

Shubham Paliwal, Vishwanath D , Rohit Rahul, Monika Sharma, Lovekesh Vig

TCS Research , New Delhi

{ shubham.p3 , vishwanath.d2 , monika.sharmal , rohit.rahul , lovekesh.vig } @ tcs.com

Abstract— With the widespread use of mobile phones and of 1) table detection and 2) table structure recognition, and
 2020 scanners to photograph and upload documents, the need for attempt to solve each sub-problem independently. While table
 arXiv:2009.04489v1 [cs.CV] 16 Jan 2020 extracting the information trapped in unstructured document detection involves detection of the image pixel coordinates
 images such as retail receipts, insurance, claim forms and financial invoices is becoming more acute. A major hurdle to this objective containing the tabular sub-image, tabular structure recognition
 is that these images often contain information in the form of involves segmentation of the individual rows and columns in
 Jan tables and extracting data from tabular sub-images presents a the detected table .
 unique set of challenges. This includes accurate detection of the In this paper, we propose TableNet, a novel end-to-end
 6 tabular region within an image, and subsequently detecting and deep learning model that exploits the inherent interdependence
 extracting information from the rows and columns of the detected between the twin tasks of table detection and table structure
 table. While some progress has been made in table detection, identification. The model utilizes a base network that is
 [cs.CV] extracting the table contents is still a challenge since this involves initialized with pre-trained VGG-19 features . This is followed
 more fine grained table structure(rows & columns) recognition . Prior approaches have attempted to solve the table detection and by two decoder branches for 1) Segmentation of the table
 structure recognition problems independently using two separate region and 2) Segmentation of the columns within a table models. In this v
 paper,
 end deep learning model for both table detection and structure region. Subsequently, rule based row extraction is employed
 recognition. The model exploits the interdependence between the to extract data in individual table cells .
 2001.01469v1 twin tasks of table detection and table structure recognition to A multi-task approach is used for the training of the deep
 segment out the table and column regions. This is followed by model. The model takes a single input image and produces
 semantic rule-based row extraction from the identified tabular two different semantically labelled output images for tables
 sub-regions. The proposed model and extraction approach was and columns . The model shares the encoding layer of VGG
 evaluated on the publicly available ICDAR 2013 and Marmot 19 for both the table and column detectors, while the decoders
 Table datasets obtaining state of the art results . Additionally, for the two tasks are separate . The shared common layers are
 we demonstrate that feeding additional semantic features further repeatedly trained from the gradients received from both the
 arXiv: improves model performance and that the model exhibits transfer learning across datasets. Another contribution of this paper is to
 provide additional table structure annotations for the Marmot table and column detectors while the decoders are trained
 data, which currently only has annotations for table detection . independently. Semantic information about elementary data
 types is then utilized to further boost model performance . The
 Index Terms— Table Detection, Table Structure Recognition, utilization of the VGG-19 as a base network, which is pre
 Scanned Documents, Information Extraction trained on the ImageNet dataset allows for exploitation of prior
 knowledge in the form of low level features learnt via training

I. INTRODUCTION

With the proliferation of mobile devices equipped with over ImageNet.
 cameras, an increasing number of customers are uploading We have evaluated TableNet's performance on the ICDAR
 documents via these devices, making the need for information 2013 dataset, demonstrating that our approach marginally
 extraction from these images more pressing. Currently, these outperforms other deep models as well as other state-of-the
 document images are often manually processed resulting in art methods in detecting and extracting tabular information in
 high about costs and inefficient data processing times . Fre- image documents . We further demonstrate that the model can
 quently, these documents contain data stored in tables with generalize to other datasets with minimal fine tuning, thereby
 multiple variations in layout and visual appearance . A key enabling transfer learning. Furthermore, the Marmot dataset
 component of information extraction from these documents which has previously been annotated for table detection was
 therefore involves digitizing the data present in these tabular also manually annotated for column detection, and these new
 sub-images . The variation in the table structure, and in the annotations will be publicly released to the community for
 graphical elements used to visually separate the tabular com- future research .
 ponents make extraction from these images a very challenging In summary, the primary contributions made in this paper
 problem . Most existing approaches to tabular information ex- are as follows:
 traction divide the problem into the two separate sub-problems 1) We propose TableNet: a novel end-to-end deep multi

task architecture for both table detection and structure text block is compared with the average height and if satisfies recognition yielding state of the art performance on the a series of rules, the ROI is regarded as a table .

public benchmark ICDAR and Marmot datasets .

T-Recs (11) was one of the earliest works to extract tabular

- 2) We demonstrate that adding additional spatial semantic data based on clustering of given word segments and overlap features to TableNet during training further boosts model of the text inside the table . Y. Wang et al . [12] estimates prob abilities from geometric measurements made on the various
- 3) We show that using a pre-trained TableNet model and entities in a given document . fine tuning it on an another new dataset will boost the Ashwin et al. (13) exploit the formatting cues from semi performance of the model on the new dataset, thereby structured HTML tables to extract data from web pages. Here allowing for transfer learning. the cells are already demarcated by tags since they are in
- 4) We have manually annotated the Marmot dataset for table HTML tables . Singh et al . [14] use object detection techniques data extraction and will release the annotations to the for Document Layout understanding. community

The rest of the paper is organized as follows: Section II provides an overview of the related work on tabular information extraction . Section III provides a detailed description of the TableNet model. Section IV outlines the extraction process with TableNet. Section V provides details about the datasets , preprocessing steps and training . Section VI outlines the experiment details and results . Finally, the conclusions and future work are presented in Section VII

II. RELATED WORK

There is significant prior work on identifying and extracting both tables and columns have common regions). Therefore, if the tabular data inside a document . Most of these have reported convolutional filters utilized to detect tables, can be reinforced results on table detection and data extraction separately (1) by column detecting filters, this should significantly improve

Before the advent of deep learning, most of the work on the performance of the model . Our proposed model , exploits table detection was based on heuristics or metadata . TINTIN this intuition and is based on the Long et al. (15) , encoder [2] exploited structural information to identify tables and their decoder model for semantic segmentation . The encoder of the component fields. [3] used hierarchical representations based model is common across both tasks, but the decoder emerges on the MXY tree for table detection and was the first attempt at as two different branches for tables and columns . Concretely, using Machine Learning techniques for this problem . T Kasar we enforced the encoding layers to use the ground truth of both et al. [4] identified intersecting horizontal , vertical lines and tables and columns of document for training. However , the low-level features and used an SVM classifier to classify an decoding layers are separated for table and column branches. image region as a table region or not .

Thus, there are two computational graphs to train .

Probabilistic graphical models were also used to detect ta- The input image for the model, is first transformed into bles; Silva et al. (5) modelled the joint probability distribution an RGB image and then, resized to 1024 * 1024 resolution . over sequential observations of visual page elements and the This modified image is processed using tesseract OCR (16) hidden state of a line (HMM) to merge potential table lines as described in the previous section . Since a single model into tables resulted in a high degree of completeness. Jing Fang produces both the output masks for the table and column et al. [6] used the table header as a starting point to detect regions, these two independent outputs have binary target pixel the table region and decompose its elements. Raskovic et al. values, depending on whether the pixel region belongs to the [7] made an attempt to detect borderless tables. They utilized table/column region or background respectively. whitespaces as a heuristic rather than content for detection . The problem of detecting tables in documents is similar to

Recently, DeepDeSRT (8) was proposed which uses deep the problem of detecting objects in real world images . Similar learning for both table detection and table structure recog- to the generic object detection problem, visual features of the nition, i.e. identifying rows, columns, and cell positions in tables can be used to detect tables and columns . The difference the detected tables . This work achieves state-of-the-art per- is that the tolerance for noise in table/column detection is much formance on the ICDAR 2013 table competition dataset . smaller than in object detection. Therefore , instead of regress After this, 19] combined deep convolutional neural networks, ing for the boundaries of tables and columns , we employed a graphical models and saliency concepts for localizing tables method to predict table and column regions pixel-wise . Recent and charts in documents . This technique was applied on an work on semantic segmentation based on pixel wise prediction , extended version of ICDAR 2013 table competition dataset has been very successful. FCN architecture, proposed by Long and outperforms existing models . (10) locates the text compo- et al. (15) , has demonstrated the accuracy of encoder-decoder nents and extracts text blocks. After that , the height of each network architectures for semantic segmentation . The FCN

III. TABLENET: DEEP MODEL FOR TABLE AND COLUMN DETECTION

In all prior deep learning based approaches, table detection and column detection are considered separately as two different problems, which can be solved independently. However, intuitively if all the columns present in a document are table be determined known apriori, the region can easily.

by definition, columns are vertically aligned word/numerical blocks. Thus , independently searching for columns can produce a lot of false positives and knowledge of the tabular region can greatly improve results for column detection (since