

# Tablext: A Combined Neural Network And Heuristic Based Table Extractor

Zach Colter  
Department of EECS  
University of Michigan  
Ann Arbor, US  
zcolter@umich.edu

Morteza Fayazi  
Department of EECS  
University of Michigan  
Ann Arbor, US  
fayazi@umich.edu

Zineb Benameur-El  
Department of EECS  
University of Michigan  
Ann Arbor, US  
zinebbe@umich.edu

Serafina Kamp  
Department of EECS  
University of Michigan  
Ann Arbor, US  
serafibk@umich.edu

Shuyan Yu  
Department of EECS  
University of Michigan  
Ann Arbor, US  
shuyan@umich.edu

Ronald Dreslinski  
Department of EECS  
University of Michigan  
Ann Arbor, US  
rdreslin@umich.edu

**Abstract**—A significant portion of the data available today is found within tables. Therefore, it is necessary to use automated table extraction to obtain thorough results when data-mining. Today’s popular state-of-the-art methods for table extraction struggle to adequately extract tables with machine-readable text and structural data. To make matters worse, many tables do not have machine-readable data, such as tables saved as images, making most extraction methods completely ineffective. In order to address these issues, a novel, open-source, general format table extractor tool, Tablext, is proposed. This tool uses a combination of computer vision techniques and machine learning methods to efficiently and effectively identify and extract data from tables. Tablext begins by using a custom Convolutional Neural Network (CNN) to identify and separate all potential tables. The identification process is optimized by combining the custom CNN with the YOLO object detection network. Then, the high-level structure of each table is identified with computer vision methods. This high-level, structural meta-data is used by another CNN to identify exact cell locations. As a final step, Optical Characters Recognition (OCR) is performed on every individual cell to extract their content without needing machine-readable text. This multi-stage algorithm allows for the neural networks to focus on completing complex tasks, while letting image processing methods efficiently complete the simpler ones. This leads to the proposed approach to be general-purpose enough to handle a large batch of tables regardless of their internal encodings or their layout complexity. Additionally, it becomes accurate enough to outperform competing state-of-the-art table extractors on the ICDAR 2013 table dataset.

**Index Terms**—Table extraction, CNN, deep learning, image processing, OCR

## I. INTRODUCTION

Tables are widely used in documents, articles, web-pages, etc. as they can concisely show complex information in a way that is suitable for human readers [1]. Autonomous table extraction enables translating this abundant amount of information to a machine-readable format which has broad applications in data-mining and information-retrieval [2], [3]. The complexity of some tables’ layouts and the various formats in which they may be found are current problems that prevent the existence of a universal, highly-accurate,

automatic table extraction tool. The purpose of this paper is to introduce a highly accurate extractor that can extract previously unobtainable information.

Formats such as Hypertext Markup Language (HTML), Portable Document Format (PDF), and Portable Network Graphics (PNG) are common encodings that may contain tables. Because of the many different ways to represent a table, designing separate extraction methods for every encoding is time-consuming and prone to errors. Moreover, an environment may have multiple formats that are not encoded uniformly. For instance, a web-page may contain an image that is not encoded in HTML. Unifying internal encodings of files is the first step to overcome universality issues.

In this work, we present an open-source<sup>1</sup> table extractor tool, Tablext, which addresses all the above-mentioned problems. Tablext solves the universality issue by extracting data from images. There are several Java and Python open-source libraries to convert specific file formats to images such as pdf2image [4] (PDF to PNG), GrabzIt [5] (HTML to PNG), etc.

Tablext uses separate Convolutional Neural Networks (CNNs) [6], [7] to identify and extract information from tables, but it does not solely rely on them. Before data extraction, Tablext uses common conventional computer vision techniques, such as line identification, to extract positional information from tables before another CNN is used to fix small mistakes. This approach can be both faster and more accurate than having a single neural network.

By first identifying important features within the tables, each neural network only has to solve a small sub-task. This means that the size required for each network is drastically reduced. The smaller size of the neural networks reduces the computation requirement as well as the risk of overfitting.

Tablext begins by using an efficient CNN to identify the

<sup>1</sup>Source code for the Tablext can be downloaded from [https://github.com/idea-fasoc/datasheet-scrubber/tree/master/src/table\\_extraction](https://github.com/idea-fasoc/datasheet-scrubber/tree/master/src/table_extraction)

exact location of the tables within the input images. Using an ensemble of this CNN and an object detection network, YOLO [8], allows for the extraction algorithm to be confined to relevant areas. By isolating the relatively expensive extraction algorithm to areas found by the identification CNN, Tablext is efficient in handling massive databases.

Tablext achieves a high extraction accuracy by both identifying the lines separating various cells, as well as looking at the positions of the text within the tables. Together, these methods can recognize complex table structures. It can generally be assumed that the identified lines correctly split the table into cells, but when no lines exist or they cannot be correctly identified, the positional data is used to separate the cells. After correctly segmenting the table into cells, the data is cleaned by a CNN before the open-source Optical Character Recognition (OCR) software, Tesseract [9], identifies the text in every individual cell.

The main contributions of this paper can be summarized as follows.

- A new method to identify cells that utilizes the positions of the data and lines in order to obtain high precision preliminary results at a low computational cost.
- A hybrid pipeline that allows the neural network to focus its attention on complex problems, resulting in a high accuracy.
- A novel, general format, open-source tool that identifies and extracts all complex layout tables with a high accuracy within various internal encodings files.

The rest of the paper is organized as follows. Section II briefly presents related works. Section III describes the proposed method for identifying the location of the tables. Section IV explains Tablext's approaches for data extraction. Section V shows the final results and comparisons to other table extractors. Finally, the paper is concluded in Section VI. Throughout this paper, four different tables (Table A, Table B, Table C, Table D) are used to explain various concepts. For clarity, figures will declare which original table they are referencing.

## II. RELATED WORK

Several methods have been studied table identification and extraction. However, some of them describe various approaches to identify tables from different types of documents and do not focus on data extraction [10]–[16].

Several systems have been proposed for the identification and extraction of tables within PDF files [2], [17]–[20]. Oro and Ruffolo [2] find positional data for the text and they use an agglomerative hierarchical clustering algorithm [21] to combine words into lines. The main limitation of this approach is that it only supports unlocked, machine-readable PDFs. However, a large amount of tables are found in locked PDFs or images. Perez-Arriaga et al. [17] find columns and rows by comparing the locations of the text-boxes and combining text-boxes below a distance threshold. Although the method has a decent recall, its false-positive rate is high. Liu et al. [20] propose a method in which, tables are found by grouping text

with similar positions and font size. This approach works with HTML format in addition to PDF files. However, it is only optimized for research papers and its extraction performance is poor especially on general documents. Similar to the other mentioned studies, this method cannot work with locked PDF files or images.

Some other studies focus on table extraction that are not related to PDFs [1], [3], [22]. Tengli et al. [1] propose a machine-learning-based technique for table extraction from HTML files. Nishida et al. [22] focus on HTML extraction too. They use a Recurrent Neural Network (RNN) model to extract data from HTML tags. These such approaches are not applicable to formats other than HTML and are therefore, limited in scope. Pinto et al. [3] propose a method to identify the horizontal lines within tables that are formatted using white space and fixed fonts to align columns in Web pages, however they do not concentrate on table extraction.

Koci et al. and Puha et al. [23], [24] propose methods that can handle tables in more diverse formats. Koci et al. [23] propose a graph representation of tables so that they can use heuristics approaches in order to extract information. Their proposed methodology is optimized to extract data from various spreadsheets. Puha et al. [24] propose a heuristic method for general table extraction. They recognize tables in images by first using OCR to identify data within cells, then using the locations of the data.

Many state-of-the-art general table extraction methods solely use neural networks. Qasim et al. [25] use a Graph Neural Network (GNN) to identify cell locations within tables. DeepDeSRT [26] uses separate, specially made neural networks to identify tables and extract their data, while TableNet [27] uses a single network for both identification and extraction. They argue this is an efficient approach since these tasks are interdependent. However, this also means that the network is less optimized for both tasks. Tablext outperformed both of these methods when tested on a common dataset. A comparison between the related works has been summarized in Table I.

By utilizing Machine Learning (ML), modern table extractors can handle more diverse tables than earlier works. However, identifying and extracting data is a difficult undertaking that requires many sub-tasks to be completed. Many of these tasks, such as identifying lines that can make up a table's structure, are good candidates for conventional computer vision methods. Instead of forcing a neural network to learn and solve every one of these tasks, Tablext introduces a new approach that combines both conventional computer vision and ML techniques. This allows the neural network to be smaller in size and focused on its particular task of identifying abnormal cells. This is especially important given the limited amount of labeled training data available for tables.

## III. TABLE LOCATION IDENTIFICATION

While the main focus of this paper is table extraction, a novel algorithm is created and used to identify tables before extraction. It is imperative to find the table locations before

TABLE I  
COMPARISON OF RELATED WORKS.

| Work             | Table identification | Table extraction | File format  | Approach  |
|------------------|----------------------|------------------|--------------|---|
| [10], [12], [16] | ✓                    | ✗                | All          | Image processing w/o ML                         |
| [11]             | ✓                    | ✗                | PDF          | PdfExtra paradigm                               |
| [13]             | ✓                    | ✗                | All          | DNN   |
| [14]             | ✓                    | ✗                | PDF          | CNN   |
| [15]             | ✓                    | ✗                | All          | SVM   |
| [1]              | ✗                    | ✓                | HTML         | ML  |
| [22]             | ✗                    | ✓                | HTML         | RNN   |
| [23]             | ✗                    | ✓                | Spreadsheets | Graph representation                            |
| [24]             | ✗                    | ✓                | All          | Image processing w/o ML                         |
| [25]             | ✗                    | ✓                | All          | GNN   |
| [2]              | ✓                    | ✓                | PDF          | Agglomerative hierarchical clustering algorithm |
| [17]             | ✓                    | ✓                | PDF          | Comparing the locations of the text-boxes       |
| [20]             | ✓                    | ✓                | PDF & HTML   | Grouping text                                   |
| [26]             | ✓                    | ✓                | All          | NN  |
| [27]             | ✓                    | ✓                | All          | NN  |
| Tablext          | ✓                    | ✓                | All          | Heuristic-based + CNN                           |

extraction because doing so increases both the performance and accuracy of the algorithm. The boost in performance comes from reduced search space for the complex information extraction. The increased accuracy is a result of the table identification network separating tables from one another and removing irrelevant text and lines that could potentially confuse the table extractor. To further increase the accuracy of the table identification, it is later combined with YOLO.

The identification algorithm needs to find both the minimum and maximum  $X$  and  $Y$  coordinates for each table that exists within each page. To efficiently calculate these coordinates, two relatively simple CNNs are used in succession. These two networks identify at which  $Y$  and  $X$  coordinates, respectively, a table exists.

All input images, after being converted from any file format, are first converted to gray-scale and resized. Preliminary tests show that a width of 800 pixels allows for accurate table extraction without a significant performance penalty. The height at this stage is scaled proportionally to avoid distortion. This resized image cannot directly be sent to a conventional CNN because its height is variable. Additionally, the large size of the input data would require a significant amount of memory for each batch, slowing down the training process. In order to reduce the input size and make the input regular, the image is sliced into horizontal strips, each with a height of 64 pixels. This number is chosen to reduce the training time without significantly degrading accuracy.

While cutting the image into strips reduces the input and

output size of the network, this process introduces a new problem. The accuracy of identifying a table at the top and the bottom of the slices is lower than identifying a table near the middle. Intuitively this is because the network only has information in one direction from these edge outputs. To combat this issue, each strip only detects the existence of a table within its innermost 32 rows of pixels. To cover the entire table, a moving window approach is used where each input strip partially overlaps its adjacent strips.

The neural network has 4 output nodes for every horizontal strip. Each node determines if there is a table at a particular location within the centermost 32 pixels. In the innermost 32 pixels, node  $i, 0 \leq i \leq 3$ , of the neural network is trained to identify if a table exists between any of the  $[8 * i, 8 * i + 7]$  pixels. The outputs of all the slices are concatenated and stored in an array that identifies in which rows a table exists. With this approach, the top and bottom of each page cannot be searched for tables, because only the inner pixels are represented in the outputs. To fix this problem, the input image is extended by 16 pixels at the top and bottom with a white border.

Fig. 1 shows the concatenated output of all the slices from the neural network. In this figure, the thin horizontal boxes represent every row the network checked for a table. These boxes show the  $Y$  locations where the neural network predicts a table exists. While this particular example shows a table that takes up the entire width of the document, this network can also identify tables that are within one of several columns on the page.

The next step is to clean the output by concatenating rows of positively identified tables into groups. If no groups are created, because none of the rows contain a table, Tablext outputs that no table is found in the image. Otherwise, the groups are prepared and sent to the second neural network in order to find the correct  $X$  coordinates. This process is simpler than having to find the correct  $Y$  coordinates as the amount of area to search is smaller. For this reason, a CNN is used without a sliding window approach.

Each group of rows is resized into a  $400 \times 400$  pixel square to efficiently achieve high accuracy. This input size is obtained by preliminary testing. The image is sent into a CNN that identifies which columns contain a table. Fig. 2 shows this technique being applied to the bottom group of rows in Fig. 1. The actual result produced by the network looks distorted to a human because the table is scaled into a square. Instead a representative, recreated image is shown. Here the vertical boxes show where the network predicts that a table exists. If the table is only in a single column of a multi-column image, the lines would cover only that section. This network can also identify two separate adjacent tables. This is achieved by not identifying a table in between the two table groups. Each vertical group is separately sent to this second stage. As an example, Fig. 1 would send three separate images into the second network.

The columns are concatenated into groups, in a similar manner to the rows. The scaling that converted the group of rows into the square, shown in Fig. 2, is inverted and applied

to the  $X$  coordinates. Additionally, the scaling used to convert the original image into the resized image, shown in Fig. 1, is inverted and applied to the  $Y$  coordinates. Doing this gives accurate coordinates of each table's start and end position regarding the original image. Fig. 3 shows a cutout of the first table, using the found  $X$  and  $Y$  coordinates.

The CNN by itself tends to over propose table regions because the existence of lines on the page are very influential to the model. This causes the CNN to predict areas around these lines as tables whether these lines correspond to tables or not. This will cause the extraction portion of the algorithm to take more time as it has to run OCR on these extra regions to figure out they are not tables. To counteract this, a YOLO [8] network is trained on table images to see if this more popular object detection method could outperform the CNN. YOLO typically identifies all the right table regions, but the bounding boxes that it proposes do not cover the entire table. For this reason, an ensemble of the CNN and YOLO is used to create an optimal identification model. Both networks first independently propose table regions. When these models propose regions that overlap, the union of these two regions is taken. This ensures both that the final regions are likely to contain tables and that the regions will likely contain the entirety of the tables.

While this approach can perfectly identify tables in multiple columns, there is an extremely rare case where two adjacent tables in separate columns have vastly different heights. In this case, extra data could be identified as a table. This rare issue is rectified later because the proposed extraction method can tolerate some extra area identified as a table.

#### IV. TABLE EXTRACTION

After all of the tables are identified, the extraction method starts independently on each table. In order to process high resolution images efficiently, a resized copy of each individual table is created with a width of 800 pixels and a proportional height. This copy is used for a majority of the cell identification steps, however the original image is kept because it is used later during OCR.

Several methods are used to identify and isolate the cells within each table. The table extraction begins by identifying the high-level, general structure of the table. It identifies the vertical and horizontal lines, both visible and implied, that correctly divide a majority of the cells. There are several ways to locate this information within the tables. One approach is to identify the visible lines separating various cells. Another method is to look at the positions of the text within the tables and extrapolate how the cells should be positioned in accordance to each other. Tabletext uses both of these approaches in order to obtain the structure of the table. Afterwards, a neural network is used to clean up any cells that do not conform to the general structure. Fig. 4 shows a table that already went through table identification. This will be used as an example to explain the steps of the table extraction.

## MCP3021

### DC ELECTRICAL SPECIFICATIONS (CONTINUED)

| Electrical Characteristics: Unless otherwise noted, all parameters apply at $V_{DD} = 5.0V$ , $V_{SS} = GND$ , $R_{PU} = 2 k\Omega$ , $T_A = -40^\circ C$ to $+85^\circ C$ , $PC$ Fast Mode Timing, $f_{SCL} = 400 kHz$ (Note 3) |            |      |       |      |         |   |
|--|------------|------|-------|------|---------|---|
| Parameter  | Sym.       | Min. | Typ.  | Max. | Units   | Conditions                                  |
| Input Leakage Current  | $I_{LJ}$   | -1   | —     | +1   | $\mu A$ | $V_{IN} = V_{SS}$ to $V_{DD}$               |
| Output Leakage Current   | $I_{LO}$   | -1   | —     | +1   | $\mu A$ | $V_{OUT} = V_{SS}$ to $V_{DD}$              |
| Pin Capacitance (all inputs/outputs)   | $C_{IN}$   | —    | —     | 10   | pF      | $T_A = 25^\circ C$ , $f = 1 MHz$ , (Note 2) |
| Bus Capacitance  | $C_B$      | —    | —     | 400  | pF      | SDA drive low, 0.4V                         |
| <b>Power Requirements</b>  |            |      |       |      |         |   |
| Operating Voltage  | $V_{DD}$   | 2.7  | —     | 5.5  | V       |   |
| Conversion Current   | $I_{DD}$   | —    | 175   | 250  | $\mu A$ |   |
| Standby Current  | $I_{DDS}$  | —    | 0.005 | 1    | $\mu A$ | SDA, SCL = $V_{DD}$                         |
| Active Bus Current   | $I_{DDA}$  | —    | —     | 120  | $\mu A$ | Note 4                                      |
| <b>Conversion Rate</b>   |            |      |       |      |         |   |
| Conversion Time  | $t_{CONV}$ | —    | 0.96  | —    | $\mu s$ | Note 5                                      |
| Analog Input Acquisition Time  | $t_{ACQ}$  | —    | 1.12  | —    | $\mu s$ | Note 5                                      |
| Sample Rate  | $f_{SAMP}$ | —    | —     | 22.3 | ksps    | $f_{SCL} = 400 kHz$ (Note 1)                |

- Note 1: Sample time is the time between conversions after the address byte has been sent to the converter. Refer to Figure 5-6.  
 Note 2: This parameter is periodically sampled and not 100% tested.  
 Note 3:  $R_{PU}$  = Pull-up resistor on SDA and SCL.  
 Note 4: SDA and SCL =  $V_{SS}$  to  $V_{DD}$  at 400 kHz.  
 Note 5:  $t_{ACQ}$  and  $t_{CONV}$  are dependent on internal oscillator timing. See Figure 5-5 and Figure 5-6 in relation to SCL.

### TEMPERATURE SPECIFICATIONS

| Electrical Characteristics: Unless otherwise noted, all parameters apply at $V_{DD} = 5.0V$ , $V_{SS} = GND$ |               |      |      |      |              |            |
|--|---------------|------|------|------|--------------|------------|
| Parameter  | Sym.          | Min. | Typ. | Max. | Units        | Conditions |
| <b>Temperature Ranges</b>  |               |      |      |      |              |            |
| Operating Temperature Range  | $T_A$         | -40  | —    | +125 | $^\circ C$   |            |
| Extended Temperature Range   | $T_A$         | -40  | —    | +125 | $^\circ C$   |            |
| Storage Temperature Range  | $T_A$         | -65  | —    | +150 | $^\circ C$   |            |
| <b>Thermal Package Resistances</b>   |               |      |      |      |              |            |
| Thermal Resistance, SOT-23   | $\theta_{JA}$ | —    | 256  | —    | $^\circ C/W$ |            |

Fig. 1. (Table A) Each horizontal boxes represents a row where a table is identified. The bottom group is highlighted with a thick box (red) for clarification.

| Electrical Characteristics: Unless otherwise noted, all parameters apply at $V_{DD} = 5.0V$ , $V_{SS} = GND$ |               |      |      |      |              |            |
|--|---------------|------|------|------|--------------|------------|
| Parameter  | Sym.          | Min. | Typ. | Max. | Units        | Conditions |
| <b>Temperature Ranges</b>  |               |      |      |      |              |            |
| Operating Temperature Range  | $T_A$         | -40  | —    | +125 | $^\circ C$   |            |
| Extended Temperature Range   | $T_A$         | -40  | —    | +125 | $^\circ C$   |            |
| Storage Temperature Range  | $T_A$         | -65  | —    | +150 | $^\circ C$   |            |
| <b>Thermal Package Resistances</b>   |               |      |      |      |              |            |
| Thermal Resistance, SOT-23   | $\theta_{JA}$ | —    | 256  | —    | $^\circ C/W$ |            |

Fig. 2. (Table A) Vertical boxes represent every column where a table is identified.

| Electrical Characteristics: Unless otherwise noted, all parameters apply at $V_{DD} = 5.0V$ , $V_{SS} = GND$ . |               |      |      |      |              |            |
|--|---------------|------|------|------|--------------|------------|
| Parameter  | Sym.          | Min. | Typ. | Max. | Units        | Conditions |
| <b>Temperature Ranges</b>  |               |      |      |      |              |            |
| Operating Temperature Range  | $T_A$         | -40  | —    | +125 | $^\circ C$   |            |
| Extended Temperature Range   | $T_A$         | -40  | —    | +125 | $^\circ C$   |            |
| Storage Temperature Range  | $T_A$         | -65  | —    | +150 | $^\circ C$   |            |
| <b>Thermal Package Resistances</b>   |               |      |      |      |              |            |
| Thermal Resistance, SOT-23   | $\theta_{JA}$ | —    | 256  | —    | $^\circ C/W$ |            |

Fig. 3. (Table A) A cutout of the original table using the coordinates identified.

| Based on the Offer Price of HK\$11.10 (being the Minimum Offer Price) |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| Investor  | Investment Amount<br>(US\$ in million) <sup>f</sup> | Number of Offer Shares<br>(rounded down to nearest whole board lot of 500 Shares) | Approximate % of total number of Offer Shares       |   | Approximate % of total Shares in issue immediately following the completion of Capitalization Issue and Global Offering |   |
|   |   |   | Assuming the Over-allotment Option is not exercised | Assuming the Over-allotment Option is exercised in full | Assuming the Over-allotment Option is not exercised   | Assuming the Over-allotment Option is exercised in full |
| Indus Funds GIC   | 15  | 10,603,500  | 5.69%   | 4.95%   | 1.08%   | 1.05%   |
| Tetrad Ventures Pte Ltd   | 30  | 21,207,500  | 11.38%  | 9.89%   | 2.16%   | 2.10%   |
| GIC Private Limited   | 10  | 7,069,000   | 3.79%   | 3.30%   | 0.72%   | 0.70%   |

Fig. 4. (Table B) A table sample identified by the proposed table identification method.

### A. Real Line Identification

Vertical and horizontal lines are commonly used to separate and organize the contents of tables. To identify these lines, the second-derivative gradient field is approximated for the gray-scale image of the table. Assuming that  $f(x, y)$  represents the brightness of a pixel at the location  $(x, y)$ . The values of  $\frac{\partial^2}{\partial x^2} f(x, y)$  and  $\frac{\partial^2}{\partial y^2} f(x, y)$  are approximated for each pixel by looking at their adjacent pixels and the gradients of them. These gradients are then stored into separate matrices. Then a  $3 \times 3$  max pooling is independently applied to these matrices, to account for noise in the image and imperfect lines. The max pooling has a stride of 1 in both  $X$  and  $Y$  directions to retain as much precision as possible when identifying the lines.

With this information, the vertical lines are found by searching for vertical segments where all the pixels in the segment have high  $\frac{\partial^2}{\partial x^2} f(x, y)$  values relative to the pixels in rest of the image and similar  $\frac{\partial^2}{\partial y^2} f(x, y)$  values to one another. Similarly, horizontal lines are found by searching for horizontal segments where all the pixels in the segment have high  $\frac{\partial^2}{\partial y^2} f(x, y)$  values relative to the pixels in rest of the image and similar  $\frac{\partial^2}{\partial x^2} f(x, y)$  values to one another.

The segment length has to be a small fraction of the table's height or width for the vertical and horizontal lines respectively. However, once a line is found, that line is extended to cover the entire width or height of the table. Identifying lines in this way helps define the structure of the table, but some cells that do not conform to the general format might be split by lines that do not cover the entire width or height of the table. This issue is solved later with a neural network that stitches together improperly cut cells. Fig. 5 depicts a debug image, created by the real line identification algorithm, that shows the few real lines in dark blue. The rest of the lines in this figure are created with another process that will be explained in the following section.

### B. Inferred Line Identification

Real line identification alone is sufficient for regular tables that have all of their data properly split into cells with defined borders. However, as Fig. 4 shows, many tables only have

| Based on the Offer Price of HK\$11.10 (being the Minimum Offer Price) |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| Investor  | Investment Amount<br>(US\$ in million) <sup>f</sup> | Number of Offer Shares<br>(rounded down to nearest whole board lot of 500 Shares) | Approximate % of total number of Offer Shares       |   | Approximate % of total Shares in issue immediately following the completion of Capitalization Issue and Global Offering |   |
|   |   |   | Assuming the Over-allotment Option is not exercised | Assuming the Over-allotment Option is exercised in full | Assuming the Over-allotment Option is not exercised   | Assuming the Over-allotment Option is exercised in full |
| Indus Funds GIC   | 15  | 10,603,500  | 5.69%   | 4.95%   | 1.08%   | 1.05%   |
| Tetrad Ventures Pte Ltd   | 30  | 21,207,500  | 11.38%  | 9.89%   | 2.16%   | 2.10%   |
| GIC Private Limited   | 10  | 7,069,000   | 3.79%   | 3.30%   | 0.72%   | 0.70%   |

Fig. 5. (Table B) Inferred lines shown with thick bright (yellow) lines and real lines shown with thin dark (blue) lines.

| Parameter   | Symbol            | Min  | Typ | Max | Unit |
|---|-------------------|------|-----|-----|------|
| CONVERSION TIME—CNV RISING EDGE TO DATA AVAILABLE | t <sub>conv</sub> | 270  | 290 | 320 | ns   |
| ACQUISITION PHASE <sup>1</sup>                    | t <sub>acq</sub>  |      |     |     |      |
| AD4002  |                   | 290  |     |     | ns   |
| AD4006  |                   | 790  |     |     | ns   |
| AD4010  |                   | 1790 |     |     | ns   |
| TIME BETWEEN CONVERSIONS                          | t <sub>tw</sub>   |      |     |     |      |
| AD4002  |                   | 500  |     |     | ns   |
| AD4006  |                   | 1000 |     |     | ns   |
| AD4010  |                   | 2000 |     |     | ns   |
| CNV PULSE WIDTH (CS MODE) <sup>2</sup>            | t <sub>cnv</sub>  | 10   |     |     | ns   |
| SCK PERIOD (CS MODE) <sup>3</sup>                 | t <sub>sck</sub>  |      |     |     |      |
| VIO > 2.7 V                                       |                   | 9.8  |     |     | ns   |

Fig. 6. (Table C) Initial sparse table before inferred vertical lines. A red box has been placed around the, difficult to handle, referenced sparse column for clarity.

proper borders for some of the cells or do not have any well defined lines at all. To correctly identify cells that real line identification alone could not, Tablext uses the positions of the data. Finding the locations of the data is trivial. All of the high contrast pixels, not already identified as a real line, belong to the data. With the locations of the data known, Tablext attempts to segment the pieces of data into columns and rows. To achieve this, Tablext creates straight lines that do not intersect the data. These such lines will be referred to “inferred lines” as the table requires these non-existent lines to properly convey its information. Fig. 5 shows when all inferred and real lines of Fig. 4 are drawn.

The task of finding inferred lines is made harder by the existence of cells that do not conform to the general table's structure. In Fig. 4, multiple cells span two columns. If inferred lines are defined as lines that do not overlap any pieces of data, the columns underneath the cells that span multiple columns would not be separated. In order to solve this issue, Tablext starts by defining the term “threshold distance”.

For vertical lines, the threshold distance is the proportion of the table's height an inferred line must not touch any data in order to be considered valid. If the table is dense, meaning that most cells contain data, a small, static threshold distance would produce accurate results. However if the table is sparse, meaning that many cells do not contain data, this approach

| Parameter   | Symb              | Min  | Typ | Max | Un |
|---|-------------------|------|-----|-----|----|
| CONVERSION TIME—CNV RISING EDGE TO DATA AVAILABLE | t <sub>conv</sub> | 270  |     |     | ns |
| ACQUISITION PHASE <sup>1</sup>                    | t <sub>acc</sub>  |      |     |     |    |
| AD4002  |                   | 290  |     |     | ns |
| AD4006  |                   | 790  |     |     | ns |
| AD4010  |                   | 1790 |     |     | ns |
| TIME BETWEEN CONVERSIONS                          | t <sub>bc</sub>   |      |     |     |    |
| AD4002  |                   | 500  |     |     | ns |
| AD4006  |                   | 1000 |     |     | ns |
| AD4010  |                   | 2000 |     |     | ns |
| CNV PULSE WIDTH (CS MODE) <sup>2</sup>            | t <sub>cnv</sub>  | 10   |     |     | ns |
| SCK PERIOD (CS MODE) <sup>3</sup>                 | t <sub>sck</sub>  | 9.8  |     |     | ns |
| V <sub>IO</sub> > 2.7 V                           |                   |      |     |     |    |

Fig. 7. (Table C) Sparse table with initialized inferred vertical lines.

would not work. A small, static threshold distance would combine adjacent empty cells, into wide inferred lines. Fig. 6 and Fig. 7 show the inferred vertical lines initialization. Fig. 7 illustrates that a few pieces of data within these columns or rows could be overwritten by the inferred lines.

To accommodate both sparse and dense tables, the threshold distance is adaptable. Algorithm 1 describes a simplified procedure of acquiring the final inferred lines with a variable threshold distance. In the algorithm, possible inferred lines refers to inferred lines found at the given threshold distance,  $\Delta$  is a small learned constant, and number of groups refers to the amount of separable groups of inferred lines that exist within the possible inferred lines.

#### Algorithm 1: Setting the threshold distance

```

1: Variables: num_group1, num_group2, thresh-
   old_distance, temp_lines
2: Output: final_lines
3: Start:
4:   threshold distance = 0.6; num_group1 = 0;
   final_lines = [];
5:   Loop:
6:     temp_lines = possible inferred lines;
7:     num_group2 = number of groups;
8:     if num_group2  $\geq$  num_group1:
9:       num_group1 = num_group2;
10:      final_lines = temp_lines;
11:      threshold_distance +=  $\Delta$ ;
12:      go to Loop;
13:   return final_lines;

```

For a visual representation of this algorithm, images of this process are produced by the code. Fig. 8 shows the final output after applying this algorithm to the sparse table in Fig. 7. At the start, the vertical inferred lines, represented as black columns, completely cover the column beginning with “Typ”. If the threshold distance was a low static number, this column would be merged with one of its neighboring columns. However, by using Algorithm 1, the lines, over several iterations, split around the “Typ” column as shown by Fig. 8.

During this process, every inferred line is assigned a “quality score”. The quality score is the percentage of the maximum possible line length that a particular inferred line

| Parameter   | Symb              | Min  | Typ | Max | Un |
|---|-------------------|------|-----|-----|----|
| CONVERSION TIME—CNV RISING EDGE TO DATA AVAILABLE | t <sub>conv</sub> | 270  | 290 | 320 | ns |
| ACQUISITION PHASE <sup>1</sup>                    | t <sub>acc</sub>  |      |     |     |    |
| AD4002  |                   | 290  |     |     | ns |
| AD4006  |                   | 790  |     |     | ns |
| AD4010  |                   | 1790 |     |     | ns |
| TIME BETWEEN CONVERSIONS                          | t <sub>bc</sub>   |      |     |     |    |
| AD4002  |                   | 500  |     |     | ns |
| AD4006  |                   | 1000 |     |     | ns |
| AD4010  |                   | 2000 |     |     | ns |
| CNV PULSE WIDTH (CS MODE) <sup>2</sup>            | t <sub>cnv</sub>  | 10   |     |     | ns |
| SCK PERIOD (CS MODE) <sup>3</sup>                 | t <sub>sck</sub>  | 9.8  |     |     | ns |
| V <sub>IO</sub> > 2.7 V                           |                   |      |     |     |    |

Fig. 8. (Table C) Sparse table with incremented inferred vertical lines.

covers without intersecting a piece of data. These quality scores are used later for combining the various lines.

Inferred horizontal lines are created in a similar manner as inferred vertical lines, except the threshold distance is a percentage of the table’s width instead of the table’s height. To better show this complete process, both vertical and horizontal inferred line identification are used on Fig. 4. The resulting inferred lines are used in conjunction with previously calculated real lines to create the debug image shown in Fig. 5.

#### C. Final Structural Line Identification

Now that all of the real and inferred lines have been located, Tablext begins to group all of this information so that it has a high level understanding of the particular table’s layout. Inferred lines are combined with their neighbors, if they are within close proximity, to create different groups of inferred lines. If a real line is within a group of inferred lines, or it is only several pixels away from a group, only the real line is kept. The real lines are considered to be the ground truth, so nearby inferred line groups are redundant. If there are no real lines near an inferred line group, the quality scores of the inferred lines are used to determine which line in the group is the best line to split the cells. A Simple Moving Average (SMA) of the quality scores, looking both forwards and backwards two pixels, is calculated for each group. Lines that do not have a valid inferred line are given a zero quality score. Using an SMA helps reduce the impact of noise in the image, while also fixing corner cases where a single pixel thick inferred line might sneak between characters in a piece of data.

The index that contains the maximum SMA value is taken for each inferred group which is not near a real line. Then, the locations of the real lines and the maximum SMA indexes are concatenated together. These combined lines are the final lines for the high level table identification. This process is completed independently for both the vertical and horizontal lines. Fig. 9 is created by applying this process on the lines found in Fig. 5.

#### D. Neural Network Correction

Since the final lines describe the regular structure of the table, cells that do not conform to that structure might be sliced into two or more pieces. Fig. 9 shows several cells that are cut in half by the final lines. In order to fix this issue, a CNN is used to recombine improperly sliced cells. A simpler



|                         |                                | Based on the Offer Price of HK\$11.10 (being the Minimum Offer Price) |   |   |   |   |
|-------------------------|--------------------------------|---|---|---|---|---|
|                         |                                | Number of Offer Shares  |   | Approximate % of total Shares in issue immediately following the completion of Capitalization Issue and Global Offering |   |   |
|                         |                                | (rounded down to nearest whole board lot of 500 Shares)               | Assuming the Over-allotment Option is not exercised | Assuming the Over-allotment Option is exercised in full   | Assuming the Over-allotment Option is not exercised | Assuming the Over-allotment Option is exercised in full |
| Cornerstone Investor    | Investment Amount              |   |   |   |   |   |
|                         | (US\$ in million) <sup>f</sup> |   |   |   |   |   |
| Indus Funds GIC         | 15                             | 10,603,500  | 5.69%   | 4.95%   | 1.08%   | 1.05%   |
| Tetrad Ventures Pte Ltd | 30                             | 21,207,500  | 11.38%  | 9.89%   | 2.16%   | 2.10%   |
| GIC Private Limited     | 10                             | 7,069,000   | 3.79%   | 3.30%   | 0.72%   | 0.70%   |

Fig. 9. (Table B) The merged lines (thin green) created by the final line identification. Vertically split cells highlighted with thick boxes (red).

solution would be to check if any high contrast pixels exist underneath an inferred line. However, this method has several problems. One issue is if characters in a cell are spaced out significantly, a vertical inferred line could exist in between two characters and the cell's contents would be split into two cells.

While these concerns can be partially alleviated by looking a certain distance to the left and right for high contrast pixels, this introduces a new risk of merging cells unnecessarily. Additionally, merging with that method can be triggered by noise in the tables. For these reasons, a neural network has been created to handle the possible merging of cells.

The CNN takes in two adjacent cells as inputs and resizes them to be  $100 \times 100$  pixels each. Testing shows that this size allows the OCR to accurately extract data. The cells are resized independently so that the line that divides them is always in the center of the two cells. If instead, the cells are merged and then resized, the CNN would need to predict which data belongs to which cell or it would require the location of the dividing line as an additional input. These approaches result in significantly lower accuracy and higher complexity, respectively.

The proposed CNN outputs the chance that these two cells are merged and if one or both of the cells are empty. By using a neural network to compute this information, the effect of imperfections in the image can be greatly reduced. It is important to know which cells are empty to avoid running OCR on those cells. This increases performance, as the OCR is the most time consuming step. The following sections talk about the implementation details of the CNN, the training data, and the post CNN concatenation.

1) *CNN Construction:* The two cells that have been sent into the CNN are merged and rotated so that the potential improper line separating these two be always a vertical line in the center of the merged image. The merged image, as the CNN sees it, always has a dimension of  $200 \times 100$ .

The CNN has three branches that are concatenated into a dense network. The first branch has three 2D convolution

|                         |                                | Based on the Offer Price of HK\$11.10 (being the Minimum Offer Price) |   |   |   |   |
|-------------------------|--------------------------------|---|---|---|---|---|
|                         |                                | Number of Offer Shares  |   | Approximate % of total Shares in issue immediately following the completion of Capitalization Issue and Global Offering |   |   |
|                         |                                | (rounded down to nearest whole board lot of 500 Shares)               | Assuming the Over-allotment Option is not exercised | Assuming the Over-allotment Option is exercised in full   | Assuming the Over-allotment Option is not exercised | Assuming the Over-allotment Option is exercised in full |
| Cornerstone Investor    | Investment Amount              |   |   |   |   |   |
|                         | (US\$ in million) <sup>f</sup> |   |   |   |   |   |
| Indus Funds GIC         | 15                             | 10,603,500  | 5.69%   | 4.95%   | 1.08%   | 1.05%   |
| Tetrad Ventures Pte Ltd | 30                             | 21,207,500  | 11.38%  | 9.89%   | 2.16%   | 2.10%   |
| GIC Private Limited     | 10                             | 7,069,000   | 3.79%   | 3.30%   | 0.72%   | 0.70%   |

Fig. 10. (Table B) The table after the CNN merged problematic cells. Rows kept are separated for accurate OCR.

layers with 16 filters with a size of  $3 \times 3$ , each layer is directly followed by a max pooling layer of size  $2 \times 2$ . This branch's main purpose is to give low fidelity information around the possibly improper line that splits the cells, but it also helps identify which cells contain data.

The second branch begins with an asymmetric average pooling layer of size  $1 \times 100$ . The output of this layer has the dimension of  $200 \times 1$ . This layer is followed by five 1D convolution layers with four filters each with a size of 3. The average pooling gives a good noise resistant indication if any data exists within the particular columns. This information helps identify the existence of data within one or both of the cells, but also helps confirm the presence of an improper split of the two cells.

The third and final branch does not take in the full merged image like the first two, instead it takes in the centermost  $20 \times 100$  pixels. The exclusive purpose of this branch is to identify an improper splitting of the cells. Intuitively, the pixels where the cells are merged and the surrounding pixels contain the most crucial information regarding whether the cells should be merged. More complex neural network layers are used in this branch due to the importance of the information and the limited amount of inputs. Five convolutions layers, each with 64 filters of size  $3 \times 3$ , are stacked upon one another.

All of these branches are individually passed through a dropout layer with a 20% dropout rate before they are flattened and concatenated. The concatenated data is sent into a dense layer of size 256, which is followed by another dropout layer, with a dropout rate of 50% and another dense layer of equal size. This second dense layer is then connected to the final output which contains 3 nodes. The first two nodes signify if the two cells contain data and the third node is high when two cells should be merged. The first dropout layers help guide the network away from relying purely on a single branch. Additionally, all of the dropout layers help avoid overfitting that would otherwise be present due to the limited amount of training data.

| Based on the Offer Price of HK\$11.10 (being the Minimum Offer Price) |                                      | < EXTEND   | < EXTEND  | < EXTEND  | < EXTEND  |   |
|---|--------------------------------------|--|---|---|---|---|
|   |                                      |  |   |   | Approximate % of total Shares in issue immediately following the completion of Capitalization Issue and Global Offering | < EXTEND  |
|   |                                      | Number of Offer Shares (rounded down to nearest whole board lot of 500 Shares) | Approximate % of total number of Offer Shares       | < EXTEND  | < EXTEND  | < EXTEND  |
|   |                                      |  | Assuming the Over-allotment Option is not exercised | Assuming the Over-allotment Option is exercised in full | Assuming the Over-allotment Option is not exercised   | Assuming the Over-allotment Option is exercised in full |
| Cornerstone Investor  | Investment Amount (US\$ in million)* |  |   |   |   |   |
| Indus Funds GIC   |                                      | 15   | 10,603,500 5.69%                                    |   | 4.95%   | 1.08%   |
| Tetrad Ventures Pte Ltd   |                                      | 30   | 21,207,500  | 11.38%  | 9.89%   | 2.16%   |
| GIC Private Limited   |                                      | 10   | 7,069,000   | 3.79%   | 3.30%   | 0.72%   |

Fig. 11. (Table B) Tablext output of the table in Fig. 4, in CSV format. Thick boxes (red) show the horizontally merged cells. Cells with multiple rows are kept separated for readability.

2) *CNN Training*: To train this CNN, a dataset that includes over 1000 tables from various domains is procured and manually annotated. The ground truth of every table’s location, along with its respective cells’ locations, is stored within Extensible Markup Language (XML) files. To translate this dataset into useful data for training purposes, Tablext’s line identification techniques are used to split up the cells. All possible pairs of adjacent cells are then merged and these merged cells are used as the training data.

To create correct labels for training, the cells’ locations are used. If a cell that contains data, according to the XML data, exists on either half of the merged cell, the respective half is considered to have data. When the XML data signifies that a single cell covers some part of both halves of the merged cell, the label for concatenation is high.

3) *Post CNN Concatenation*: The outputs of the CNN fill in matrices that record which cells should be merged and which cells contain data. Multiple cells in both the vertical and horizontal directions can be combined, if necessary, into a single large merged cell. Fig. 10 is a debug image that shows this result. Cells with multiple independent horizontal lines are not combined at this stage, so that the OCR is fed a single line of data for having a higher accuracy. These horizontal lines can be merged later to allow for a more concise table, or left separated to better represent a table’s proportions. With the cells properly defined, the data is almost ready for OCR. However, the pixel resolution used by the table extractor, 800 pixels, is generally not sufficient to obtain accurate results from OCR. To address this problem, the lines defining the cell boundaries are scaled and used to cut cells into the original image.

#### E. OCR

All the non-empty, individual cells are sent to Tesseract, an open-source OCR tool. The returned text is then stored in a Column Separated Values (CSV) format with the location provided by Tablext. The output is shown in Fig. 11. All the data within an extended cell is placed within a single cell in the CSV. To allow for accurate data interpretation, all of the merged cells point back to that original data by using the keyword “EXTEND” and a directional arrow. These cells are

emphasized with boxes in Fig. 11. For clarity, cells that take up multiple horizontal lines are not merged in the CSV.

## V. EVALUATION

Unlike Tablext, most state-of-the-art table extraction papers do not have their code available for download. However, both DeepDeSRT [26] and TableNet [27] have tested their design on the ICDAR 2013 table competition dataset [28]. To compare Tablext quantitatively to these works, Tablext is also tested against this dataset. Additionally, Tablext’s extraction capability is compared against the popular, open-source, PDF extractor Tabula [29] on a manually procured diverse dataset with over 100 tables. Although table identification is not the focus of this paper, the proposed method produces high quality results, while only using a small fraction of the total runtime.

### A. Table identification results

The table identifier is tested on a dataset with over 400 tables from various domains. This dataset includes many abstract tables without any lines that are much harder to identify than most tables.

Let AP be the area predicted in an image and AL be the true occupied area. Precision and recall, for an image, can be calculated by  $\text{precision} = \frac{AP \cap AL}{AP}$  and  $\text{recall} = \frac{AP \cap AL}{AL}$ . The results of the custom CNN, the YOLO model, and a combination of the custom CNN and YOLO model have been summarized in Table II. After averaging the precision and recall across the dataset, with each image weighted equally, final results of 86% and 87% are achieved for precision and recall respectively for the combined model.

Although the YOLO only model achieves higher precision, the higher recall of the combination of the custom CNN and YOLO is preferable since it is better to over propose regions rather than to miss any tables. This is because regions that are not tables can be discarded later in the extraction process while completely missed tables cannot be recovered.

### B. Table extraction result comparison

1) *ICDAR comparison*: Both DeepDeSRT and TableNet use 34 randomly selected images from the ICDAR dataset for their testing set. Then, they calculate precision and recall in



| AD7683  |                        |         |      |      |         |      |      |      |
|---|------------------------|---------|------|------|---------|------|------|------|
| VDD = 5 V; VREF = VDD; TA = -40°C to +85°C, unless otherwise noted. |                        |         |      |      |         |      |      |      |
| Table 3.  |                        |         |      |      |         |      |      |      |
| Parameter   | Conditions             | A Grade |      |      | B Grade |      |      | Unit |
|   |                        | Min     | Typ  | Max  | Min     | Typ  | Max  |      |
| <b>ACCURACY</b>   |                        |         |      |      |         |      |      |      |
| No Missing Codes  |                        | 15      |      |      | 16      |      |      | Bits |
| Integral Linearity Error  |                        | -6      | ±3   | +6   | -3      | ±1   | +3   | LSB  |
| Transition Noise  |                        |         | 0.5  |      |         | 0.5  |      | LSB  |
| Gain Error <sup>1</sup> , Twin to Tmax                              |                        |         | ±2   | ±24  |         | ±2   | ±15  | LSB  |
| Offset Error, Twin to Tmax  |                        |         | ±0.7 | ±1.6 |         | ±0.4 | ±1.6 | mV   |
| <b>AC ACCURACY</b>  |                        |         |      |      |         |      |      |      |
| Total Harmonic Distortion   | f <sub>m</sub> = 1 kHz |         | -100 |      |         | -106 |      | dB   |

<sup>1</sup> See the Terminology section. These specifications include full temperature range variation, but do not include the error contribution from the external reference.  
<sup>2</sup> All specifications in dB are referred to a full-scale input, FS. Tested with an input signal at 0.5 dB below full scale, unless otherwise specified.

VDD = 2.7 V; VREF = 2.5V; TA = -40°C to +85°C, unless otherwise noted.

| Table 4.                               |                        |         |      |      |         |      |      |      |
|--|------------------------|---------|------|------|---------|------|------|------|
| Parameter                              | Conditions             | A Grade |      |      | B Grade |      |      | Unit |
|  |                        | Min     | Typ  | Max  | Min     | Typ  | Max  |      |
| <b>ACCURACY</b>                        |                        |         |      |      |         |      |      |      |
| No Missing Codes                       |                        | 15      |      |      | 16      |      |      | Bits |
| Integral Linearity Error               |                        | -6      | ±3   | +6   | -3      | ±1   | +3   | LSB  |
| Transition Noise                       |                        |         | 0.85 |      |         | 0.85 |      | LSB  |
| Gain Error <sup>1</sup> , Twin to Tmax |                        |         | ±2   | ±30  |         | ±2   | ±15  | LSB  |
| Offset Error, Twin to Tmax             |                        |         | ±0.7 | ±3.5 |         | ±0.7 | ±3.5 | mV   |
| <b>AC ACCURACY</b>                     |                        |         |      |      |         |      |      |      |
| Total Harmonic Distortion              | f <sub>m</sub> = 1 kHz |         | -94  |      |         | -98  |      | dB   |

Fig. 12. (Table D) A text-based PDF page in our testbench with potentially problematic columns highlighted with boxes (red).

TABLE II  
TABLE IDENTIFICATION YOLO COMPARISON.

|                   | Recall | Precision | F1-Score |
|-------------------|--------|-----------|----------|
| Custom CNN Only   | 0.5093 | 0.3802    | 0.4354   |
| YOLO Only         | 0.8654 | 0.8989    | 0.8818   |
| Custom CNN + YOLO | 0.8716 | 0.8663    | 0.8689   |

the specific way mandated by ICDAR. For every cell within a table, “adjacency relations” defined by ICDAR, are found with its nearest horizontal and vertical neighbors. In order to get precision and recall, the adjacency relations are compared with the ground truth. Precision and recall are computed for each document then the average is taken across all the documents.

To compare Tablertext to the other papers, the same method to calculate precision and recall is used. It should be noted that both DeepDeSRT and TableNet fine tune their networks by training on the remaining ICDAR dataset. In order to prove its ability to extract data from tables in general formats, Tablertext does not use any ICDAR data to train with. The results comparing the different methods can be found in Table III. Tablertext has both the highest precision and overall F1-Score out of the methods.

2) *Diverse Dataset comparison:* Without the source code for DeepDeSRT or TableNet, Tablertext cannot compare to them qualitatively or quantitatively on a large diverse dataset. Therefore, a comparison is made with the popular PDF extractor Tabula. Tabula can only extract non-scanned, text-based PDFs. So a testbench, that exclusively contains text-based PDFs, is created by randomly selecting over 100 tables within PDFs from various domains. Tabula is given the PDF pages with the text and meta-information, while Tablertext is given a scanned image of each PDF page. Fig. 12 is an example of one such page within our testbench.

Fig. 12 is extracted by both Tablertext and Tabula. With this input, Fig. 13 and Fig. 14 respectively show Tablertext’s

TABLE III  
ICDAR DATA EXTRACTION RESULTS COMPARISON. BOTH DEEPDESRT AND TABLENET ARE TRAINED ON A SUBSET OF ICDAR DATASET WHILE TABLETEXT DOES NOT USE ANY ICDAR DATA TO TRAIN WITH WHICH PROVES TABLETEXT’S ABILITY TO EXTRACT DATA FROM TABLES IN GENERAL FORMATS.

| Work       | Recall | Precision | F1-Score |
|------------|--------|-----------|----------|
| Tablertext | 0.9091 | 0.9221    | 0.9156   |
| DeepDeSRT  | 0.8736 | 0.9593    | 0.9144   |
| TableNet   | 0.9001 | 0.9307    | 0.9151   |

TABLE IV  
DIVERSE DATA EXTRACTION RESULTS COMPARISON.

| Work       | Recall | Precision | F1-Score |
|------------|--------|-----------|----------|
| Tablertext | 0.9192 | 0.9437    | 0.9313   |
| Tabula     | 0.7110 | 0.734     | 0.7223   |

and Tabula’s final CSV output. Both have a high accuracy when reproducing the cells’ contents, though Tabula is simply reading text provided to it, while Tablertext is using OCR. One issue with the Tablertext output is that Tesseract appears to not recognize the plus-minus character. Tabula also did not perfectly print the data within the cells, despite reading straight from the PDF. The minus character could not be recognized by Tabula. This problem appears because this PDF is using an abnormal character instead of the standard ASCII symbol.

Several severe structural errors appear in Tabula’s output and these errors have been highlighted in Fig. 14. Looking back at the original image in Fig. 12, it becomes clear what caused these errors. The cells “A Grade” and “B Grade” overlap both the columns beginning with “Typ” and “Max”. Errors like this are common for conventional text-based extractors [30]. Tablertext meanwhile, with its dynamic threshold ratio, easily separates these two columns then merges the split “A Grade” and “B Grade” cells back into a single cell.

The highlighted cells in Fig. 14 emphasize the importance of accurate cell location identification. Even with clever post-processing, it would be impossible to tell which column the merged cells with a single piece of data belong to. For

| Table 3.                               |                        |                  |      |     |                  |      |     |      |
|--|------------------------|------------------|------|-----|------------------|------|-----|------|
| Parameter                              | Conditions             | A Grade < EXTEND |      |     | B Grade < EXTEND |      |     | Unit |
|  |                        | Min              | Typ  | Max | Min              | Typ  | Max |      |
| <b>ACCURACY</b>                        |                        |                  |      |     |                  |      |     |      |
| No Missing Codes                       |                        | 15               |      |     | 16               |      |     | Bits |
| Integral Linearity Error               |                        | -6               | 3    | 6   | -3               | 1    | 3   | LSB  |
| Transition Noise                       |                        |                  | 0.5  |     |                  | 0.5  |     | LSB  |
| Gain Error <sup>1</sup> , Twin to Tmax |                        |                  | 2    | 24  |                  | 2    | 15  | LSB  |
| Offset Error, Twin to Tmax             |                        |                  | 0.7  | 1.6 |                  | 0.4  | 1.6 | mV   |
| <b>AC ACCURACY</b>                     |                        |                  |      |     |                  |      |     |      |
| Total Harmonic Distortion              | f <sub>m</sub> = 1 kHz |                  | -100 |     |                  | -106 |     | dB   |

  

| Table 4.                               |                        |                  |      |     |                  |      |     |      |
|--|------------------------|------------------|------|-----|------------------|------|-----|------|
| Parameter                              | Conditions             | A Grade < EXTEND |      |     | B Grade < EXTEND |      |     | Unit |
|  |                        | Min              | Typ  | Max | Min              | Typ  | Max |      |
| <b>ACCURACY</b>                        |                        |                  |      |     |                  |      |     |      |
| No Missing Codes                       |                        | 15               |      |     | 16               |      |     | Bits |
| Integral Linearity Error               |                        | -6               | 3    | 6   | -3               | 1    | 3   | LSB  |
| Transition Noise                       |                        |                  | 0.85 |     |                  | 0.85 |     | LSB  |
| Gain Error <sup>1</sup> , Twin to Tmax |                        |                  | 2    | 30  |                  | 2    | 15  | LSB  |
| Offset Error, Twin to Tmax             |                        |                  | 0.7  | 3.5 |                  | 0.7  | 3.5 | mV   |
| <b>AC ACCURACY</b>                     |                        |                  |      |     |                  |      |     |      |
| Total Harmonic Distortion              | f <sub>m</sub> = 1 kHz |                  | -94  |     |                  | -98  |     | dB   |

Fig. 13. (Table D) Tablertext’s final output in CSV format. Boxes (red) are showing the correct handling of problematic cells.

| Parameter                   | Conditions              | A Grade |           | B Grade |           | Unit |
|-----------------------------|-------------------------|---------|-----------|---------|-----------|------|
|                             |                         | Min     | Typ Max   | Min     | Typ Max   |      |
| ACCURACY                    |                         |         |           |         |           |      |
| No Missing Codes            |                         | 15      |           | 16      |           | Bits |
| Integral Linearity Error    |                         | ?6      | ±3 +6     | ?3      | ±1 +3     | LSB  |
| Transition Noise            |                         |         | 0.5       |         | 0.5       | LSB  |
| Gain Error1, TMIN to TMAX   |                         |         | ±2 ±24    |         | ±2 ±15    | LSB  |
| Offset Error1, TMIN to TMAX |                         |         | ±0.7 ±1.6 |         | ±0.4 ±1.6 | mV   |
| AC ACCURACY                 |                         |         |           |         |           |      |
| Total Harmonic Distortion   | f <sub>IN</sub> = 1 kHz |         | ?100      |         | ?106      | dB   |

  

| Parameter                   | Conditions              | A Grade |           | B Grade |           | Unit |
|-----------------------------|-------------------------|---------|-----------|---------|-----------|------|
|                             |                         | Min     | Typ Max   | Min     | Typ Max   |      |
| ACCURACY                    |                         |         |           |         |           |      |
| No Missing Codes            |                         | 15      |           | 16      |           | Bits |
| Integral Linearity Error    |                         | ?6      | ±3 +6     | ?3      | ±1 +3     | LSB  |
| Transition Noise            |                         |         | 0.85      |         | 0.85      | LSB  |
| Gain Error1, TMIN to TMAX   |                         |         | ±2 ±30    |         | ±2 ±15    | LSB  |
| Offset Error1, TMIN to TMAX |                         |         | ±0.7 ±3.5 |         | ±0.7 ±3.5 | mV   |
| AC ACCURACY                 |                         |         |           |         |           |      |
| Total Harmonic Distortion   | f <sub>IN</sub> = 1 kHz |         | ?94       |         | ?98       | dB   |

Fig. 14. (Table D) Tabula’s final output in CSV format. Boxes (red) are showing the incorrect handling of problematic cells

instance, looking at the row beginning with “Transition noise”, there is no way to know if the data belongs to the “Typ” or “Max” column.

The results comparing the two methods can be found in Table IV. Tablext’s ability to handle complex tables greatly surpasses Tabula’s, despite all of the extra meta-information that Tabula has access to.

## VI. CONCLUSION

This paper introduces a novel, general approach for table extraction. By utilizing both deep learning and computer vision techniques, the neural network can focus its attention on complex problems, while allowing more conventional methods to handle the simpler tasks. This focus allows Tablext to beat competing state-of-the-art neural-network-based table extraction methods as well as a popular open-source tool that requires table meta-data.

## ACKNOWLEDGEMENT

This material is based on research sponsored by Air Force Research Laboratory (AFRL) and Defense Advanced Research Projects Agency (DARPA) under agreement number FA8650-18-2-7844. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon.

## REFERENCES

- [1] A. Tengli, Y. Yang, and N. L. Ma, “Learning table extraction from examples,” in *Proceedings of the 20th international conference on Computational Linguistics*. Association for Computational Linguistics, 2004, p. 987.
- [2] E. Oro and M. Ruffolo, “Trex: An approach for recognizing and extracting tables from pdf documents,” in *2009 10th International Conference on Document Analysis and Recognition*. IEEE, 2009, pp. 906–910.
- [3] D. Pinto, A. McCallum, X. Wei, and W. B. Croft, “Table extraction using conditional random fields,” in *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, 2003, pp. 235–242.
- [4] E. Belval, “pdf2image: Open source scientific tools for python,” <https://github.com/Belval/pdf2image>, last accessed 2019-07-28.
- [5] “Grabzit: Open source scientific tools for python,” <https://github.com/GrabzIt/grabzit>, last accessed 2019-07-29.

- [6] Y. LeCun, P. Haffner, L. Bottou, and Y. Bengio, “Object recognition with gradient-based learning,” in *Shape, contour and grouping in computer vision*. Springer, 1999, pp. 319–345.
- [7] M. Fayazi, Z. Colter, E. Afshari, and R. Dreslinski, “Applications of artificial intelligence on the modeling and optimization for analog and mixed-signal circuits: A review,” *IEEE Transactions on Circuits and Systems I: Regular Papers*, 2021.
- [8] J. Redmon and A. Farhadi, “Yolov3: An incremental improvement,” *arXiv*, 2018.
- [9] R. S. *et al.*, “Tesseract ocr: Open source scientific tools for python,” <https://github.com/UB-Mannheim/tesseract>, last accessed 2019-07-29.
- [10] F. Cesarini, S. Marinai, L. Sarti, and G. Soda, “Trainable table location in document images,” in *Object recognition supported by user interaction for service robots*, vol. 3. IEEE, 2002, pp. 236–240.
- [11] M. Fan and D. S. Kim, “Table region detection on large-scale pdf files without labeled data,” *CoRR*, *abs/1506.08891*, 2015.
- [12] B. Gatos, D. Danatsas, I. Pratikakis, and S. J. Perantonis, “Automatic table detection in document images,” in *International Conference on Pattern Recognition and Image Analysis*. Springer, 2005, pp. 609–618.
- [13] A. Gilani, S. R. Qasim, I. Malik, and F. Shafait, “Table detection using deep learning,” in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 1. IEEE, 2017, pp. 771–776.
- [14] L. Hao, L. Gao, X. Yi, and Z. Tang, “A table detection method for pdf documents based on convolutional neural networks,” in *2016 12th IAPR Workshop on Document Analysis Systems (DAS)*. IEEE, 2016, pp. 287–292.
- [15] T. Kasar, P. Barlas, S. Adam, C. Chatelain, and T. Paquet, “Learning to detect tables in scanned document images using line information,” in *2013 12th International Conference on Document Analysis and Recognition*. IEEE, 2013, pp. 1185–1189.
- [16] F. Shafait and R. Smith, “Table detection in heterogeneous documents,” in *Proceedings of the 9th IAPR International Workshop on Document Analysis Systems*, 2010, pp. 65–72.
- [17] M. O. Perez-Arriaga, T. Estrada, and S. Abad-Mota, “Tao: system for table detection and extraction from pdf documents,” in *The Twenty-Ninth International Flairs Conference*, 2016.
- [18] T. Hassan and R. Baumgartner, “Table recognition and understanding from pdf files,” in *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, vol. 2. IEEE, 2007, pp. 1143–1147.
- [19] B. Yildiz, K. Kaiser, and S. Miksch, “pdf2table: A method to extract table information from pdf files,” in *IICAI*, 2005, pp. 1773–1785.
- [20] Y. Liu, K. Bai, P. Mitra, and C. Giles, “Searching for tables in digital documents,” in *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, vol. 2. IEEE, 2007, pp. 934–938.
- [21] A. K. Jain, M. N. Murty, and P. J. Flynn, “Data clustering: a review,” *ACM computing surveys (CSUR)*, vol. 31, no. 3, pp. 264–323, 1999.
- [22] K. Nishida, K. Sadamitsu, R. Higashinaka, and Y. Matsuo, “Understanding the semantic structures of tables with a hybrid deep neural network architecture,” in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [23] E. Koci, M. Thiele, W. Lehner, and O. Romero, “Table recognition in spreadsheets via a graph representation,” in *2018 13th IAPR International Workshop on Document Analysis Systems (DAS)*, 2018, pp. 139–144.
- [24] A. Puha, O. Rinciog, and V. Posea, “Enhancing open data knowledge by extracting tabular data from text images,” in *DATA*, 2018.
- [25] S. R. Qasim, H. Mahmood, and F. Shafait, “Rethinking table recognition using graph neural networks,” in *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 2019, pp. 142–147.
- [26] S. Schreiber, S. Agne, I. Wolf, A. Dengel, and S. Ahmed, “Deepdesrt: Deep learning for detection and structure recognition of tables in document images,” in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 1. IEEE, 2017, pp. 1162–1167.
- [27] S. S. Paliwal, D. Vishwanath, R. Rahul, M. Sharma, and L. Vig, “Tablenet: Deep learning model for end-to-end table detection and tabular data extraction from scanned document images,” in *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 2019, pp. 128–133.
- [28] M. Göbel, T. Hassan, E. Oro, and G. Orsi, “Icdar 2013 table competition,” in *2013 12th International Conference on Document Analysis and Recognition*, 2013, pp. 1449–1453.

- [29] M. Aristarán, T. Mike, J. B. Merrill, J. Das, D. Frackman, and T. Swicegood, "Tabula: Open source scientific tools for python," <https://github.com/tabulapdf/tabula>, last accessed 2019-07-29.
- [30] R. Dreslinski *et al.*, "Fully-autonomous soc synthesis using customizable cell-based synthesizable analog circuits," University of Michigan Ann Arbor United States, Tech. Rep., 2019.