

# Handling big tabular data of ICT supply chains: a multi-task, machine-interpretable approach

Bin Xiao, Murat Simsek, Burak Kantarci  
School of Electrical Engineering and Computer Science  
University of Ottawa  
Ottawa, Canada  
{bxiao103, murat.simsek, burak.kantarci}@uottawa.ca

Ala Abu Alkheir  
Directorate, Analytics  
Lytica Inc.  
Ottawa, Canada  
ala\_abualkheir@lytica.com

**Abstract**—Due to the characteristics of Information and Communications Technology (ICT) products, the critical information of ICT devices is often summarized in big tabular data shared across supply chains. Therefore, it is critical to automatically interpret tabular structures with the surging amount of electronic assets. To transform the tabular data in electronic documents into a machine-interpretable format and provide layout and semantic information for information extraction and interpretation, we define a Table Structure Recognition (TSR) task and a Table Cell Type Classification (CTC) task. We use a graph to represent complex table structures for the TSR task. Meanwhile, table cells are categorized into three groups based on their functional roles for the CTC task, namely Header, Attribute, and Data. Subsequently, we propose a multi-task model to solve the defined two tasks simultaneously by using the text modal and image modal features. Our experimental results show that our proposed method can outperform state-of-the-art methods on ICDAR2013 and UNLV datasets.

**Index Terms**—Big Data Analytics, Supply Chain Optimization, Image Processing, Table Structure Recognition, Table Cell Type Classification

## I. INTRODUCTION

In Information and Communications Technology (ICT) supply chains, electronic devices often contain various parameters, units, or other critical information formatted in tables, making it vital to extract and interpret tables from electronic documents. Even though tables in electronic documents are user-friendly for human readers, they are often not structured and not machine-interpretable. It is also not practical for humans to read, extract and interpret tables from millions of electronic documents. Therefore, to deal with the vast amount of electronic documents in the ICT supply chain and make unstructured tabular data machine-interpretable, we define a Table Structure Recognition (TSR) problem, which can recover a complex table structure with a graph, and a Table Cell Type Classification (CTC) problem based on the cell's functional roles. More specifically, for the TSR problem, each cell in a table is represented with a vertex in a graph, and three types of cell associations, namely vertical connection, horizontal connection, and no connection, are defined to represent the locational relations among table cells,

which can be represented by the edges in a graph. Meanwhile, similar to some existing studies [1], [2], [3] discussing tabular cell classification problems with different taxonomies of cell types, we define three types of cells, namely Header, Attribute, and Data. Figure 1 shows a sample table with defined types in the CTC problem and its graph representation of the table numbered part. Typically, headers in a table can express the meaning of their corresponding columns, attributes represent the meaning of their corresponding rows, and data cells are used to present the exact information of the specific header and attribute. In other words, the facts and information contained in a table can become accessible and easily located and extracted based on the defined three types of cells and the machine-readable table structure. It is worth mentioning that most of the existing studies [1], [2], [3] on CTC problems focus on tables in spreadsheets which is a much easier problem definition because spreadsheets can provide more meta-information and the default units in the spreadsheets are cells. There are some studies [4], [5] trying to extract entities and information from images, which is similar with our defined CTC problem, but they are not focusing on tables in document images, meaning that their problem definitions are not suitable to be solved together with TSR problem in a multi-task manner.

The main contribution of this study is three-fold: 1) This study extends the CTC problem to the tables in Portable Document Format (PDF) documents and image documents and builds a benchmark for this problem. 2) We propose a multi-task approach that simultaneously solves the defined TSR and CTC problems, achieving state-of-the-art results. Experimental results show that our proposed method can increase the F1 score from 70.76% to 88.17%, from 71.90 to 83.74 for the CTC task on the ICDAR2013 and UNLV datasets, respectively. For the TSR task, the proposed method can increase the F1 score from 92.30% to 93.04% on the ICDAR2013 dataset and from 87.24% to 89.51% on the UNLV dataset. 3) We discuss the different aspects of the proposed method and show their effectiveness with experiments.

This paper is organized as follows: Section II discusses the related studies. Section III presents the full process of transforming the tabular data into the machine-interpretable format and presents the proposed multi-task method. Section IV presents and discusses the experimental results. We

This work was supported in part by Mathematics of Information Technology and Complex Systems (MITACS) Accelerate Program, Smart Computing for Innovation Program (SOSCIP) and Lytica Inc.

MAXIMUM RATINGS AND ELECTRICAL CHARACTERISTICS (T <sub>J</sub> =25°C unless otherwise noted)				
PARAMETER	SYMBOL	UNIT	TSF34U60C	
Maximum repetitive peak reverse voltage	V <sub>RRM</sub>	V	80	90
Maximum average forward rectified current	I <sub>FAV</sub>	A	14.30	16
Peak forward surge current, 8.3 ms single half sine-wave superimposed on rated load per diode	I <sub>FSM</sub>	A	250	
Voltage rate of change (Rated V <sub>dv</sub> )	dV/dt	V/μs		
			TYP.	MAX.
Instantaneous forward voltage per diode (Note 1)	V <sub>F</sub>	V	0.48	0.57
	I <sub>F</sub> = 15A	T <sub>J</sub> = 25°C		
	I <sub>F</sub> = 15A	T <sub>J</sub> = 125°C	0.43	0.52
Instantaneous reverse current per diode at rated reverse voltage	I <sub>R</sub>	μA	-	500
				60
Typical thermal resistance per diode	R <sub>θJC</sub>	°C/W	4	
Operating junction temperature range	T <sub>J</sub>	°C	-55 to +150	
Storage temperature range	T <sub>STG</sub>	°C	-55 to +150	

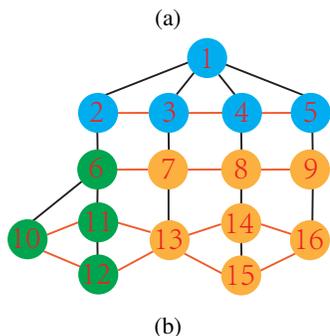


Fig. 1: Figure (a) shows three types of table cells, namely header, attribute, and data of an ICT device. Figure (b) is part of the table cells' graph representation, in which each vertex represents a cell, red lines stand for horizontal connection, and black lines stand for vertical connection.

draw our conclusion and possible directions in section V.

## II. RELATED WORK

To extract and interpret critical information from the surging amount of electronic documents in the global ICT supply chain, it is vital but often challenging to transform unstructured tabular data into structured and machine-interpretable format, because it needs several separate steps to finish this task. Since the dominant studies in recent years first transform PDF files into document images and then use deep learning based methods, we only focus on deep learning based approaches.

### A. Tabular Structure Recognition

TSR aims to identify table cells in identical rows and columns, which is challenging because of the complex table structures, such as merging cells and non-explicit border lines. Using a graph to represent the complex table structure is a popular design, in which table cells are represented by the graph nodes, and the associations of cells are defined in three types: vertical connection, horizontal connection and no connection, and graph edges are used to represent the defined cell associations. Many studies [6], [7] follow this problem formulation and propose bottom-up approaches. In contrast, some work, such as DeepDeSRT [8], CascadeTabNet [9] and TableDet [10], using top-down approaches would define the structural recognition as object detection or segmentation problem, often together with table detection problem. Similarly, these methods often employ and extend Cascade R-

CNN [11], Mask-RNN [12] and utilize transfer learning and data augmentation methods to improve the performance.

### B. Table Cell Type Classification

Many studies employ text embeddings and stylistic features for the spreadsheet's CTC problem. Elvis et al., in their work [1] use Weka [13] to select features and try various tree-based classifiers. Majid et al., in their work [2], propose to incorporate pre-trained cell embedding and stylistic features. They design a neighbor-based approach to generate cell contexts, propose an embedding model built on the InferSent model [14], and build an LSTM based classifier. Language models also have been used to the tabular cell classification problem. Training a large language model requires a large corpus and many computation resources. In order to apply a language model to the CTC task, the pre-trained language model is firstly used to generate the feature embeddings and fine-tune a cell type classification model. TUTA [3] is a large language model for a generally structured table trained with a large spreadsheets corpus. Besides, other language models, such as BERT [15] also can be used for the CTC problem. All these studies focus on Spreadsheets, which have inherent structural information that can help to train and fine-tune context-based language models. Conversely, tabular data in PDF documents do not contain this type of internal structural information, making it a more challenging task.

## III. PROPOSED METHOD

As illustrated in Fig. 2, manufacturers generate "big tabular data" with key information in tables but are hard to be shared with other participators in the supply chain. Our proposed method transforms unstructured tabular information into structured and machine-interpretable format, paving the way to efficient information sharing in the whole supply chain.

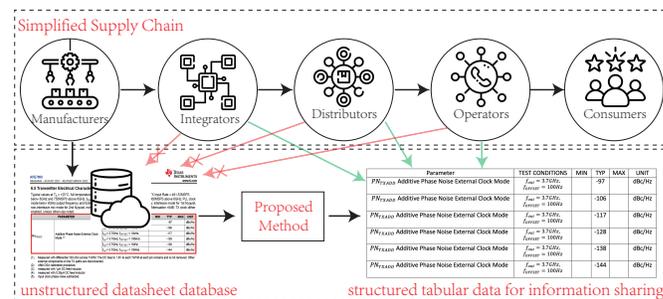


Fig. 2: A simplified ICT supply chain sample.

### A. Problem Definition

TSR aims at recovering complex table structures into a structured and machine-interpretable format. Following the problem definition in some existing bottom-up approaches [7], we also use a graph to represent a complex table structure, in which graph vertexes are used to represent table cells and graph edges are use to denote the associations between two

cells in a table, as shown in Figure 1. Assuming that each table cell's bounding box is given, which can be denoted by  $\{x_i^1, x_i^2, y_i^1, y_i^2\}$ , where  $i$  is the cell number. Then, the  $k$ th table  $t_k$  in a table set  $\mathbf{T} = \{t_k; k \in K\}$  containing  $K$  tables can be represented by a set of cells  $\mathbf{C}_k = \{c_i^k; i \in N_k\}$ , where  $c_i^k$  stands for the  $i$ th cell in the table  $t_k$ ,  $N_k$  is the number of cells in the table  $t_k$ . Thus, for the table structure recognition problem, given a training set consists of  $K$  tables, and the table  $t_k$ 's cell association set  $\mathbf{R}_k = \{r_{\{i,j\}}^k; i \neq j, i \in N_k, j \in N_k, k \in K\}$ , where  $r_{\{i,j\}}^k = \{\{c_i^k, c_j^k\}; i \neq j, i \in N_k, j \in N_k, k \in K\}$ , and the corresponding label set  $\mathbf{Y}_k^{tsr} = \{y_{\{i,j\}}^k; y \in \{0, 1, 2\}, i \in N_k, j \in N_k, k \in K\}$  of the association set, where 0, 1, 2 means the two cells in the association are horizontal connected, vertical connected and not connected. Therefore, the task of table structure recognition is to train a predictive model with given training data that can determine the probability of the two cells relation in a cell association, namely  $P_\theta(y_{\{i,j\}}^k = \hat{y}_{\{i,j\}}^k | r_{\{i,j\}}^k), \hat{y}_{\{i,j\}}^k \in \{0, 1, 2\}$ .

TABLE I: A summary of notations used in this paper.

$\mathbf{T}$	A table set
$\mathbf{C}_k$	The cell set of the $k$ th table
$\mathbf{R}_k$	The cell association set of the $k$ th table
$\mathbf{Y}_k^{tsr}$	The label set of the $k$ th table for the TSR
$\mathbf{Y}_k^{ctc}$	The label set of the $k$ th table for the CTC
$\{x_i^1, x_i^2, y_i^1, y_i^2\}$	The coordinate of the $i$ th cell in a table
$c_i$	The $i$ th cell in a table
$t_k$	The $k$ th table in a table set $\mathbf{T}$ containing $K$ tables
$r_{\{i,j\}}^k$	The association of the $i$ th, $j$ th cell in the $k$ th table
$\mathcal{L}$	The loss function
$\theta, \phi, \omega, \sigma$	The trainable parameters of each function
$\mathcal{F}$	A fully connected layer
$\mathbf{f}$	The output of a fully connected layer
$CAT$	The function that can represent the fusion layer
$\mathcal{C}$	The classifier
$\mathcal{E}$	The embedding network
$\mathbf{x}$	The input of a function
$\mathbf{y}$	The output of the proposed method
$\mathbb{R}^{ch*h*w}$	A real value set with $ch$ channel, $h$ height, $w$ width

CTC problem aims at identifying their functional roles in a complex table structure. For the table cell classification program, given a training set consists of  $K$  tables, each table  $t_k$  has a set of cells  $\mathbf{C}_k = \{c_i^k; i \in N_k\}$  and their corresponding types  $\mathbf{Y}_k^{ctc} = \{y_i^k; i \in N_k\}$ . We define three types of cells, namely Header, Attribute, Data, meaning that  $y_i^k \in \{0, 1, 2\}$ . Therefore, the task of table cell classification is to train a predictive model with given training data that can determine the probability of a cell belonging a predefined cell types, namely  $P_\phi(y_i^k = \hat{y}_i^k | c_i^k), \hat{y}_i^k \in \{0, 1, 2\}$ . Notably, the proposed method is a bottom-up approach and requires that each cell's bounding boxes are known. In our implementation, we use MMOCR [16] to detect table cells and extract the content from table images. The notations used in this paper are summarized in Table I.

### B. Multi-task Table Structure Recognition and Table Cell Type Classification

To transform unstructured tabular data into a structured and machine-interpretable format, we need three steps: table

detection, table structure recognition, and table cell type classification. Luckily, there have been some studies that achieved promising results on the table detection, as discussed in section II. In practice, we first convert all the PDF documents into document images and then use TableDet [10] to extract all the tables from the document images such that we can obtain the table set  $T = \{t_k; k \in K\}$ . To generate the cell association set required by the table structure recognition task, we follow the popular KD-tree based K-nearest method, which is also used in existing studies [7], [17]. More specifically, we can calculate the distance between two cells by a pre-defined distance metric, usually Euclidean Distance, using their bounding boxes. Thus, for the  $i$ th cell in the  $k$ th table  $t_k$ , we can find its nearest top  $M$  neighbor cells to form association pairs, and a KD-Tree model can easily implement this process. Notably,  $M$  is a hyperparameter set to 20 in our implementation.

Even though PDF documents and document images can only provide very limited meta-information, they can inherently provide visual and text information, meaning that features can be extracted from both text and image modalities. Following the typical design of multi-task approaches, there are two branches in the proposed model, one is for the TSR task, and another is for the CTC task, as shown in Figure 3. The TSR branch only utilizes the features from the image modal, which can be extracted by an embedding network. In contrast, the CTC branch uses features from both image modal, text modal, and cells' coordinates. More specifically, in the TSR branch, we follow the design of study [17], crop both cells' images in each cell association pair and their context image, and then use the cropped cell images and the context image as the input in the table structure recognition branch, as shown in the input sample part of Figure 3. Then the generated images are fed into an embedding network, denoted by the function  $\mathcal{E}_\phi$ , to extract visual feature maps. Then the extracted feature maps are fed into a fully connected layer which can be represented by a function  $\mathcal{F}_{\theta_1}^{tsr}$ . At last, the outputs of  $\mathcal{F}_{\theta_1}^{tsr}$  are further fed into the classifier  $\mathcal{C}_{\sigma_1}^{tsr}$  with a softmax function to get the probability of each class. The full process can be represented by the Equation 1, where  $ch, h, w$  are the number of channels, image height, and image width, respectively, and all of  $\theta, \sigma, \phi$  are trainable parameters.

$$\begin{aligned} \mathbf{y}_{tsr} &= \mathcal{C}_{\sigma_1}^{tsr}(\mathbf{f}_{tsr}), \mathbf{f}_{tsr} = \mathcal{F}_{\theta_1}^{tsr}(\mathcal{E}_\phi(\mathbf{x}^{tsr})), \\ \mathbf{x}_{tsr} &\in \mathbb{R}^{ch*h*w}, \mathbf{f}_{tsr} \in \mathbb{R}^{l_1}, \mathbf{y}_{tsr} \in \mathbb{R}^3 \end{aligned} \quad (1)$$

Meanwhile, in the CTC branch, together with the cells' coordinates, the cropped cell images and their corresponding texts are used as the inputs of the CTC branch. Assuming that the InferSent [14] module can be represented by function  $\mathcal{I}_\omega$ , the fully connected layer for the text modal can be denoted by function  $\mathcal{F}_{\theta_2}^{ctc}$ , then the text features can be calculated by the Equation 2, where  $\mathbf{f}_2$  is the output feature and  $l_{text}$  is the number of tokens in the input sentence.

$$\mathbf{f}_{text} = \mathcal{F}_{\theta_2}^{ctc}(\mathcal{I}_\omega(\mathbf{x}_{text})), \mathbf{x}_{text} \in \mathbb{R}^{l_{text}}, \mathbf{f}_{text} \in \mathbb{R}^{l_2} \quad (2)$$

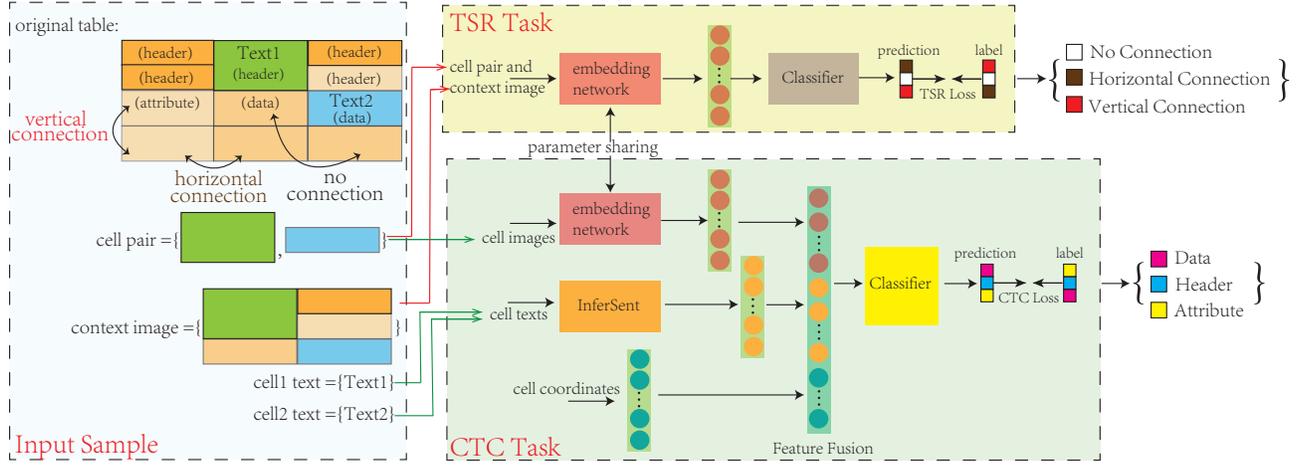


Fig. 3: Overall architecture of the proposed multi-task model. Notably, header, attribute and data are three types of cell in the CTC task, and horizontal connection, vertical connection and no connection are three possible outputs in the TSR task.

Similarly, we can obtain the image features and the coordinate features in the CTC branch using the Equation 3 and 4. At last, all features are fused together with fusion function  $\mathcal{CAT}$  and fed into the classifier  $\mathcal{C}_{\sigma_2}^{ctc}$ , which can be achieved by Equation 5. It is worth mentioning that  $l_1, l_2, l_3$  and  $l_4$  are determined by the number of neurons in each fully connected layer, and  $\mathcal{E}_\phi$  shares parameters in both branches.

$$\mathbf{f}_{img} = \mathcal{F}_{\theta_2}^{ctc}(\mathcal{E}_\phi(\mathbf{x}_{img})), \mathbf{x}_{img} \in \mathbb{R}^{c*h*w}, \mathbf{f}_{img} \in \mathbb{R}^{l_3} \quad (3)$$

$$\mathbf{f}_{coord} = \mathcal{F}_{\theta_3}^{ctc}(\mathbf{x}_{coord}), \mathbf{x}_{coord} \in \mathbb{R}^4, \mathbf{f}_{coord} \in \mathbb{R}^{l_4} \quad (4)$$

$$\mathbf{y}_{ctc} = \mathcal{C}_{\sigma_2}^{ctc}(\mathcal{CAT}(\mathbf{f}_{img}, \mathbf{f}_{text}, \mathbf{f}_{coord})), \mathbf{y}_{ctc} \in \mathbb{R}^3 \quad (5)$$

Since the proposed method employs multi-task architecture, meaning that all the parameters should be trained together, the loss function  $\mathcal{L}_{mt}$  can be defined as Equation 6, which contains a TSR loss denoted by  $\mathcal{L}_{tsr}$ , a CTC loss denoted by  $\mathcal{L}_{ctc}$  and a hyper parameter  $\lambda$  to balance  $\mathcal{L}_{tsr}$  and  $\mathcal{L}_{ctc}$ . We use cross-entropy loss for both  $\mathcal{L}_{tsr}$  and  $\mathcal{L}_{ctc}$  in our implementation, which can be defined as Equation 7, where  $y_i$  means the  $i$ th prediction and  $\hat{y}_i$  is its corresponding ground truth.

$$\mathcal{L}_{mt} = (1 - \lambda) * \mathcal{L}_{tsr} + \lambda * \mathcal{L}_{ctc} \quad (6)$$

$$\mathcal{L}_{ce}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^N \hat{y}_i \log y_i \quad (7)$$

#### IV. EXPERIMENTS AND ANALYSIS

ICDAR2013 and UNLV datasets are popular for table detection and TSR problems, containing 156 and 540 tables, respectively. We further annotate these two datasets with "Header," "Attribute," or "Data," as discussed in section III-A,

and exclude all the empty cells in the CTC task labeling because empty cells can be easily identified without using any machine learning models. In the experiments, the ICDAR2013 dataset is randomly split into a validation set with 33 tables and a testing set with 123 tables. The UNLV dataset is also randomly split into a training, validation, and testing set with 323, 107, 110 tables, respectively. Notably, we remove 18 tables from the UNLV dataset because of the ambiguity of their labels, and we use the training set of the UNLV dataset when evaluating the performance on the ICDAR2013 dataset, following the experimental setting of study [17]. We use MMOCR [16] to extract text contents from the table cells for both ICDAR2013 and UNLV datasets.

#### A. Implementation Details and Experimental Results

As discussed in section III-B, the proposed method contains an embedding network  $\mathcal{E}_\phi$ , which is used in both the TSR task and CTC task. In our implementation, we use a simple ConvNet-4 [18], which consists of four convolution layers. Following the method described in the study [17], the input images of the embedding network are firstly resized to the dimension  $3 * 84 * 84$  with zero padding to keep the original height-width ratio. Then the resized images are fed into the embedding network  $\mathcal{E}_\phi$  followed by a fully connected layer  $\mathcal{F}$  to obtain the features. In our implementation, we set the number of neurons in both  $\mathcal{F}_{\theta_1}^{tsr}$  and  $\mathcal{F}_{\theta_2}^{ctc}$  to 128, meaning that both  $l_1$  and  $l_3$  equal 128. Similarly,  $l_2$  and  $l_4$  are set to 128 and 32, respectively, and  $\lambda$  in the loss function is set to 0.3 for the ICDAR2013 dataset and 0.6 for the UNLV dataset. Each classifier, as shown in Figure 3, is implemented by a fully connected layer with three neurons followed by a softmax function to transform the logits into a distribution. We use a simple concatenation function to implement the fusion function  $\mathcal{CAT}$  which does not have any trainable parameters. Notably, since the datasets used in the experiments are unbal-

anced, we use a cost-sensitive method in the implementation of the loss functions.

For the TSR task, we list TabbyPDF [19], GraphTSR [7], DeepDeSRT [8] and CATT-Net [17] as benchmark models. Considering that our proposed multi-task model relies on features from both image modal and text modal, we follow the popular pre-trained fine-tune paradigm, implement four models by fine-tuning the Glove [20], FastText [21], BERT [15] to compare the state-of-the-art methods in the NLP. In terms of image modal, we implement two image classification models, namely ConvNeXt [22] and ResNet50 [23] without using pre-trained weights. BERT is a popular context-based language model, meaning that the performance of these two models heavily relies on long dependency in a context, whereas Glove and FastText are two typical non-context-based models. ResNet50 is a classic convolution network that is widely used in image classification problems, while ConvNeXt [22] is the state-of-the-art model for the image classification problem. Since we only focus on the tables and use TableDet [10] to extract all the tables from the document images without considering the context information of the tables, fine-tuning BERT shows very poor performance for the CTC problem; hence not included in Table III.

Table II and Table III show the overall Precision, Recall, and F1-score of the proposed method for the TSR task and CTC task, respectively. The performance scores of TabbyPDF, GraphTSR, and DeepDeSRT on the ICDAR2013 dataset come from the study [7]. "-" in Table II means the score is not reported in related studies. The experimental results show that the proposed method can outperform the benchmark models in the TSR and CTC tasks.

TABLE II: Experimental results for the TSR task.

Method	ICDAR2013			UNLV		
	Prec	Recall	F1	Prec	Recall	F1
TabbyPDF	78.90	84.50	81.60	—	—	—
GraphTSR	81.90	85.50	83.70	—	—	—
DeepDeSRT	57.30	56.40	56.80	—	—	—
CATT-Net	<b>94.10</b>	90.70	92.30	86.28	<b>88.31</b>	87.24
Ours	92.85	<b>93.29</b>	<b>93.04</b>	<b>92.66</b>	86.78	<b>89.52</b>

TABLE III: Experimental results for the CTC task.

Method	ICDAR2013			UNLV		
	Prec	Recall	F1	Prec	Recall	F1
Glove	56.47	60.88	57.74	55.00	62.03	57.33
FastText	67.85	64.16	65.46	63.72	68.85	65.85
ConvNeXt	70.15	71.42	70.76	65.88	66.72	66.21
ResNet50	72.81	69.03	70.16	72.62	71.83	71.90
Ours	<b>87.94</b>	<b>88.41</b>	<b>88.17</b>	<b>82.19</b>	<b>85.51</b>	<b>83.74</b>

## B. Discussion and Analysis

1) *Ablation Study*: In this section, we conduct extra experiments to demonstrate each component’s effectiveness in our proposed method. Firstly, we implement two single-task models for the TSR task and CTC, respectively, using the corresponding branch design of the proposed method. Since features from the image modal, text modal, and coordinate

information are utilized in the CTC task, we also implement three models that only use features from each modal. The experimental results are shown in Table IV and Table V, in which  $S(I)$ ,  $S(T)$ ,  $S(C)$ ,  $S(I+T+C)$  means the single-task model with image feature, text feature, coordinate feature, and the combined three types of features, respectively, ours means the proposed multi-task method.

The experimental results in Table IV show that features from text modal and coordinate information are not efficient and can hardly bring any benefits when combined with features from the image modal. In contrast, visual features from the image modal are more efficient and can lead to the best performance compared with other benchmarks. Therefore, we only use visual features in the TSR task in our proposed multi-task method. Meanwhile, the experimental results in Table V show that visual features, text features, and coordinate information can be helpful for the CTC problem, and the combination of both types of features can contribute to a significant performance improvement. Therefore, we use visual features, text features, and coordinate information in the CTC branch of the proposed multi-task method.

TABLE IV: Ablation study results for the TSR task.

Method	ICDAR2013			UNLV		
	Prec	Recall	F1	Prec	Recall	F1
S(T)	36.36	34.56	34.32	40.03	35.91	35.94
S(I)	89.42	91.71	90.52	84.39	<b>87.52</b>	85.75
S(C)	50.81	77.69	45.30	60.38	83.75	61.17
S(T+I+C)	88.71	93.53	90.98	90.57	84.16	87.11
Ours	<b>92.85</b>	<b>93.29</b>	<b>93.04</b>	<b>92.66</b>	86.78	<b>89.52</b>

TABLE V: Ablation study results for the CTC task.

Method	ICDAR2013			UNLV		
	Prec	Recall	F1	Prec	Recall	F1
S(T)	68.11	62.46	64.48	72.047	67.85	69.41
S(I)	70.12	67.52	68.46	63.98	76.49	68.62
S(C)	86.07	<b>90.84</b>	87.65	76.98	<b>85.55</b>	80.74
S(T+I+C)	85.44	88.68	86.53	77.75	81.59	78.70
Ours	<b>87.94</b>	88.41	<b>88.17</b>	<b>82.19</b>	85.51	<b>83.74</b>

2) *The impact of hyper parameter lambda*: Our proposed method contains a TSR branch and a CTC branch, both of which lead to a loss, as shown in Equation 6. The hyper-parameter  $\lambda$  is used to balance the TSR task’s loss and the CTC task’s loss. We conducted extra experiments to discuss the influence of this hyperparameter on the proposed model, and the results are shown in Figure 4. The figure shows that the performance of the proposed method can be influenced by the value of  $\lambda$ , but can sustain an overall stable performance.

3) *The impact of feature fusion*: As shown in Figure 3, the features used in the CTC task come from different modalities, including text modal, image modal, and cells’ coordinates, and the feature fusion layer fuse these features together before they are fed into the classifier. Therefore, we conduct extra experiments to discuss the impact of these values. The number of neurons in a fully connected layer can be any integer larger than 0. Considering the scores reported in Table V, and simplifying the problem, we assume  $l_1 = l_2 = l_3 = l_4$  in

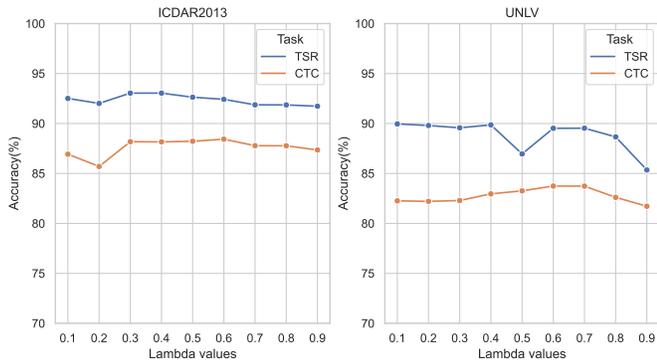


Fig. 4: Experimental results with different  $\lambda$  values.

this section. Table VI and Table VII present the experimental results when  $l_1$  equals 16, 32, 64, 128 respectively, and the results show that increasing the number of neurons can lead to better performance to some degree, especially for TSR.

TABLE VI: Experimental results for the the TSR task.

Method	ICDAR2013			UNLV		
	Prec	Recall	F1	Prec	Recall	F1
$l_{16}$	81.24	92.57	85.86	81.65	82.25	81.65
$l_{32}$	84.75	<b>93.12</b>	88.49	79.77	82.86	80.55
$l_{64}$	<b>85.69</b>	92.51	<b>88.75</b>	79.85	<b>85.74</b>	82.53
$l_{128}$	80.13	93.60	85.54	<b>87.41</b>	85.66	<b>86.45</b>

TABLE VII: Experimental results for the CTC task.

Method	ICDAR2013			UNLV		
	Prec	Recall	F1	Prec	Recall	F1
$l_{16}$	79.76	82.48	78.83	<b>82.84</b>	81.94	82.26
$l_{32}$	83.54	83.50	82.29	81.02	85.04	82.92
$l_{64}$	89.47	81.22	84.72	80.76	87.09	<b>83.57</b>
$l_{128}$	<b>90.83</b>	<b>84.33</b>	<b>87.22</b>	79.36	<b>87.18</b>	82.86

## V. CONCLUSION AND FUTURE WORK

Participants in the ICT supply chains generate "big tabular data" that is difficult to manage but often contains valuable information. This paper introduced a multi-task model that simultaneously solves the TSR and CTC problems, which are critical steps in transforming tabular data into a machine-interpretable format. Experimental results show that both TSR and CTC tasks can benefit from the multi-task design compared to their single-task counterparts. Besides, the proposed method can outperform the state-of-the-art benchmark models by a large margin, increasing the F1 score from 70.76% to 88.17%, and from 71.90 to 83.74 for the CTC task on the ICDAR2013 and UNLV dataset respectively. For the TSR task, increase the F1 score from 92.30% to 93.04% on the ICDAR2013 dataset and from 87.24% to 89.51% on the UNLV dataset. We tackle the TSR and CTC simultaneously with a single multi-task model, whereas the two problems can also be addressed in sequence in future work.

## REFERENCES

- [1] E. Koci, M. Thiele, O. Romero, and W. Lehner, "Cell classification for layout recognition in spreadsheets," in *International Joint Conference on Knowledge Discovery, Knowledge Engineering, and Knowledge Management*. Springer, 2016, pp. 78–100.
- [2] M. G. Gol, J. Pujara, and P. Szekely, "Tabular cell classification using pre-trained cell embeddings," in *2019 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2019, pp. 230–239.
- [3] Z. Wang, H. Dong, R. Jia, J. Li, Z. Fu, S. Han, and D. Zhang, "Tuta: Tree-based transformers for generally structured table pre-training," in *Proc. of the 27th ACM SIGKDD Conference*, 2021, pp. 1780–1790.
- [4] D. Lohani, A. Belaïd, and Y. Belaïd, "An invoice reading system using a graph convolutional network," in *Asian Conference on Computer Vision*. Springer, 2018, pp. 144–158.
- [5] G. Tang, L. Xie, L. Jin, J. Wang, J. Chen, Z. Xu, Q. Wang, Y. Wu, and H. Li, "Matchvie: Exploiting match relevancy between entities for visual information extraction," *arXiv preprint arXiv:2106.12940*, 2021.
- [6] D. A. et al, "Table structure recognition based on cell relationship, a bottom-up approach," in *RANLP*, 2019.
- [7] Z. Chi, H. Huang, H.-D. Xu, H. Yu, W. Yin, and X.-L. Mao, "Complicated table structure recognition," *ArXiv*, vol. abs/1908.04729, 2019.
- [8] S. S. et al, "Deepdesrt: Deep learning for detection and structure recognition of tables in document images," *14th IAPR Intl. Conf. on Document Analysis and Recognition*, vol. 01, pp. 1162–1167, 2017.
- [9] C. Tensmeyer, V. I. Morariu, B. Price, S. Cohen, and T. Martinez, "Deep splitting and merging for table structure decomposition," in *Intl Conf on Document Analysis and Recognition*. IEEE, 2019, pp. 114–121.
- [10] J. Fernandes, M. Simsek, B. Kantarci, and S. Khan, "Tabledet: An end-to-end deep learning approach for table detection and table image classification in data sheet images," *Neurocomputing*, 2021.
- [11] Z. Cai and N. Vasconcelos, "Cascade r-cnn: Delving into high quality object detection," *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6154–6162, 2018.
- [12] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *IEEE Intl. Conf. on computer vision*, 2017, pp. 2961–2969.
- [13] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, "The weka data mining software: an update," *ACM SIGKDD explorations newsletter*, vol. 11, no. 1, pp. 10–18, 2009.
- [14] A. Conneau, D. Kiela, H. Schwenk, L. Barrault, and A. Bordes, "Supervised learning of universal sentence representations from natural language inference data," in *Proc. of the Conf. on Empirical Methods in Natural Language Processing*, Sep. 2017, pp. 670–680.
- [15] J. D. et al., "Bert: Pre-training of deep bidirectional transformers for language understanding," in *NAACL*, 2019.
- [16] Z. Kuang, H. Sun, Z. Li, X. Yue, T. H. Lin, J. Chen, H. Wei, Y. Zhu, T. Gao, W. Zhang et al., "Mmocr: a comprehensive toolbox for text detection, recognition and understanding," in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 3791–3794.
- [17] B. Xiao, M. Simsek, B. Kantarci, and A. A. Alkheir, "Table structure recognition with conditional attention," *Preprint: https://arxiv.org/abs/2203.03819*, 2022.
- [18] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra, "Matching networks for one shot learning," *arXiv preprint arXiv:1606.04080*, 2016.
- [19] A. Shigarov, A. Altaev, A. Mikhailov, V. Paramonov, and E. Cherkashin, "Tabbypdf: web-based system for pdf table extraction," in *Intl Conf. on Information and Software Technologies*, 2018, pp. 257–269.
- [20] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proc. of the conf. on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [21] T. Mikolov, E. Grave, P. Bojanowski, C. Puhres, and A. Joulin, "Advances in pre-training distributed word representations," in *Proc. of the Intl. Conference on Language Resources and Evaluation*, 2018.
- [22] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," *arXiv preprint arXiv:2201.03545*, 2022.
- [23] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.