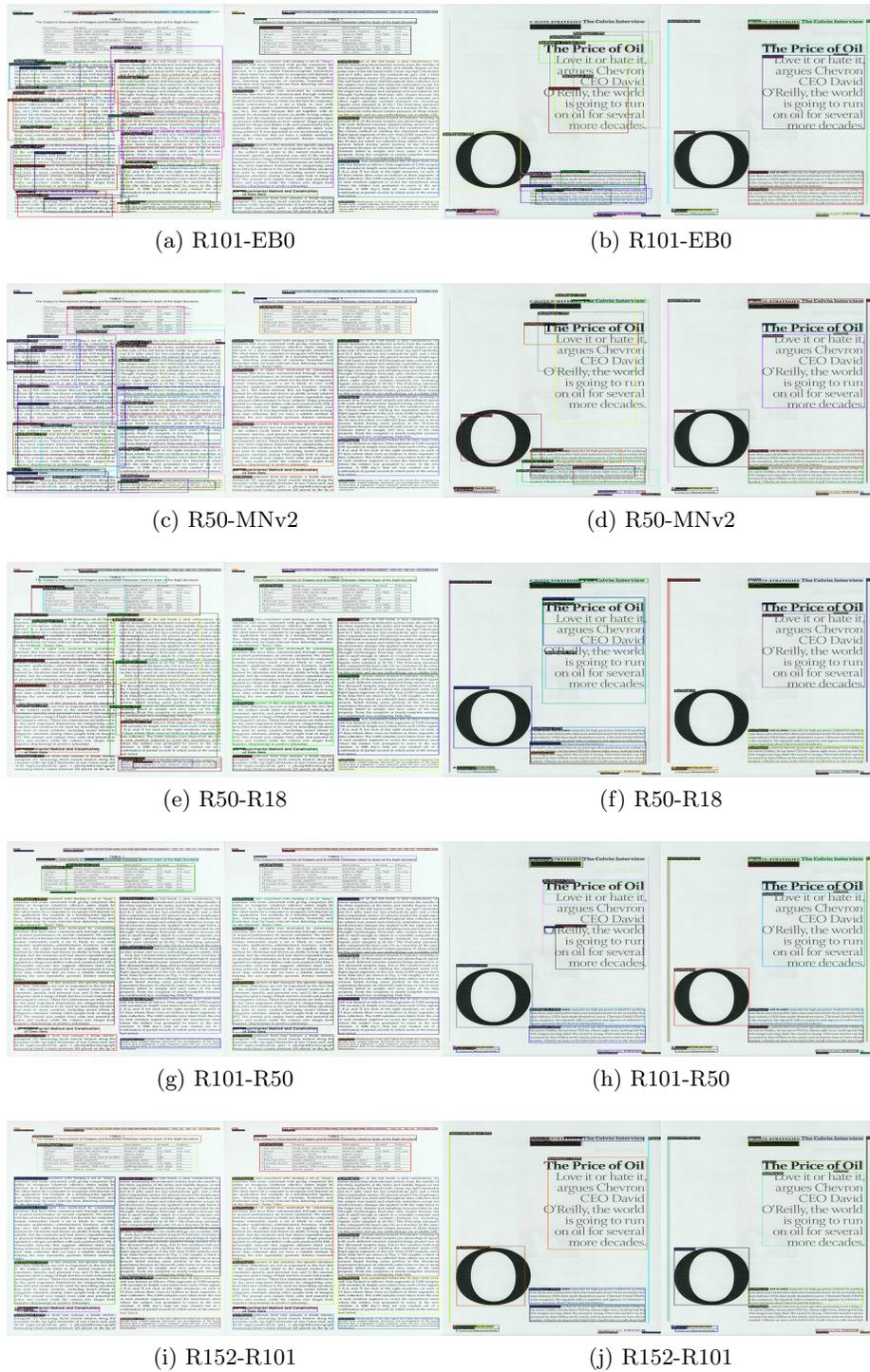


Fig. 3: Qualitative analysis with various distilled networks on PRIMA dataset (left: predicted; right: ground-truth).

a larger and more powerful model (e.g. ResNet152) that has learned to capture intricate patterns in the data. By distilling its knowledge into a smaller model, we effectively compress the information, making it easier to transfer and generalize to new data.

6 Failure case: cross architecture distillation

The motivation for using graphs is to transfer knowledge more effectively by reducing the feature alignment, layer dimension check, and so on. Similarly, we also aim to perform cross-architecture distillation from transformers (ViT-B) to ResNet (R101). The performance has been reported in Table 2.

Table 2: Performance of knowledge distillation from ViT-Base to ResNet50

Dataset	backbone	AP	AP@50	AP@75	APs	APm	API
PublayNet [75]	ViTB-R50	11.7	29.8	20.7	12.5	19.2	20.1
	ViT-B	91.7	98.9	96.7	40.2	59.8	62.3
PRIMA [19]	ViTB-R50	2.4	11.2	9.7	0.3	3.9	4.7
	ViT-B	46.1	62.6	47.3	31.3	33.4	50.5
HJ	ViTB-R50	9.8	22.8	16.7	8.2	15.1	16.7
	ViT-B	81.2	89.7	84.1	36.8	55.7	56.8
DoclayNet [51]	ViTB-R50	5.7	24.6	13.2	1.2	7.8	8.2
	ViT-B	65.6	84.7	73.6	37.8	55.4	59.2

It has been observed that the performance gap between the teacher and student model is quite large due to feature misalignment between the transformer layer and the convolution layer. Also, the transformer used a self-attention mechanism and in Resnet we used Relu activation so the Region proposal is completely different for these two backbones (feature compression is not possible), so we cannot use shared RPN which leads to this feature misalignment and we cannot perform node indexing based on the teacher classification loss during distillation.