# Split, Embed and Merge: An accurate table structure recognizer

Zhenrong Zhang, Jianshu Zhang, Jun Du and Fengren Wang

**Abstract**

Table structure recognition is an essential part for making machines understand tables. Its main task is to recognize the internal structure of a table. However, due to the complexity and diversity in their structure and style, it is very difficult to parse the tabular data into the structured format which machines can understand, especially for complex tables. In this paper, we introduce Split, Embed and Merge (SEM), an accurate table structure recognizer. SEM is mainly composed of three parts, splitter, embedder and merger. In the first stage, we apply the splitter to predict the potential regions of the table row/column separators, and obtain the fine grid structure of the table. In the second stage, by taking a full consideration of the textual information in the table, we fuse the output features for each table grid from both vision and text modalities. Moreover, we achieve a higher precision in our experiments through providing additional textual features. Finally, we process the merging of these basic table grids in a self-regression manner. The corresponding merging results are learned through the attention mechanism. In our experiments, SEM achieves an average F1-Measure of 97.11% on the SciTSR dataset which outperforms other methods by a large margin. We also won the first place of complex tables and third place of all tables in Task-B of ICDAR 2021 Competition on Scientific Literature Parsing. Extensive experiments on other publicly available datasets further demonstrate the effectiveness of our proposed approach.

*Keywords:* Table structure recognition, Self-regression, Attention mechanism, Encoder-decoder

## 1. Introduction

In this age of knowledge and information, documents are a very important source of information for many different cognitive processes such as knowledge database creation, optical character recognition (OCR), graphic understanding, document retrieval, etc. Automatically processing the information embedded in these documents is crucial. Numerous efforts have been made in the past to automatically extract the relevant information from documents [1, 2, 3]. As a particular entity, the tabular structure is very commonly encountered in documents. These tabular structures convey some of the most important information in a very concise form. Therefore, they are extremely prevalent in domains like finance, administration, research, and even archival documents. Table structure recognition (TSR) aims to recognize the table internal structure to the machine readable data mainly presented in two formats: logical structure and physical structure [4]. More concretely, logical structure only contains every cell's row and column spanning information, while the physical one additionally contains bounding box coordinates of cells. As a result, table structure recognition as a precursor to contextual table understanding will be useful in a wide range of applications [1, 2, 3].

Table structure recognition is a challenging problem due to the complex structure and high variability in table layouts. A spanning cell is a table cell that occupies at least two rows or columns. If a table contains spanning cells, it is called a complex table, as shown in Figure 1. Although significant efforts have been made in the past to recognize the internal structure of tables through an automated process [3, 5, 6], most of these methods [3, 7] only focus on simple tables and are hard to accurately recognize the structure of complex tables. The spanning cells usually contain more important semantic information than other simple cells, because they are more likely to be table headers in a table. The table header is crucial to understand the table. Therefore, more needs to be done for recognizing the structure of complex tables.

Recently, many works [8, 9, 1] have demonstrated the significant impact of

| System | Simple | Complex | All |
| --- | --- | --- | --- |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

(a)

| System | Simple | Complex | All |
| --- | --- | --- | --- |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

(b)

| System | TEDS | | |
| --- | --- | --- | --- |
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

(c)

| System | TEDS | | |
| --- | --- | --- | --- |
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

(d)

Figure 1: An intuitive comparison between simple and complex tables. The example of the simple table is shown in (a), and (b) is its real structure. The example of the complex table is shown in (c), and (d) is its real structure. Note that in (d), the cells with the contents of "System" and "TEDS" occupy multiple rows or multiple columns, so it is a complex table.

using visual and textual representations in a joint framework. However, most previous methods [3, 6, 10] in table structure recognition only use the spatial or visual features without considering the textual information of each table cell. The structures of some tables have a certain ambiguity from the visual appearance, especially for table cells which contain multi-line contents. Therefore, to recognize the table structure accurately, it is inevitable to take advantage of the cross-modality nature of visually-rich table images, where visual and textual information should be jointly modeled. In our work, we design vision module and text module in our embedder to extract visual features and textual features, respectively, and achieve a higher recognition accuracy.

Most existing literature [10, 11, 12] on table structure recognition depends on extraction of meta-information from the pdf document or the OCR models to extract low-level layout features from the image. Nevertheless, these methods fail to extend to scanned documents due to the absence of meta-information or errors made by the OCR, when there is a wide variety in table layouts and text organization. In our work, we address the problem of table structure recognition by directly operating over table images with no dependency on meta-information

or OCR.

In this study, we introduce Split, Embed and Merge (SEM), an accurate table structure recognizer as shown in Figure 2. Considering that the table is composed of a set of table cells and each table cell is composed of one or more basic table grids, we deem table grids as the basic processing units in our framework. Therefore, we design the pipeline of SEM as follows: 1) divide table into basic table grids 2) merge them to recover the table cells. The final table structure can be obtained by parsing all table cells. As a consequence, SEM mainly has three components: splitter, embedder and merger. The splitter, which is actually a fully convolutional network (FCN) [13], is first applied to predict the fine grid structure of the table as shown in the upper-right of Figure 2. The embedder as a feature extractor embeds vision and plain text contained in a table grid into a feature vector. More specifically, we use the RoIAlign [14] to extract the visual features from the output of the backbone, and extract textual features using the off-the-shelf recognizer [15] and the pretrained BERT [16] model. Finally, the merger which is a Gated Recurrent Unit (GRU) decoder will predict the gird merged results step by step based on the grid-level features extracted by the embedder. For each predicted merged result, the attention mechanism built into the merger scans the entire grid elements and predicts which grids should be merged at the current step. The proposed method can not only process simple tables well, but also complex tables. The ambiguity problem of the table structure recognition based on visual appearance can be alleviated through our embedder. Moreover, SEM is able to directly operate over table images, which enhances the applicability of the system (to both PDFs and images).

The main contributions of this paper are as follows:

- We present an accurate table structure recognizer, Split, Embed and Merge (SEM), to recognize the table structure. The designed merger can accurately predict table structure based on the fine grid structure of the table. This proposed new method can not only process simple tables well, but also complex tables.

- We demonstrate that jointly modeling the visual and textual information in the table will further boost model performance. Through visualization in experiments, the ambiguity problem of the table structure recognition can be alleviated based on our multimodality features.

- Based on our proposed method, we won the first place of complex tables and the third place of all tables in Task-B of ICDAR 2021 Competition on scientific literature parsing. In addition, we also achieved the results with an average F1-Measure of 97.11% and 95.72% in SciTSR and SciTSR-COMP datasets, respectively, demonstrating the effectiveness of our method.

## 2. Related Work

### 2.1. Table Structure Recognition

Analyzing tabular data in unstructured documents mainly focuses on three problems: i) table detection: localizing the bounding boxes of tables in documents [17, 18], ii) table structure recognition: parsing only the structural (row and column layout) information of tables [3, 6, 19], and iii) table recognition: parsing both the structural information and content of table cells [5]. In this study, we mainly focus on table structure recognition. Most early proposed methods [20, 21, 22] are based on heuristics. While these methods were primarily dependent on hand-crafted features and heuristics (horizontal and vertical ruling lines, spacing and geometric analysis).

Due to the rapid development of deep learning and the massive amounts of tabular data in documents on the Web, many deep learning-based methods [3, 5, 6, 10, 7], which are robust to the input type (whether being scanned images or native digital), have also been presented to understand table structures. These also do not make any assumptions about the layouts, are data-driven, and are easy to fine-tune across different domains. [3, 7] utilize recently published insights from semantic segmentation [13] research for identifying rows, columns, and cell positions within tables to recognize table structures. However,

[3, 7] do not consider the complex tables containing spanning cells, so that they cannot handle the structure recognition of complex tables well. GraphTSR [10] proposes a novel graph neural network for recognizing the table structure in PDF files and can recognize the structure of complex tables. GraphTSR takes the table cells as input which means that it fails to generalize well due to the absence of meta-information or errors made by the OCR. EDD [5] treats table structure recognition as a task similar to img2latex [15, 23]. EDD directly generates the HTML tags that define the structure of the table through an attention-based structure decoder. [6] presents the TabStructNet for table structure recognition that combines cell detection and interaction modules to localize the cells and predict their row and column associations with other detected cells. Compared with the aforementioned methods, our method SEM not only takes table images as input, but also can recognize the structure of complex tables well.

*2.2. Attention Mechanisms*

Given a query element and a set of key elements, an attention function can adaptively aggregate the key contents according to attention weights, which measure the compatibility of query-key pairs. The attention mechanisms as an integral part of models enable neural networks to focus more on relevant elements of the input than on irrelevant parts. They were first studied in natural language processing (NLP), where encoder-decoder attention modules were developed to facilitate neural machine translation [24, 25, 26, 27]. In particular, self-attention, also called intra-attention, is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence. Self-attention has been used successfully in a variety of tasks including reading comprehension, abstractive summarization, and textual entailment. The landmark work, Transformer [27], presents the transduction model relying entirely on self-attention to compute representations of its input and output, and substantially surpasses the performance of past work.

The success of attention modeling in NLP [24, 26, 27] has also led to its adoption in computer vision such as object detection [28, 29], semantic seg-

6

mentation [30, 31], image captioning [32] and text recognition [15, 33], etc. DETR [28] completes the object detection by adoptting an encoder-decoder architecture based on transformers [27] to directly predict a set of object bounding boxes. In order to capture contextual information, especially in the long range, [31] proposes the point-wise spatial attention network (PSANet) to aggregate long-range contextual information in a flexible and adaptive manner. Mask TextSpotter v2 [33] applies a spatial attentional module for text recognition, which alleviates the problem of character-level annotations and improves the performance significantly. In our work, we apply the transformers to capture the long-range dependencies on grid-level feature and build attention mechanisms into our merger to predict which gird elements should be merged together to recover table cells.

### 2.3. Multimodality

Several joint learning tasks such as image captioning [34, 9], visual question answering [35, 36, 8], and document semantic structure extraction [1] have demonstrated the significant impact of using visual and textual representations in a joint framework. [9] aligned parts of visual and language modalities through a common, multimodal embedding, and used the inferred alignments to learn to generate novel descriptions of image regions. [8] proposed a novel model, Multimodal Multi-Copy Mesh (M4C), for the TextVQA task based on a multimodal transformer architecture accompanied by a rich representation for text in images and achieved the state-of-the-art. [1] considered document semantic structure extraction as a pixel-wise segmentation task, and presented a unified model, Multimodal Fully Convolutional Network (MFCN). MFCN classifies pixels based not only on their visual appearance, as in the traditional page segmentation task, but also on the content of underlying text. In our work, we take a full consideration of the plain text contained in table images, and design the embedder to extract both visual and textual features at the same time. The experiments also prove that more accurate results will be obtained when providing additional textual information on visual clues.
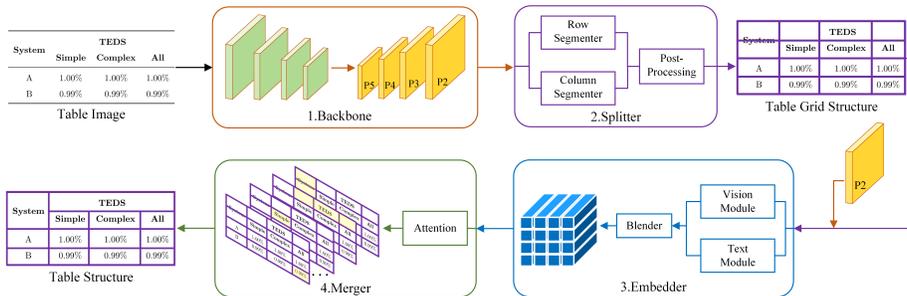
7

Figure 2: **SEM pipeline** The backbone is applied to extract the feature maps from the table image. The splitter uses the backbone features to predict a set of basic table grids. The embedder extracts the region features corresponding to each basic table gird. The merger predicts which grid elements need to be merged to recover the table cells.

## 3. Method

The overall pipeline of our system is shown in Figure 2. The modified ResNet-34 [37] with FPN [38] as our backbone is first applied to the input table image to extract multi-level feature maps. The splitter takes the output of the backbone as input and predicts the fine grid structure of the table. The table grid structure is in the form of row and column separators that span the entire image as shown in the upper-right of Figure 2. The following embedder extracts the feature representation of each basic table grid. Finally, based on the grid-level features extracted by the embedder, the merger with the attention mechanism will predict which grids should be merged step by step. The table structure can be recovered based on the merged results. In the following subsections, three main components in our system, namely, the splitter, the embedder and the merger, will be elaborated.

### 3.1. Splitter

Different from the method in [5], performing table structure prediction on the image-level features, we believe that using table grids as the basic processing units will be more reasonable, so we design the splitter to predict the basic table grid pattern. Inspired by the segmentation-based methods [39, 40] in

8

Figure 3: The illustration of the splitter. The splitter takes a feature map as input and predicts the potential regions of the table row/column separators, which are the green masks in the table images. The following post-processing is used to extract the basic table grids according to the segmentation results from the segmenters.

the field of text detection and the FCN [13] in image segmentation, we refer to the potential regions of the table row/column separators as the foreground and design the splitter which contains two separate row/column segmenters to predict the table row/column separator maps $\hat{\mathbf{S}}^{\mathrm{row}}/\hat{\mathbf{S}}^{\mathrm{col}}$ as shown in Figure 3. $\hat{\mathbf{S}}^{\mathrm{row}}/\hat{\mathbf{S}}^{\mathrm{col}} \in \mathbb{R}^{H \times W}$ and $H \times W$ is the size of the input image.

Each segmenter is actually the fully convolutional network which contains a convolutional layer, ReLU and a convolutional layer. As some table row/column separator regions are quite slender, it is important to ensure segmentation results have a high resolution. The kernel size and the stride of each convolutional layer in the segmenter is set to $3 \times 3$ and 1, respectively, which keeps the same spatial resolution of the input and the output. Moreover, we modify the ResNet-34 by setting the stride of the first convolutional layer with $7 \times 7$ kernel size to 1, and remove the adjacent pooling layer to guarantee the resolution of the lowest-level feature map is consistent with the input image. We strongly believe that rich semantics extracted by deeper layers can help with obtaining more accurate segmentation results, so we add a top-down path [38] in our backbone to enrich semantics in feature maps. Finally, the backbone generates a feature pyramid with four feature maps {P2, P3, P4, P5}, whose final output strides are 1, 2, 4, 8, respectively. The number of channels in the feature maps is $D$. We take P2 as the input of the splitter.

9

The loss function is defined as follows:

$$\mathcal{L}_s^{\text{row}} = \sum_{j=1}^{H} \sum_{i=1}^{W} \frac{L(\hat{S}_{i,j}^{\text{row}}, S_{i,j}^{\text{row}})}{\sum_{j=1}^{H} \sum_{i=1}^{W} S_{i,j}^{\text{row}}} \tag{1}$$

$$\mathcal{L}_s^{\text{col}} = \sum_{j=1}^{H} \sum_{i=1}^{W} \frac{L(\hat{S}_{i,j}^{\text{col}}, S_{i,j}^{\text{col}})}{\sum_{j=1}^{H} \sum_{i=1}^{W} S_{i,j}^{\text{col}}} \tag{2}$$

in which

$$L(x,y) = -(y\log(\sigma(x)) + (1-y)\log(1-\sigma(x))) \tag{3}$$

where $\mathbf{S}^{\text{row}}/\mathbf{S}^{\text{col}}$ denotes the ground-truth of the table row/column separator map. $S_{i,j}^{\text{row}}/S_{i,j}^{\text{col}}$ is 1 if the pixel in $i^{th}$ column and $j^{th}$ row belongs to the table row/column separator region, otherwise 0. The $\sigma$ is the sigmoid function.

The goal of our post-processing is to extract table row/column lines from the table row/column separator map as shown in Figure 3. Specifically, we first apply the sigmoid function to the predicted segmentation map $\hat{\mathbf{S}}^{\text{row}}/\hat{\mathbf{S}}^{\text{col}}$ and average them by column/row size to obtain the $\bar{\mathbf{S}}^{\text{row}}/\bar{\mathbf{S}}^{\text{col}}$ as illustrated in Eq. 4 5, where $\bar{\mathbf{S}}^{\text{row}} \in \mathbb{R}^{H \times 1}$ and $\bar{\mathbf{S}}^{\text{col}} \in \mathbb{R}^{1 \times W}$. Then we binarize the $\bar{\mathbf{S}}^{\text{row}}/\bar{\mathbf{S}}^{\text{col}}$ into $\tilde{\mathbf{S}}^{\text{row}}/\tilde{\mathbf{S}}^{\text{col}}$, $\tilde{S}_j^{\text{row}}/\tilde{S}_i^{\text{col}} = 1$ indicating this row/column is a potential table line. For the block that is equal to 1 in $\tilde{\mathbf{S}}^{\text{row}}/\tilde{\mathbf{S}}^{\text{col}}$, we select the row/column with the maximum value of the corresponding block in $\bar{\mathbf{S}}^{\text{row}}/\bar{\mathbf{S}}^{\text{col}}$ as the final table row/column line.

$$\bar{S}_j^{\text{row}} = \frac{1}{W} \sum_{i}^{W} \sigma\left(\hat{S}_{i,j}^{\text{row}}\right) \tag{4}$$

$$\bar{S}_i^{\text{col}} = \frac{1}{H} \sum_{j}^{H} \sigma\left(\hat{S}_{i,j}^{\text{col}}\right) \tag{5}$$

We can easily obtain a set of bounding boxes $\mathbf{G}$ of table grids from the table row/column lines. $\mathbf{G} \in \mathbb{R}^{(M \times N) \times 4}$, where $M$, $N$ are the number of rows and columns occupied by the table grid structure, respectively. More specifically, each bounding box can be precisely defined by $(x_1, y_1, x_2, y_2)$, where $(x_1, y_1)$ corresponds to the position of the upper left in the bounding box, and $(x_2, y_2)$ represents the position of the lower right.
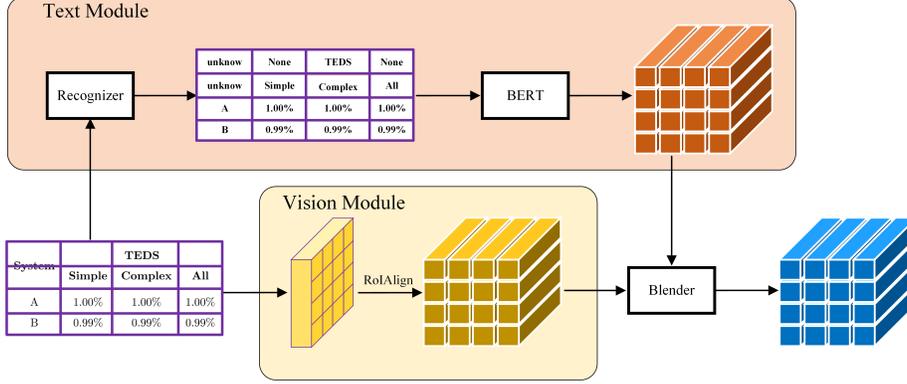
10

## 3.2. Embedder



Figure 4: The illustration of the embedder. It is composed of vision module (VM), text module (TM) and blender module (BM). The embedder extracts the gird-level visual and textual features from VM and TM, respectively. Finally, the BM fuses the both features.

The embedder aims to extract the feature representations of each grid. [9, 8] have demonstrated the effectiveness of taking advantage of the multi-modality. Different from the previous table structure recognition methods [5, 6, 19] which mostly recover the table structure based on the visual modality, we fuse the output features for each basic table grid from both visual and textual modalities. Therefore, we design the vision module and text module in the embedder to extract visual features $\mathbf{E}^v$ and textual features $\mathbf{E}^t$, respectively, and fuse both features to produce the final grid-level features $\mathbf{E}$ through the blender module. $\mathbf{E}^v \in \mathbb{R}^{(M \times N) \times D}$, $\mathbf{E}^t \in \mathbb{R}^{(M \times N) \times D}$ and $\mathbf{E} \in \mathbb{R}^{(M \times N) \times D}$, where $D$ represents the number of feature channels.

As shown in Figure 4, the vision module takes the image-level feature map P2 from the FPN and the well-divided table grids $\mathbf{G}$ obtained from the splitter as input. It applies the RoIAlign [14] to pool a fixed size $R \times R$ feature map $\hat{\mathbf{E}}_i^v$ for each table grid.

$$\hat{\mathbf{E}}_i^v = \mathrm{RoIAlign}_{R \times R}(\mathrm{P2}, \mathbf{G}_i) \quad \forall i = \{1, ..., M \times N\} \tag{6}$$

11

where $\hat{\mathbf{E}}_i^v \in \mathbb{R}^{R \times R \times D}$. The final visual features $\mathbf{E}_i^v$ are obtained according to:

$$\mathbf{E}_i^v = \text{FFN}(\hat{\mathbf{E}}_i^v) \quad \forall i = \{1, ..., M \times N\} \tag{7}$$

in which

$$\text{FFN}(\mathbf{x}) = \max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2 \tag{8}$$

where FFN [27] is actually two linear transformations with a ReLU activation in between. $\mathbf{x} \in \mathbb{R}^{d_{\text{in}}}$, $\mathbf{W}_1 \in \mathbb{R}^{d_{\text{in}} \times d_{\text{ff}}}$, $\mathbf{b}_1 \in \mathbb{R}^{d_{\text{ff}}}$, $\mathbf{W}_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{out}}}$, $\mathbf{b}_2 \in \mathbb{R}^{d_{\text{out}}}$. The dimensionality of input and output is $d_{\text{in}}$ and $d_{\text{out}}$, and the inner-layer has dimensionality $d_{\text{ff}}$. Here we set $d_{\text{ff}} = d_{\text{out}}$ in default.

The table image is both visually-rich and textual-rich, so it is necessary to make full use of the textual information in the table to achieve a more accurate table structure recognizer. As shown in the text module of Figure 4, we apply the off-the-shelf recognizer [15] to obtain a sequence of $M \times N$ contents for all table grids, and embed contents into corresponding feature vectors $\hat{\mathbf{E}}^t$ using a pretrained BERT model [16]. $\hat{\mathbf{E}}^t \in \mathbb{R}^{(M \times N) \times B}$, where $B$ is the feature vector dimension of the BERT. It's worth noting that both the recognizer and the BERT do not update the parameters during the training phase. The final textual features $\mathbf{E}^t$ are obtained by applying FFN again to fine-tune the extracted textual features $\hat{\mathbf{E}}^t$ to make it more suitable for our network.

$$\mathbf{E}_i^t = \text{FFN}(\hat{\mathbf{E}}_i^t) \quad \forall i = \{1, ..., M \times N\} \tag{9}$$

The blender module in Figure 4 is to fuse the visual features $\mathbf{E}^v$ and textual features $\mathbf{E}^t$, and its specific process is as follows:

1) For each basic table grid, we first obtain the intermediate results $\hat{\mathbf{E}}_i$ according to :

$$\hat{\mathbf{E}}_i = \text{FFN}\left( \begin{bmatrix} \mathbf{E}_i^v \\ \mathbf{E}_i^t \end{bmatrix} \right) \quad \forall i \in [1, ..., M \times N] \tag{10}$$

where $[\cdot]$ is the concatenation operation. The input and output dimensionality of the FFN is $2D$ and $D$, respectively.

2) So far, the features of each basic table grid are still independent of each other, especially for textual features. Therefore, we introduce the transformer [27] to capture long-range dependencies on table grid elements. We take the features $\hat{\mathbf{E}}$ as query, key and value which are required by the transformer. The output of the transformer as final grid-level features $\mathbf{E}$ have a global receptive field.
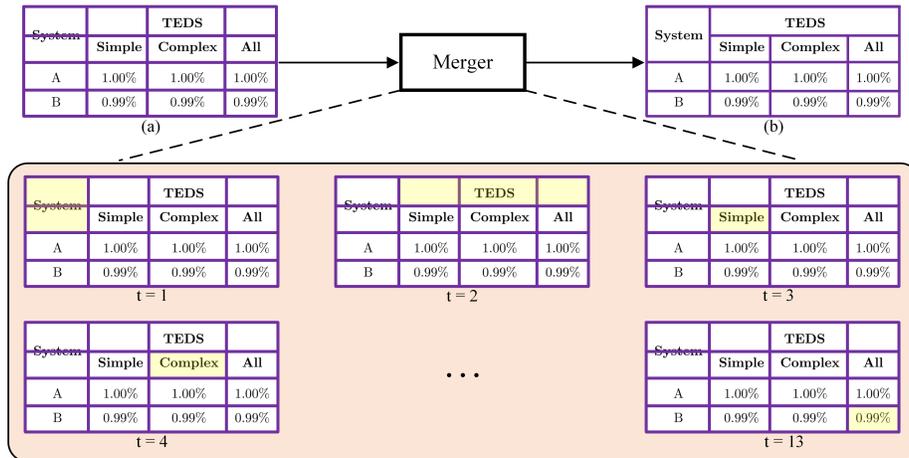
| System | TEDS | | |
|---|---|---|---|
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

(a)

Merger

| System | TEDS | | |
|---|---|---|---|
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

(b)

| System | TEDS | | |
|---|---|---|---|
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

t = 1

| System | TEDS | | |
|---|---|---|---|
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

t = 2

| System | TEDS | | |
|---|---|---|---|
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

t = 3

| System | TEDS | | |
|---|---|---|---|
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

t = 4

· · ·

| System | TEDS | | |
|---|---|---|---|
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

t = 13

Figure 5: The illustration of the merger. The yellow masks in lower part indicate which table grid elements should be merged in each time step.

### 3.3. Merger

The merger is an RNN that takes the grid-level features $\mathbf{E}$ as input and produces a sequence of merged maps $\mathbf{M}$ as shown in Figure 5. Here we choose the Gated Recurrent Unit (GRU) [41], an improved version of simple RNN.

$$\mathbf{M} = \{\mathbf{m}_1, \mathbf{m}_2, ..., \mathbf{m}_C\} \tag{11}$$

where $C$ is the length of a predicted sequence. Each merged map $\mathbf{m}_t$ is a $(M \times N)$-dimension vector, the same size as the element of $\mathbf{E}$, and the value of each grid element $m_{ti}$ is 1 or 0, indicating whether the $i^{th}$ grid element belongs to the $t^{th}$ cell or not. The cells that span multiple rows or columns can be recovered according to $\mathbf{M}$.

13

Figure 6: The illustration of the attention mechanism. The prediction of current hidden state $\hat{\mathbf{h}}_t$ and the grid-level features $\mathbf{E}$ is used as query and key, respectively.

Inspired by the successful applications of attention mechanism in img2latex [15, 42], text recognition [43, 44], machine translation [27], etc., we build the attention mechanism into our merger and achieve promising results. For the merged map $\mathbf{m}_t$ decoding, we compute the prediction of current hidden state $\hat{\mathbf{h}}_t$ from previous context vector $\mathbf{c}_{t-1}$ and its hidden state $\mathbf{h}_{t-1}$:

$$\hat{\mathbf{h}}_t = \text{GRU}(\mathbf{c}_{t-1}, \mathbf{h}_{t-1}) \tag{12}$$

Then we employ an attention mechanism with $\hat{\mathbf{h}}_t$ as the query and grid-level features $\mathbf{E}$ as both key and the value:

$$\mathbf{m}_t = f_{att}(\mathbf{E}, \hat{\mathbf{h}}_t) \tag{13}$$

$$\mathbf{c}_t = \frac{\mathbf{m}_t}{\|\mathbf{m}_t\|_1} \mathbf{E} \tag{14}$$

where $\|\cdot\|_1$ is the vector 1-norm. As shown in Figure 6, we design $f_{att}$ function

14

as follows:

$$\mathbf{F} = \mathbf{Q} * \sum_{l=1}^{t-1} \mathbf{m}_l \tag{15}$$

$$\hat{m}_{ti} = \boldsymbol{\nu}^{\mathrm{T}} \tanh(\mathbf{W}_{att}\hat{\mathbf{h}}_t + \mathbf{U}_{att}\mathbf{e}_i + \mathbf{U}_F\mathbf{f}_i) \tag{16}$$

$$m_{ti} = \mathrm{Binarize}(\hat{m}_{ti}) \tag{17}$$

in which

$$\mathrm{Binarize}(x) = \begin{cases} 1 & \text{if } \sigma(x) > 0.5 \\ 0 & \text{otherwise} \end{cases} \tag{18}$$

where $*$ denotes a convolution layer, $\sum_{l=1}^{t-1} \mathbf{m}_l$ denotes the sum of past determined grids, $\hat{m}_{ti}$ denotes the output energy, $\mathbf{f}_i$ denotes the element of $\mathbf{F}$, which is used to help append the history information into standard attention mechanism. It's worth noting that the attention mechanism is completed on the grid-level features. For each cell, it is quite clear which grid elements belong to it. Therefore, unlike the previous methods [15, 23] using the softmax to obtain the attention probability, we use the Binarize Eq. 18 to calculate. Moreover, we find that the model is difficult to converge when using the softmax.

With the context vector $\mathbf{c}_t$, we compute the current hidden state:

$$\mathbf{h}_t = \mathrm{GRU}(\mathbf{c}_t, \hat{\mathbf{h}}_t) \tag{19}$$

The training loss of the merger is defined as follows:

$$\mathcal{L}_{\mathrm{m}} = \sum_t \sum_i \frac{L(\hat{m}_{ti}, y_{ti})}{C\|\mathbf{y}_t\|_1} \tag{20}$$

where function $L$ has been defined in Eq. 3, $C$ is the length of a predicted sequence and $y_{ti}$ denotes the ground-truth of cell's grid elements. $y_{ti}$ is 1 if the $i^{th}$ grid element belong to the cell of time step $t$, otherwise 0.

### 3.4. Post-Processing

Through the merger, we can obtain the spanning of each table cell along the rows and columns. By combining the spanning information and the prediction

results of the splitter, which contains the table row/column lines information, the bounding box coordinates of each table cell can be obtained. We match the text content with position to the table cells according to the IOU. The output for every table image finally contains coordinates of predicted cell bounding boxes along with cell spanning information and its content.

## 4. Experiment

### 4.1. Dataset

We use the publicly available table structure datasets — SciTSR [10], SciTSR-COMP [10] and PubTabNet [5] to evaluate the effectiveness of our model. Statistics of these datasets are listed in Table 1.

Table 1: Statistics of the datasets used for our experiments.

| Dataset | SciTSR | SciTSR-COMP | PubTabNet |
|---|---|---|---|
| Train | 12k | - | 500k |
| Validation | - | - | 9k |
| Test | 3k | 716 | 9k |

1) **SciTSR** [10] is a large-scale table structure recognition dataset, which contains 15,000 tables in PDF format as well as their corresponding high quality structure labels obtained from LaTeX source files. SciTSR splits $12,000$ for training and $3,000$ for testing. Furthermore, to reflect the model's ability of recognizing complex tables, [10] extracts all the 716 complex tables from the test set as a test subset, called SciTSR-COMP. It's worth noting that SciTSR provides the text contents with positions for each table image, but not with being aligned with the table cells. However, in our model, we need the text position in each table cell to generate the labels of splitter. Therefore, we apply the data preprocessing [1] to align the text information with the table cells.

---

[1] https://github.com/ZZR8066/SciTSR

2) **PubTabNet** [5] contains over 500k training samples and 9k validation samples. PubTabNet [5] annotates each table image with information about both the structure of table and the text content with position of each non-empty table cell. Moreover, nearly half of them are complex tables which have spanning cells in PubTabNet.

### 4.2. Label Generation

**Label of Splitter** We use the annotation, namely the text content with position being aligned to each table cell, to generate the ground-truth of the table row/column separator map $\mathbf{S}^{\text{row}}/\mathbf{S}^{\text{col}}$ for the splitter. The $\mathbf{S}^{\text{row}}/\mathbf{S}^{\text{col}}$ is designed to maximize the size of the separator regions without intersecting any non-spanning cell content, as shown in Figure 7. Different from traditional notion of cell separators, which for many tables are thin lines with only a few pixels. Predicting small regions is more difficult than predicting large regions. In the case of unlined tables, the exact location of the cell separator is ill-defined.

**Label of Merger** Since we obtain the label of the splitter, we could divide the table into a set of basic table grids as shown in the upper-right of Figure 2. The original table structure annotation can provide which grids each table cell occupies in the basic table grid pattern. We arrange the table cells in a top-to-bottom and left-to-right way and use the grids occupied by each cell as the prediction target for a certain time step of the merger.



| System | TEDS | | |
| --- | --- | --- | --- |
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

$\mathbf{S}^{\text{row}}$

| System | | TEDS | |
| --- | --- | --- | --- |
| | Simple | Complex | All |
| A | 1.00% | 1.00% | 1.00% |
| B | 0.99% | 0.99% | 0.99% |

$\mathbf{S}^{\text{col}}$

Figure 7: Example of the ground-truth of table row/column separator map for the splitter. The red mask is the table row/column separator region.

### 4.3. Metric

We use both F1-Measure [45] and Tree-Edit-Distance-based Similarity (TEDS) metric [5], which are commonly adopted in table structure recognition literature and competitions, to evaluate the performance of our model for recognition of the table structure.

In order to use the F1-Measure, the adjacency relationships among the table cells need to be detected. F1-Measure measures the percentage of correctly detected pairs of adjacent cells, where both cells are segmented correctly and identified as neighbors.

The TEDS metric was recently proposed in [5]. While using the TEDS metric, we need to present tables as a tree structure in the HTML format. Finally, TEDS between two trees is computed as:

$$\mathrm{TEDS}(T_a, T_b) = 1 - \frac{\mathrm{EditDist}(T_a, T_b)}{\max(|T_a|, |T_b|)} \tag{21}$$

where $T_a$ and $T_b$ are the tree structure of tables in the HTML formats. EditDist represents the tree-edit distance [46], and $|T|$ is the number of nodes in $T$.

### 4.4. Implementation Details

The modified ResNet-34 [37] as our backbone is pre-trained on ImageNet [47]. The number of FPN channels is set to $D = 256$. The pool size $R \times R$ of RoIAlign in vision module is set to $3 \times 3$. The recognizer [15] is pre-trained on 35M table cell images, which are cropped from 500k table images in the PubTabNet dataset [5], and the success rate of word predictions per table cell reaches 94.1%. The BERT [16] we used is from the transformers package [2]. The hidden state dimension in the merger is set to 256.

The training objective of our model is to minimize the segmentation loss (Eq. 1, Eq. 2) and the cell merging loss (Eq. 20). The objective function for optimization is shown as follows:

$$O = \lambda_1 \mathcal{L}_s^{\mathrm{row}} + \lambda_2 \mathcal{L}_s^{\mathrm{col}} + \lambda_3 \mathcal{L}_m \tag{22}$$

---

[2]https://github.com/huggingface/transformers

In our experiments, we set $\lambda_1 = \lambda_2 = \lambda_3 = 1$. We employ the ADADELTA algorithm [48] for optimization, with the following hyper parameters: $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\varepsilon = 10^{-9}$. We set the learning rate using the cosine annealing schedule [49] as follows:

$$\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\frac{T_{cur}}{T_{max}}\pi)) \tag{23}$$

where $\eta_t$ is the updated learning rate. $\eta_{min}$ and $\eta_{max}$ are the minimum learning rate and the initial learning rate, respectively. $T_{cur}$ and $T_{max}$ are the current number of iterations and the maximum number of iterations, respectively. Here we set $\eta_{min} = 10^{-6}$ and $\eta_{max} = 10^{-4}$.

Our model SEM is trained and evaluated with table images in original size. We use the NVIDIA TESLA V100 GPU with 32GB RAM memory for our experiments and the batch-size of 8. The whole framework was implemented using PyTorch.

### 4.5. Visualization

In this section, we visualize the segmentation results of the spliter and show how the merger recovers the table cells from the table grid elements through attention visualization.

**Visualization of Splitter** We refer the potential regions of the table row (column) separators as the foreground as shown in Figure 7, and design the splitter which is actually a fully convolutional network (FCN) to predict the foreground in table images. As shown in the first two rows of Figure 8, we can obtain accurate segmentation results through the splitter. The fine grid structure of the table can be obtained by post-processing as shown in the third row of Figure 8. It is worth noting that the example table in Figure 8 (a) is the simple table, while others are complex tables. We can find that the structure of the simple table has been recovered correctly through the splitter from the third row of Figure 8. However, the structure of complex tables is not complete and still needs to be processed. Therefore, we design the following embedder

19

Figure 8: The visualization results from our system on table images of the SciTSR dataset. **First Row:** the green masks are the segmentation results of the row segmenter in the splitter. **Second Row:** the green masks are the segmentation results of the column segmenter in the splitter. **Third Row:** the blue lines indicate the boundaries of the basic table grids which are extracted through post-processing from both row and column segmentation results. **Fourth Row:** the blue lines indicate the boundaries of the table cells which are the merged results from the merger.

and merger to recover the structure of complex tables based on the outputs of the splitter.

**Visualization of Merger** In order to recover the table cells, we build the attention mechanism into our merger to predict which grid elements should be merged step by step. The merged result in each step is a binary map, and the table cell can be recovered by merging the elements that are 1 in the binary map. Taking the table of Figure 8 (b) as a example, the attention mechanism is visualized in Figure 9. The cell with the content of "Number of modules" in Figure 9 occupies the first row of basic table grids. Our merger correctly predicts the structure of this cell through the attention mechanism as shown in the first time step of Figure 9.

Figure 9: The visualization of the attention mechanism in the merger on the table images of the SciTSR dataset. The blue lines are the prediction of table grid structure from the splitter. The green mask in the table image denotes which grid elements should be merged for each time step.

## 4.6. Ablation Study

In order to investigate the effect of each component, we conduct ablation experiments through several designed systems as shown in Table 2. The model is not modified except the component being tested.

Table 2: Comparison among systems from T1 to T4. Attributes for comparison include: 1) employing the splitter; 2) using the vision module (VM) in the embedder; 3) using the text module (TM) in the embedder; 4) employing the merger.

| System | Splitter | Embedder | | Merger |
| | | VM | TM | |
| --- | --- | --- | --- | --- |
| T1 | ✓ | - | - | - |
| T2 | ✓ | - | ✓ | ✓ |
| T3 | ✓ | ✓ | - | ✓ |
| T4 | ✓ | ✓ | ✓ | ✓ |

**The Number of Transformer Layers** We measure the performance of T1-T4 with different numbers of transformer layers in the embedder. We try from 0 to 3 as shown in Figure 10. When Num = 0 in Figure 10, it means the transformer layer is removed. In the T3 configuration, only the vision module

21

(VM) in the embedder is used to extract the visiual features to represent each grid element. Also there is not much gap whether the transformer layer is added or not. Through a series of convolutional layers, the backbone features P2 already has a certain receptive field. Therefore, it is not significant to add the transformer layers while the VM has pooled each grid features from P2. It is worth noting that when the Num is greater than 0, the performance of the designed system T2 outperforms the model without the transformer layer. This is because the transformer layer here can capture the semantical dependencies among all table grid elements. As our final system, T4 achieves the best result when Num = 1, so we set Num = 1 in subsequent experiments by default.



Figure 10: Performance by varying number of transformer layers in T2, T3, T4 on the SciTSR test dataset.

**The Effectiveness of the Merger** In Table 3, we show the F1-Measure of systems T1-T4 on SciTSR and SciTSR-COMP datasets. Almost 76.3% of the tables are simple tables in SciTSR test dataset, and all are complex tables in the SciTSR-COMP dataset. The performance gap between T1 and other systems (T2-T4) is remarkable on the SciTSR dataset, but the gain is more significant on the SciTSR-COMP dataset, e.g., almost 6% in F1-Measure from T1 to T4. This is because all table cells have only one table grid in the simple table, which means that the table grid structure is the table structure. However there are

Table 3: Comparison of F1-Measure among different systems in Table 2 on SciTSR and SciTSR-COMP datasets.

| System | SciTSR | | | SciTSR-COMP | | | FPS |
|--------|--------|--------|--------|-------------|--------|--------|-------|
|        | P      | R      | F1     | P           | R      | F1     |       |
| T1     | 96.69  | 94.15  | 95.40  | 93.81       | 96.06  | 89.77  | 16.47 |
| T2     | 96.63  | 94.36  | 95.48  | 94.15       | 88.04  | 90.99  | 2.00  |
| T3     | 97.40  | 95.97  | 96.68  | 96.52       | 93.82  | 95.15  | 3.65  |
| T4     | 97.70  | 96.52  | 97.11  | 96.80       | 94.67  | 95.72  | 1.94  |

some table cells have more than one table grids in the complex table. Therefore, the designed system T1 can only process simple tables well by using splitter to predict the fine grid structure of table, and T2-T4 have the ability to recover the structure of the complex tables through the merger. The comparison of T1 with T2, T3, T4 on the SciTSR-COMP dataset demonstrates the effectiveness of the merger.

**Vision and Language Modalities** In order to evaluate the effect of each modality, we design the systems T2, T3, T4 as shown in Table 2. Each system uses vision module (VM), text module (TM) or both in the embedder. The experiment results on SciTSR and SciTSR-COMP are shown in Table 3. Compared with T4, systems T2 and T3 that only use TM or VM are sub-optimal. When both TM and VM are used, the system (T4) performance reaches the best. As shown in Figure 11, although the predictions of table grid structure from the splitter in both T3 and T4 are the same, the T3 system which only uses VM is more unstable comparing with T4 which uses both VM and TM in the embedder.

**The Efficiency of Each Component** In order to investigate the efficiency of each component, we compare the frames per second (FPS) of T1-T4 systems as shown in Table 3. From T1 to T2-T4 systems, the speed of the systems is much slower. This is because as the number of table cells increases, the decoding steps of the merger increases. The reason why T2 and T4 are slower than T3

is that the former uses a recognizer to recognize the content in the basic table grid and applies the BERT to extract the corresponding textual features.



Figure 11: The comparison results between designed system T3 and T4. The fisrt column is the results on the SciTSR dataset. The second column is the results on the PubTabNet dataset. **First Row:** the predictions of the table grid structure from the splitter. **Second Row:** the predictions of the table structure from the T3 which only uses the vision module in the embedder. **Third Row:** the predictions of the table structure from the T4 which uses both the vision module and text module in the embedder. Note that the predictions of table grid structure in systems T3 and T4 are the same, and the predictions of table structure in the third row are all totally correct. The red dash boxes denote the different predictions between T3 and T4.

## 4.7. Comparison with State-of-the-art Methods

We compare our method with other state-of-the-art methods on both SciTSR and SciTSR-COMP datasets. The results are shown in Table 4. Our model is

24

Table 4: A performance comparison between our method and other state-of-the-art methods on the SciTSR and SciTSR-COMP datasets.

| Method | SciTSR | | | SciTSR-COMP | | | FPS |
|---|---|---|---|---|---|---|---|
| | P | R | F1 | P | R | F1 | |
| Adobe [3] | 82.9 | 79.6 | 81.2 | 79.6 | 73.7 | 76.5 | - |
| TabbyPDF [50] | 91.4 | 91.0 | 91.2 | 86.9 | 84.1 | 85.5 | - |
| DeepDeSRT [3] | 89.8 | 89.7 | 89.7 | 81.1 | 81.3 | 81.2 | 20.88 |
| GraphTSR [10] | 93.6 | 93.1 | 93.4 | 94.3 | 92.5 | 93.4 | - |
| TabStruct-Net [6] | 92.7 | 91.3 | 92.0 | 90.9 | 88.2 | 89.5 | 0.77 |
| T1 | 96.69 | 94.15 | 95.40 | 93.81 | 96.06 | 89.77 | 16.47 |
| SEM | 97.70 | 96.52 | **97.11** | 96.80 | 94.67 | **95.72** | 1.94 |

trained and tested with default configuration. Comparing with other methods [10, 3, 6], our method achieves state-of-the-art. Similar to DeepDeSRT [3], T1 is actually a segmentation model. We take full consideration of the extremely unbalanced number of foreground and background pixels in the segmentation masks and design a more reasonable segmentation loss to penalize the model during training in Eq. 1 2, which makes the performance of T1 significantly better than DeepDeSRT. It is worth noting that GraphTSR [10] needs the text position in table cells during both the training and testing stage, while our method only takes table images as input during inference. Although the comparison between GraphTSR and our method is not fair, our method still outperforms it to a certain extend. The TabStruct-Net [6] applies a detection network [14] to detect individual cells in a table image, however, it fails to capture empty cells accurately due to the absence of cell content. Our SEM obtains the bounding boxes of cells based on table lines detection, which makes the prediction of empty cells less difficult. In addition, the embedder makes full use of the visual and textual modalities, and the merger enables the model to process complex tables, which ultimately makes our method achieve state-of-the-art. Some table structure recognition results of our and other methods are shown in Figure 12.

Figure 12: Some structure recognition results of our and other methods on table images of the SciTSR dataset. The blue lines denote the prediction of table structure. **First Row:** the results of the DeepDeSRT. **Second Row:** the intermediate cell detection results of the TabStruct-Net. **Third Row:** the results of our method. The predictions of table structure in the third row are all correct.

## 4.8. ICDAR 2021 Competition on Scientific Literature Parsing, Task-B

ICDAR 2021 Competition on Scientific Literature Parsing, Task-B [3] is organized by the IBM company in conjunction with IEEE ICDAR 2021. This

---

[3] https://aieval.draco.res.ibm.com/challenge/40/overview

competition aims to drive the advances in table recognition. Different from the table structure recognition task, we need to recognize not only the structure of the table, but also the content within each cell. Through our method, we can not only predict the structure of the table, but also obtain the position of each cell. Inspired by [51, 52, 33, 53], we use the RoIAlign to pool the features of table cells and append an attention-based recognizer [15] to recognize the content in table cells. Note that the modified models are trained in an end-to-end manner. The single model results of our methods are shown in Table 5.

Table 5: The performance of table recognition on PubTabNet validation set.

| Method | TEDS | | | FPS |
|---|---|---|---|---|
| | Simple | Complex | All | |
| T3 + Recognizer | 94.7 | 92.1 | 93.4 | 1.81 |
| T4 + Recognizer | 94.8 | 92.5 | 93.7 | 1.23 |

Based on the configuration of T3 with a recognizer, we divide our model into three sub-networks, splitter, merger and newly added recognizer, adopting multi-model fusion for each sub-network. Finally, we combine the training set with the validation set for training. The results of the competition are shown in Table 6. Our team is named USTC-NELSLIP, and we won the first place of complex tables and third place of all tables.

*4.9. Error Analysis*

In this section, we show some incorrect table structure recognition results of the SEM as shown in Figure 13. Our splitter occasionally misses or overcuts the basic table grids when the blank space between cells is too large. In the training phase, the merger predicts the spanning information of cells on the correct basic table grid pattern. Therefore, in the inference phase, once the splitter predicts incorrect results, it is difficult for the merger to fix them.

Table 6: Table recognition competition results on PubTabNet final evaluation data set.

| Team Name | TEDS | | |
|---|---|---|---|
| | Simple | Complex | All |
| Davar-Lab-OCR | 97.88 | 94.78 | **96.36** |
| VCGroup | **97.90** | 94.68 | 96.32 |
| **USTC-NELSLIP** | 97.60 | **94.89** | 96.27 |
| YG | 97.38 | 94.79 | 96.11 |
| DBJ | 97.39 | 93.87 | 95.66 |
| TAL | 97.30 | 93.93 | 95.65 |
| PaodingAI | 97.35 | 93.79 | 95.61 |
| anyone | 96.95 | 93.43 | 95.23 |
| LTIAYN | 97.18 | 92.40 | 94.84 |
| EDD | 91.20 | 85.40 | 88.30 |



Figure 13: Some incorrect table structure recognition results of our method. **First Column**: the row segmentation results of the splitter. **Second Column**: the column segmentation results of the splitter. **Third Column**: the predictions of the table grid structure from the splitter. **Fourth Column**: the final prediction results of our method. The red dash boxes denote the incorrect prediction results.

## 5. Conclusion

In this study, we proposed a new method for the table structure recognition, SEM. The proposed method takes images as input with no dependency on meta-information or OCR. It mainly contains three components including splitter, embedder and merger. We first split table images into a set of basic table grids. Then the embedder is used to extract the feature representations of each grid element. Finally, we use the merger with the attention mechanism to predict which grid elements should be merged to recover the table cells. The final table structure can be obtained by parsing all table cells. The method can not only process simple tables well, but also the complex tables. We demonstrate through visualization and experiment results that the attention mechanism built in the merger performs well in predicting which grid elements belong to each cell. To our best knowledge, this is the first time to take a full consideration of the textual information in table images and design the embedder to extract both the visual and the textual features. The ablation studies prove the effectiveness of our embedder. Our method achieves state-of-the-art on both SciTSR and SciTSR-COMP datasets. Based on our method, we won the first place of complex tables and third place of all tables in Task-B of ICDAR 2021 Competition on Scientific Literature Parsing.

## References

[1] X. Yang, E. Yumer, P. Asente, M. Kraley, D. Kifer, C. L. Giles, Learning to extract semantic structure from documents using multimodal fully convolutional neural network (2017).

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

[3] S. Schreiber, S. Agne, I. Wolf, A. Dengel, S. Ahmed, Deepdesrt: Deep learning for detection and structure recognition of tables in document images,

in: 2017 International Conference on Document Analysis and Recognition (ICDAR), Vol. 01, 2017, pp. 1162–1167. `doi:10.1109/ICDAR.2017.192`.

[4] R. Zanibbi, D. Blostein, R. Cordy, A survey of table recognition: Models, observations, transformations, and inferences, Int. J. Doc. Anal. Recognit. 7 (1) (2004) 1–16. `doi:10.1007/s10032-004-0120-9`.

[5] X. Zhong, E. ShafieiBavani, A. Jimeno Yepes, Image-based table recognition: Data, model, and evaluation, in: A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm (Eds.), Computer Vision – ECCV 2020, 2020, pp. 564–580.

[6] S. Raja, A. Mondal, C. V. Jawahar, Table structure recognition using top-down and bottom-up cues, in: A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm (Eds.), Computer Vision – ECCV 2020, 2020, pp. 70–86.

[7] S. A. Siddiqui, P. I. Khan, A. Dengel, S. Ahmed, Rethinking semantic segmentation for table structure recognition in documents, in: 2019 International Conference on Document Analysis and Recognition (ICDAR), 2019, pp. 1397–1402. `doi:10.1109/ICDAR.2019.00225`.

[8] R. Hu, A. Singh, T. Darrell, M. Rohrbach, Iterative answer prediction with pointer-augmented multimodal transformers for textvqa, in: 2020 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 9989–9999. `doi:10.1109/CVPR42600.2020.01001`.

[9] A. Karpathy, L. Fei-Fei, Deep visual-semantic alignments for generating image descriptions, IEEE Transactions on Pattern Analysis and Machine Intelligence 39 (4) (2017) 664–676. `doi:10.1109/TPAMI.2016.2598339`.

[10] Z. Chi, H. Huang, H.-D. Xu, H. Yu, W. Yin, X.-L. Mao, Complicated table structure recognition (2019).

[11] W. Xue, Q. Li, D. Tao, Res2tim: Reconstruct syntactic structures from table images, in: 2019 International Conference on Document Analysis and Recognition (ICDAR), 2019, pp. 749–755. `doi:10.1109/ICDAR.2019.00125`.

in: 2017 International Conference on Document Analysis and Recognition (ICDAR), Vol. 01, 2017, pp. 1162–1167. `doi:10.1109/ICDAR.2017.192`.

[4] R. Zanibbi, D. Blostein, R. Cordy, A survey of table recognition: Models, observations, transformations, and inferences, Int. J. Doc. Anal. Recognit. 7 (1) (2004) 1–16. `doi:10.1007/s10032-004-0120-9`.

[5] X. Zhong, E. ShafieiBavani, A. Jimeno Yepes, Image-based table recognition: Data, model, and evaluation, in: A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm (Eds.), Computer Vision – ECCV 2020, 2020, pp. 564–580.

[6] S. Raja, A. Mondal, C. V. Jawahar, Table structure recognition using top-down and bottom-up cues, in: A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm (Eds.), Computer Vision – ECCV 2020, 2020, pp. 70–86.

[7] S. A. Siddiqui, P. I. Khan, A. Dengel, S. Ahmed, Rethinking semantic segmentation for table structure recognition in documents, in: 2019 International Conference on Document Analysis and Recognition (ICDAR), 2019, pp. 1397–1402. `doi:10.1109/ICDAR.2019.00225`.

[8] R. Hu, A. Singh, T. Darrell, M. Rohrbach, Iterative answer prediction with pointer-augmented multimodal transformers for textvqa, in: 2020 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 9989–9999. `doi:10.1109/CVPR42600.2020.01001`.

[9] A. Karpathy, L. Fei-Fei, Deep visual-semantic alignments for generating image descriptions, IEEE Transactions on Pattern Analysis and Machine Intelligence 39 (4) (2017) 664–676. `doi:10.1109/TPAMI.2016.2598339`.

[10] Z. Chi, H. Huang, H.-D. Xu, H. Yu, W. Yin, X.-L. Mao, Complicated table structure recognition (2019).

[11] W. Xue, Q. Li, D. Tao, Res2tim: Reconstruct syntactic structures from table images, in: 2019 International Conference on Document Analysis and Recognition (ICDAR), 2019, pp. 749–755. `doi:10.1109/ICDAR.2019.00125`.

[12] S. R. Qasim, H. Mahmood, F. Shafait, Rethinking table recognition using graph neural networks (2019).

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

[14] K. He, G. Gkioxari, P. Dollar, R. Girshick, Mask r-cnn, in: Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2017.

[15] J. Zhang, J. Du, S. Zhang, D. Liu, Y. Hu, J. Hu, S. Wei, L. Dai, Watch, attend and parse: An end-to-end neural network based approach to handwritten mathematical expression recognition, Pattern Recognition 71 (2017) 196–206. `doi:https://doi.org/10.1016/j.patcog.2017.06.017`.

[16] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), 2019, pp. 4171–4186. `doi:10.18653/v1/N19-1423`.

[17] M. Göbel, T. Hassan, E. Oro, G. Orsi, Icdar 2013 table competition, in: 2013 International Conference on Document Analysis and Recognition (ICDAR), 2013, pp. 1449–1453. `doi:10.1109/ICDAR.2013.292`.

[18] L. Gao, Y. Huang, H. Déjean, J.-L. Meunier, Q. Yan, Y. Fang, F. Kleber, E. Lang, Icdar 2019 competition on table detection and recognition (ctdar), in: 2019 International Conference on Document Analysis and Recognition (ICDAR), 2019, pp. 1510–1515. `doi:10.1109/ICDAR.2019.00243`.

[19] C. Tensmeyer, V. I. Morariu, B. Price, S. Cohen, T. Martinez, Deep splitting and merging for table structure decomposition, in: 2019 International

Conference on Document Analysis and Recognition (ICDAR), 2019, pp. 114–121. `doi:10.1109/ICDAR.2019.00027`.

[20] K. Itonori, Table structure recognition based on textblock arrangement and ruled line position, in: 1993 International Conference on Document Analysis and Recognition (ICDAR), 1993, pp. 765–768. `doi:10.1109/ICDAR.1993.395625`.

[21] B. Coüasnon, A. Lemaitre, Recognition of Tables and Forms, 2014, pp. 647–677. `doi:10.1007/978-0-85729-859-1_20`.

[22] T. Kieninger, Table structure recognition based on robust block segmentation, 1998, pp. 22–32.

[23] J. Zhang, J. Du, Y. Yang, Y.-Z. Song, S. Wei, L. Dai, A tree-structured decoder for image-to-markup generation, in: Proceedings of the 37th International Conference on Machine Learning, Vol. 119, 2020, pp. 11076–11085.

[24] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning to align and translate (2016).

[25] M.-T. Luong, H. Pham, C. D. Manning, Effective approaches to attention-based neural machine translation (2015).

[26] J. Gehring, M. Auli, D. Grangier, D. Yarats, Y. N. Dauphin, Convolutional sequence to sequence learning, in: Proceedings of the 34th International Conference on Machine Learning, Vol. 70, 2017, pp. 1243–1252.

[27] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, in: Advances in neural information processing systems, 2017, pp. 5998–6008.

[28] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, S. Zagoruyko, End-to-end object detection with transformers (2020).

[29] H. Hu, J. Gu, Z. Zhang, J. Dai, Y. Wei, Relation networks for object detection (2018).

[30] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, H. Lu, Dual attention network for scene segmentation, in: 2019 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 3141–3149. `doi:10.1109/CVPR.2019.00326`.

[31] H. Zhao, Y. Zhang, S. Liu, J. Shi, C. C. Loy, D. Lin, J. Jia, Psanet: Point-wise spatial attention network for scene parsing, in: V. Ferrari, M. Hebert, C. Sminchisescu, Y. Weiss (Eds.), Computer Vision – ECCV 2018, 2018, pp. 270–286.

[32] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel, Y. Bengio, Show, attend and tell: Neural image caption generation with visual attention (2016).

[33] M. Liao, P. Lyu, M. He, C. Yao, W. Wu, X. Bai, Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes, IEEE Transactions on Pattern Analysis and Machine Intelligence 43 (2) (2021) 532–548.

[34] X. Xiao, L. Wang, K. Ding, S. Xiang, C. Pan, Dense semantic embedding network for image captioning, Pattern Recognition 90 (2019) 285–296. `doi:https://doi.org/10.1016/j.patcog.2019.01.028`.

[35] M. Farazi, S. Khan, N. Barnes, Accuracy vs. complexity: A trade-off in visual question answering models, Pattern Recognition 120 (2021) 108106. `doi:https://doi.org/10.1016/j.patcog.2021.108106`.

[36] G. KV, A. Nambiar, K. S. Srinivas, A. Mittal, Linguistically-aware attention for reducing the semantic gap in vision-language tasks, Pattern Recognition 112 (2021) 107812. `doi:https://doi.org/10.1016/j.patcog.2020.107812`.

[37] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 770–778. `doi:10.1109/CVPR.2016.90`.

[38] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, S. Belongie, Feature pyramid networks for object detection, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 936–944. `doi:10.1109/CVPR.2017.106`.

[39] Y. Zhu, J. Du, Textmountain: Accurate scene text detection via instance segmentation, Pattern Recognition 110 (2021) 107336. `doi:https://doi.org/10.1016/j.patcog.2020.107336`.

[40] P. Lyu, M. Liao, C. Yao, W. Wu, X. Bai, Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes, in: V. Ferrari, M. Hebert, C. Sminchisescu, Y. Weiss (Eds.), Computer Vision – ECCV 2018, 2018, pp. 71–88.

[41] J. Chung, C. Gulcehre, K. Cho, Y. Bengio, Empirical evaluation of gated recurrent neural networks on sequence modeling, arXiv preprint arXiv:1412.3555 (2014).

[42] J. Zhang, J. Du, L. Dai, Track, attend, and parse (tap): An end-to-end framework for online handwritten mathematical expression recognition, IEEE Transactions on Multimedia 21 (1) (2019) 221–233. `doi:10.1109/TMM.2018.2844689`.

[43] B. Shi, M. Yang, X. Wang, P. Lyu, C. Yao, X. Bai, Aster: An attentional scene text recognizer with flexible rectification, IEEE Transactions on Pattern Analysis and Machine Intelligence 41 (9) (2019) 2035–2048. `doi:10.1109/TPAMI.2018.2848939`.

[44] Z. Qiao, Y. Zhou, D. Yang, Y. Zhou, W. Wang, Seed: Semantics enhanced encoder-decoder framework for scene text recognition, in: 2020 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 13525–13534. `doi:10.1109/CVPR42600.2020.01354`.

[45] M. Hurst, A constraint-based approach to table structure derivation, in:

Proceedings of the Seventh International Conference on Document Analysis and Recognition - Volume 2, 2003, p. 911.

[46] M. Pawlik, N. Augsten, Tree edit distance: Robust and memory-efficient, Information Systems 56 (2016) 157–173. `doi:https://doi.org/10.1016/j.is.2015.08.004`.

[47] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: 2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2009, pp. 248–255. `doi:10.1109/CVPR.2009.5206848`.

[48] M. D. Zeiler, Adadelta: An adaptive learning rate method (2012).

[49] I. Loshchilov, F. Hutter, Sgdr: Stochastic gradient descent with warm restarts (2017).

[50] A. Shigarov, A. Altaev, A. Mikhailov, V. Paramonov, E. Cherkashin, Tabbypdf: Web-based system for pdf table extraction, in: R. Damaševičius, G. Vasiljevienė (Eds.), Information and Software Technologies, 2018, pp. 257–269.

[51] Y. Liu, H. Chen, C. Shen, T. He, L. Jin, L. Wang, Abcnet: Real-time scene text spotting with adaptive bezier-curve network, in: 2020 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 9806–9815. `doi:10.1109/CVPR42600.2020.00983`.

[52] P. Lyu, M. Liao, C. Yao, W. Wu, X. Bai, Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes, in: Computer Vision – ECCV 2020, 2018, pp. 67–83.

[53] M. Liao, G. Pang, J. Huang, T. Hassner, X. Bai, Mask textspotter v3: Segmentation proposal network for robust scene text spotting, in: Computer Vision – ECCV 2020, 2020.