

# Deep Splitting and Merging for Table Structure Decomposition

Chris Tensmeyer, Vlad I. Morariu, Brian Price, Scott Cohen  
*Adobe Research*  
*San Jose, USA*  
 {tensmeyer,morariu,bprice,schoen}@adobe.com

Tony Martinez  
*Dept. of Computer Science*  
*Brigham Young University*  
*Provo, USA*  
 martinez@cs.byu.edu

**Abstract**—Given the large variety and complexity of tables, table structure extraction is a challenging task in automated document analysis systems. We present a pair of novel deep learning models (Split and Merge models) that given an input image, 1) predicts the basic table grid pattern and 2) predicts which grid elements should be merged to recover cells that span multiple rows or columns. We propose projection pooling as a novel component of the Split model and grid pooling as a novel part of the Merge model. While most Fully Convolutional Networks rely on local evidence, these unique pooling regions allow our models to take advantage of the global table structure. We achieve state-of-the-art performance on the public ICDAR 2013 Table Competition dataset of PDF documents. On a much larger private dataset which we used to train the models, we significantly outperform both a state-of-the-art deep model and a major commercial software system.

**Keywords**—Table Structure Extraction; Deep Learning; Convolutional Neural Network

## I. INTRODUCTION

Tables are commonly used to present structured data in documents. Automatically determining the structure of tables, i.e., the locations of cells and how they are organized into columns and rows, is a crucial part of document understanding and information extraction.

This work focuses on documents as images and in Portable Document Format (PDF) that contain tables with no structural annotations. Images, such as document scans, are composed of pixels that visually reproduce the document, but do not have any structural information. While PDFs are composed of discrete drawing commands for text and line art, table structure is not explicitly encoded in these commands. For example, a single text drawing command may be used to render part of a cell, an entire cell, parts of multiple cells, or multiple entire cells. Such ambiguities make it non-trivial to extract table structure from PDF drawing commands.

Several factors complicate table structure parsing (see [1] for an overview). Lined tables are generally easier than unlined tables [2] because they have explicit visual cues that delimit cells, columns, and rows. Other tables are partially lined or use shading to denote boundaries. Tables may have a complex and hierarchical header structure or may have no headers. Processing tables with multi-line cells involves determining if vertically adjacent text lines are part of the

same cell or not. Additionally, cells that span multiple columns and rows add to table complexity.

Many previous works that extract table structure are rule based systems that analyze PDF drawing commands [2]–[5], i.e., they don't directly analyze the document image. While heuristic methods have previously outperformed learning methods on PDFs, they cannot be directly applied to images. Furthermore, heuristic rules have many parameters that must be tuned as tables vary widely in structure, visual appearance, and complexity.

Recently, deep learning has been proposed to learn table structure directly from images [6]. The Deep DeSRT model [6] uses a neural network originally designed for semantic segmentation of natural scenes and thus primarily uses local information to classify pixels. However, local information is often insufficient to identify table cell boundaries because both cells and cell separators contain large regions of white space. Furthermore, cell boundaries in many unlined tables are visually indicated by globally consistent text justification (e.g., all cells in a column are left justified to the column boundary). To solve these issues, our proposed deep models perform pooling over large image regions to allow discriminative information to propagate across the whole image.

In this work, we propose a pair of deep learning models that predicts table structure from well-cropped table images. The *Split* model predicts the fine grid structure of the table, and the *Merge* model predicts where cells span multiple columns or rows. Together, we refer to the combined Split and Merge models as SPLERGE. The Split model uses novel pooling regions, termed *projection pooling*, where features are pooled along rows and columns of pixels to help learned features propagate across large images. Thus, pixels in the same row (column), but on opposite sides of the image, can exchange their local information on where the row (column) boundaries are. The output of this model is the grid structure of the table, ignoring spanning cells. The Merge model takes as input the original table image and predicts which elements of the previously predicted grid need to be merged to recover spanning cells. The novel contribution of the Merge model is *grid pooling*, where pooling regions are determined by the grid output of the Split model. This allows for local pixels features to easily propagate across large cells and



Figure 1. Overview of SPLERGE. First the Split model predicts the basic grid of the table, ignoring cells that span multiple rows or columns. Then the Merge model predicts which grid elements should be merged to recover spanning cells.

allow features to be exchanged among neighboring cells.

We achieve state-of-the-art performance on the ICDAR 2013 Table Competition dataset [2], obtaining an adjacency relationship F-measure of 95.15% compared to the best published result of 94.6% [2]. We also evaluate SPLERGE on a much larger and diverse private dataset of table images, significantly outperforming a state-of-the-art deep learning model and a major commercial software system.

## II. RELATED WORKS

While much work has focused on table localization, fewer works have attempted to extract table structure.

The ICDAR 2013 Table Competition [2] compared many systems for table structure recognition, including the current state-of-the-art method by Nurminen [3]. This method performs edge detection on the rendered PDF image (filtering out edges due to text) to find junctions and rectangular areas. Then it follows a series of heuristic steps to assign text to rows, identifying likely columns and headers, revising rows, merging columns, find spanning cells, etc. Shigarov et al. [5] later proposed a simpler set of heuristics based on bottom-up merging of PDF components. The TEXUS system [4] performs simple table structure parsing by considering each text chunk to be a cell. Compared to these methods, we do not require heuristics and can operate on both PDF and image input tables.

The DeepDeSRT [6] learning model formulates table structure extraction as two independent pixel labeling tasks to respectively segment row objects and segment column objects. It finetunes the FCN-2s model [7] that has been pretrained on the PASCAL VOC semantic segmentation challenge [8] for natural scenes. While it only requires an image as input, it does not perform as well as heuristic methods on PDFs. This may be because both rows (columns) and non-row (non-column) regions often contain large amounts of white background (e.g. space between cells in a row), which makes the pixel classes difficult to discriminate from local information. In contrast, SPLERGE pools information over large regions to use more global information. Clinchat et al. [9] compare two learning models on scanned images of seventeenth century handwritten tables. They focus on this narrow problem and use template matching to recover columns and convert the images into “PDF-like” documents. To recover table rows, they use either a Conditional Random

Field (CRF) or their proposed Edge Convolutional Networks, an extension of Graph Convolutional Networks [10].

## III. SPLERGE

SPLERGE, our proposed table structure extraction method, is composed of two deep learning models that perform split and merge operations in sequence (see Fig. 1). The Split model takes an input image of a well cropped table and produces the grid structure of the table in the form of row and column separators that span the entire image. Because some tables contain spanning cells, we apply the Merge model to the grid output of the Split model to merge together adjacent grid elements to recover the spanning cells.

### A. Split Model

The Split model takes as input an image of any dimension  $H \times W$  and produces two 1-D output signals:  $r \in [0, 1]^H$  and  $c \in [0, 1]^W$ . The output signals  $r$  and  $c$  represent the probability that each row (column) of pixels is part of a logical table row (column) separator region. The Split model is composed of 3 sub-networks:

- 1) Shared Fully Convolutional Network (SFCN)
- 2) Row Projection Network (RPN)
- 3) Column Projection Network (CPN)

The SFCN computes local image features that are used by both the RPN and CPN. The RPN and CPN then take those local features and further process them to predict row and column separators ( $r$  and  $c$  respectively).

The SFCN is composed of 3 convolution layers with  $7 \times 7$  kernels, with the last layer performing a dilated convolution [11] with a dilation factor of 2. Each layer produces 18 feature maps and uses the ReLU activation function.

Dilated convolutions, like pooling, increase the receptive field of the network, but unlike pooling, they preserve the spatial resolution of the input. Preserving the input spatial resolution is important in table structure extraction because many column and row separators are only a few pixels wide. In [6], better results were obtained when resizing the initial input to make the separator regions bigger. It is also crucial to have a large receptive field because determining the locations of row and column separators may require global context. For example, text that is consistently left-justified to the same position is indicative of a column separator.

The output of the SFCN is fed as input to both the RPN and CPN. The output of the RPN is  $r$ , the probability

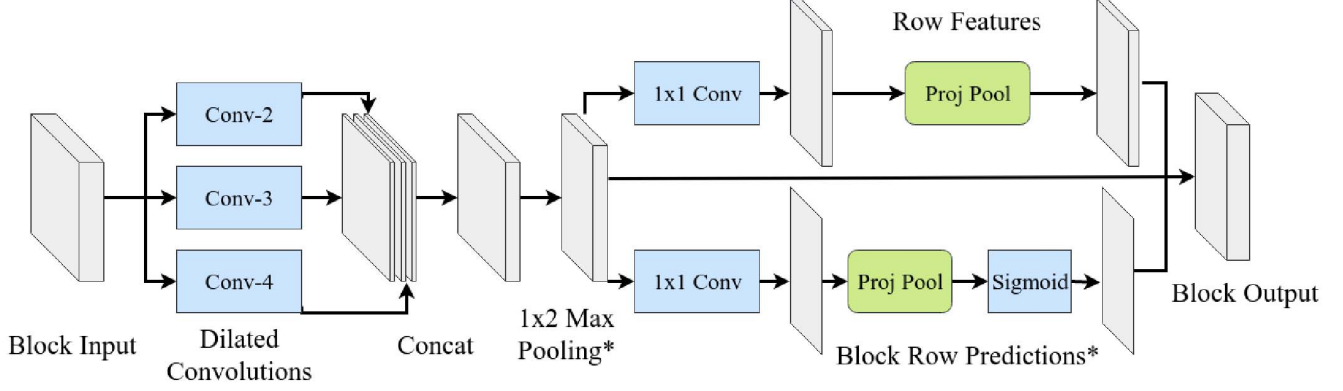


Figure 2. A single block of the Row Projection Network (RPN), which is a sub-network of the Split model. Components marked with \* indicate that they are not included in all 5 blocks that compose the RPN. Blocks in the Column Projection Network (CPN) use 2x1 Max Pooling.

that each row of pixels is part of a row separator region. Similarly, the output of the CPN is  $c$ . Because the RPN and CPN have identical structure, except for whether projection and pooling operations are over rows or columns of pixels, we focus our discussion only on the RPN.

The RPN in this work is composed of 5 blocks chained together, though any number of blocks can be used. Empirically, using more than 5 blocks did not improve our results, and a similar process was used to determine the other specific architectural choices. To simplify the discussion and notation, we use actual hyperparameter values for our experiments, but other reasonable values could be used. Changing the hyperparameters in a reasonable range did not seem to significantly impact the results in our informal experiments. Fig. 2 shows the operations performed by a single block. First, the input is (in parallel) fed to 3 convolution layers, which have dilation factors of 2/3/4, to produce 6 feature maps each. The output of each dilated convolution is concatenated to obtain 18 feature maps. Using multiple dilation factors allows the RPN to learn multi-scale features and increase its receptive field while still examining more local evidence.

Next, the RPN performs 1x2 max pooling (CPN does 2x1 max pooling). This decreases the width of the feature maps, but maintains the height, so that the output signal  $r$  is of size  $H$ . Only the first three blocks perform max pooling to make sure that the width is not downsampled too much.

Afterwards, the RPN computes row features (top branch of Fig. 2) through a 1x1 convolution operation followed by *projection pooling* (Fig. 3). Projection pooling is inspired by the projection profile operation used to find gaps of white-space in classical layout analysis. With projection pooling, we maintain the spatial size of the input (instead of collapsing to 1D as in projection profiles), and simply replace each

value in the input with its row average. Specifically,

$$\hat{F}_{ij} = \frac{1}{W} \sum_{j'=1}^W F_{ij'} \quad (1)$$

where  $i$  indexes over rows in feature map  $F$ , and  $1 \leq j \leq W$  indexes over columns of  $F$ . We call  $\hat{F}$  the row projection pooling of  $F$ , and apply this operation independently on each feature map, as is typical of pooling operations. Pooling in this fashion allows information to propagate across the entire width of the image, which may be more than 1000 pixels. These row features are concatenated to the output of the max pooling operation, so that each pixel has both local and row-global features. Note that the CPN performs *column* projection pooling, which similarly is

$$\hat{F}_{ij} = \frac{1}{H} \sum_{i'=1}^H F_{i'j} \quad (2)$$

The bottom branch of Fig. 2 shows how blocks produce row predictions, though not every block does so. A 1x1 convolution produces a single output map, over which we perform projection pooling. We then apply a sigmoid function to produce probabilities. Since each row of pixels contains a single unique probability, we can take a vertical slice to obtain a 1D signal of probabilities,  $r^n$ , where  $n$  indicates the block index. To make the intermediate prediction  $r^n$  available to block  $n+1$ , we also concatenate the probabilities (in 2D) to the output of the block.

In our implementation, only the last 3 blocks produce outputs, i.e.  $r^3, r^4, r^5$ . During training, we apply a loss to all three predictions, but after training, we only use the last prediction,  $r^5$ , for inference. This iterative prediction procedure allows the model to make a prediction and then refine that prediction. Such techniques have been successfully applied in previous work on structured keypoint detection tasks in natural scenes (e.g., human pose [12]).

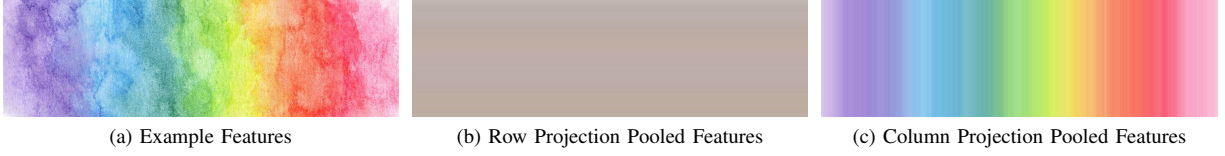


Figure 3. Illustration of projection pooling operation. Every pixel in (a) is replaced by its row average (b) or its column average (c).

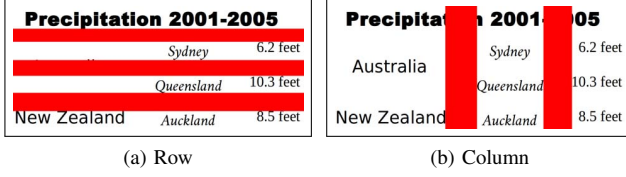


Figure 4. Example 1D GT row and column signals for training the Split model. The 1D signals are replicated to 2D and superimposed on the input image for viewing purposes.

1) *Training*: The SFCN, RPN, and CPN sub-networks are jointly trained in a typical supervised fashion on table images at 150 DPI. We assume that the images are cropped to only include the table cells and exclude table titles, captions, and footnotes that are not in the cell region.

Each table has annotated ground truth (GT) 1D signals  $r^*$  and  $c^*$ . The GT is designed to maximize the size of the separator regions without intersecting any non-spanning cell content, as shown in Fig. 4. This is contrary to the traditional notion of cell separators, which for many tables are thin lines that are only a few pixels thick. Predicting small regions is more difficult than large regions, and in the case of unlined tables, the exact location of the cell separator is ill defined. The GT separator regions may intersect cell content for cells that span multiple rows or columns. The goal of the Split model is to recover the basic grid of the table, and spanning cells are handled by the Merge model.

The loss function is the average element-wise binary cross entropy between the block predictions and the GT signals:

$$L(r, r^*) = \frac{1}{H} \sum_{i=1}^H r_i^* \log r_i + (1 - r_i^*) \log(1 - r_i) \quad (3)$$

To prevent overfitting, we modify equation 3 to clip the per-element loss to 0 when  $|r_i^* - r_i| < 0.1$ . The total loss is a weighted sum of individual output losses:

$$L_{tot} = L(r^5, r^*) + \lambda_4 L(r^4, r^*) + \lambda_3 L(r^3, r^*) + L(c^5, c^*) + \lambda_4 L(c^4, c^*) + \lambda_3 L(c^3, c^*) \quad (4)$$

where we set  $\lambda_4 = 0.25$  and  $\lambda_3 = 0.1$ .

We train the model from random initialization with the ADAM optimizer [13] for approximately  $10^6$  weight updates. We use a batch size of 1 because the table images have different spatial sizes. We use an initial learning rate of 0.00075 and decay it by a factor of 0.75 every 80K updates.

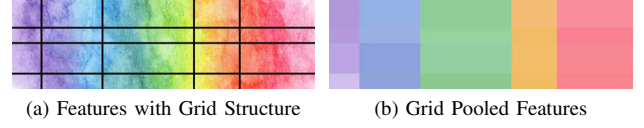


Figure 5. Illustration of grid pooling using the black grid structure visualized on top of the features in (a). Values are averaged within each grid cell to produce (b).

2) *Inference*: Once we have predicted  $r$ , we need to infer at which pixel positions the row separators occur. For simplicity, the discussion focuses on  $r$ , but the same procedure applies to  $c$  to obtain column separators. To do so, we partition the image into rows and row separator regions by performing a graphcut segmentation [14] over  $r$ . Then we choose the row pixel positions corresponding to the midpoints of each inferred separator region.

To create the graph for segmenting  $r$ , we have  $H$  nodes arranged in a linear chain, where each node is connected to its two neighbors (except the two nodes at either end). The neighbor edge weights are uniformly set to  $\lambda_{gc} = 0.75$ . Node  $i$  is connected to the source node with an edge weight of  $r_i$  and to the sink node with an edge weight of  $1 - r_i$ .

## B. Merge Model

The Merge model uses the input image and the outputs of the Split model to predict which grid elements need to be merged to recover cells that span multiple rows or columns. The input tensor is a concatenation of the table image, the output row/col probabilities ( $r, c$ ), the inferred row/col regions, and the predicted grid structure. The predicted probabilities,  $r$  and  $c$ , are transformed to a 2D image by stacking (i.e.,  $[r, r, \dots, r]$ ). The inferred row/col regions are rendered as binary mask (similar to the red regions in Fig. 4). The predicted grid structure is rendered as a binary mask where the midpoint of each row and column separator region is rendered as a 7-pixel wide line. Additionally, the grid structure is used to determine the pooling regions of the model.

If the grid structure is composed of  $M$  rows and  $N$  columns, then the model outputs two matrices:

- 1)  $D$  - probs. of up-down merges (size  $(M - 1) \times N$ )
- 2)  $R$  - probs. of left-right merges (size  $M \times (N - 1)$ )

$D_{ij}$  is the probability of merging cell  $(i, j)$  with cell  $(i + 1, j)$ , and  $R_{ij}$  is the probability of merging cells  $(i, j)$  and  $(i, j + 1)$ .  $D$  is not of size  $M \times N$  because there are only  $M - 1$  pairs of up-down merges in any column. In our



<b>Precipitation 2001-2005</b>		0	0	0	1	1
	Sydney	6.2 feet	1	0	0	0
Australia	Queensland	10.3 feet	0	0	0	0
New Zealand	Auckland	8.5 feet	0	0	0	0

(a) Image with Grid Structure

(b) GT  $D$

(c) GT  $R$

Figure 6. Example GT for training the Merge model for the structure predicted by the Split model.

formulation, all of these probabilities are independent, i.e., a single grid element may be merged in multiple directions.

The Merge model architecture is similar to the Split model. There is a set of 4 shared convolution layers (no dilation), where 2x2 average pooling occurs after layers 2 and 4. Afterwards, the model has 4 branches, where each branch predicts an  $M \times N$  matrix of probabilities that cells should merge in a particular direction, i.e. up, down, left, or right. Call these matrices  $u, d, l, r$ . While our independence assumption suggests that we multiply the two individual probabilities together in Eqs. 5, 6, this would introduce optimization difficulties when both probabilities are near 0, so we compute  $D$  and  $R$  as

$$D_{ij} = \frac{1}{2}u_{i+1,j}d_{ij} + \frac{1}{4}(u_{i+1,j} + d_{ij}) \quad (5)$$

$$R_{ij} = \frac{1}{2}l_{i,j+1}r_{ij} + \frac{1}{4}(l_{i,j+1} + r_{ij}) \quad (6)$$

Intuitively, we only predict that a pair of cells should be merged if there is consistency between the individual branch outputs.

Each branch is composed of 3 blocks that are similar to the Split model block shown in Fig. 2. The differences are that the parallel convolution layers use 1/2/3 dilation factors, no max pooling is performed, and projection pooling is replaced with *grid pooling* (Fig. 5). In grid pooling, every pixel location replaces its value with the average of all pixels within its grid element:

$$\hat{F}_{ij} = \frac{1}{|\Omega(i,j)|} \sum_{i',j' \in \Omega(i,j)} F_{i'j'} \quad (7)$$

where  $\Omega(i,j)$  is the set of coordinates for all pixels that share the same grid element as  $(i,j)$ . After grid pooling, all pixels inside the same grid element share the same value, which allows information to propagate within each cell. Subsequent convolutions allow information to propagate between neighboring cells. To produce the  $u, d, l$ , or  $r$  matrix for a given branch, we average the predicted per-pixel probabilities in each grid element and arrange them in an  $M \times N$  matrix. Like the Split model, the Merge model also performs iterative output refinement, where blocks 2 and 3 produce output predictions.

1) *Training*: Because the Split and Merge models are intended to be used in sequence, we train the Merge model using the grid structures produced by the Split model. To construct the GT  $D$  and  $R$  matrices (see Fig. 6), we

- 1) Iterate over all spanning cells in the table
- 2) Identify which grid elements intersect the GT bounding box of the spanning cell
- 3) Set the probability of merge for each of these cells to 1 for the appropriate directions

As in the Split model, the loss function for each output is the average (clipped) element-wise binary cross entropy (Eq. 4). The total loss is then

$$L_{tot} = L(D^3, D^*) + \lambda_2 L(D^2, D^*) + L(R^3, R^*) + \lambda_2 L(R^2, R^*) \quad (8)$$

Because spanning cells only occur in 15% of tables in the private dataset used to train our models, we subsample this set so that 50% of the training set for the Merge model has at least one pair of cells that require merging. The training hyperparameters are similar to those of the Split model.

Inference is performed by thresholding  $D$  and  $R$  at 0.5 probability, and merging the indicated cells. The network prediction has no constraint that the resulting merges produce only rectangular cells, so in post-processing additional merges are added to ensure that the resulting table structure has only rectangular cells. For example, 3 grid elements merged together to form an L-shaped cell would be merged with a 4<sup>th</sup> element to create a cell that spans 2 rows and 2 columns.

#### IV. EXPERIMENTAL RESULTS

We tested SPLERGE on two datasets, including the ICDAR 2013 Table Competition benchmark dataset [2], where we achieve state-of-the-art results. This dataset is composed of 156 tables with unambiguous structure that were extracted from PDFs downloaded from government websites. The entire set of PDF tables is considered test data.

The second dataset is a private collection of 93,000 tables taken from web-scraped PDFs that have been manually annotated with structure information. These tables exhibit a larger variety of complexity, visual appearance, and structure than the ICDAR 2013 dataset. For example, there are only 2 completely unlined tables in ICDAR 2013, but the private collection is 15% unlined tables. Unlike the ICDAR 2013 dataset, the private collection includes unlined tables with multi-line cells. It also contains tables with poorly aligned columns, cells with complex structure, raster images inside of cells, and complex headers (such as nested headers and headers in the middle of the table). Some example tables are shown in Fig. 7.

We randomly split this dataset into 83K/5K/5K for training, validation, and testing. Because tables from the same PDF exhibit some similarity, we ensured that the set of PDFs used to construct each data split are disjoint.

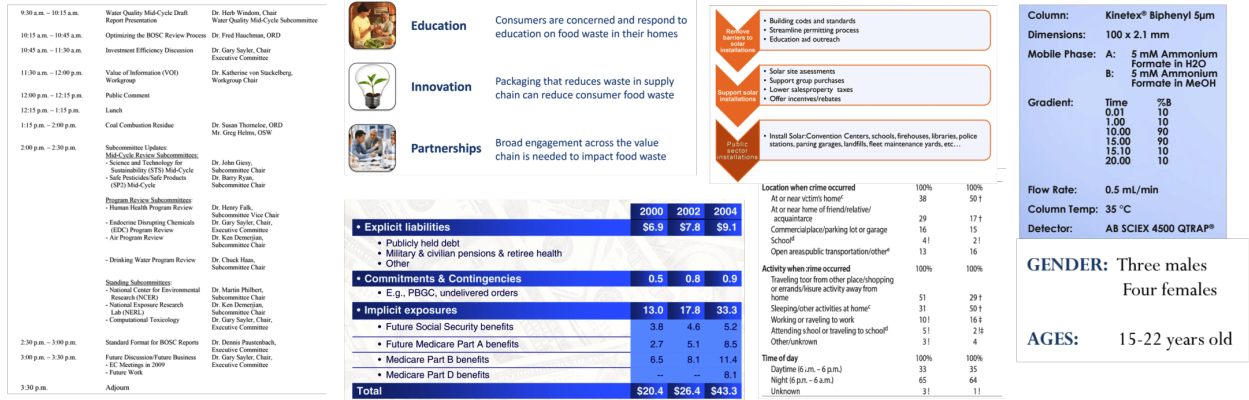


Figure 7. Example table images from the private collection. Some tables in the private collection have complex cells, such as multi-line text (in white space tables), nested tables, vector graphics, and raster images.

Table I  
ICDAR 2013 TABLE DATASET RESULTS.

Category	Method	F-measure	Recall	Precision	PDF Input	Image Input
Proposed Models	Split	86.79	86.64	86.93		✓
	Split-PDF	91.63	91.26	92.00	✓	✓
	SPLERGE-PDF	90.89	90.44	91.36	✓	✓
	SPLERGE-PDF + Heuristics	91.95	91.46	92.45	✓	✓
	Split + Heuristics	93.00	92.24	93.78	✓	✓
	Split-PDF + Heuristics	<b>95.26</b>	<b>94.64</b>	<b>95.89</b>	✓	✓
Previously Reported Results	Nurminen [3]	94.60	94.09	95.12	✓	✓
	TEXUS [4]	82.59	84.23	81.02	✓	
	Shigarov [5] $C_1$	91.50	91.21	91.80	✓	
	Shigarov [5] $C_2$	93.64	92.33	94.99	✓	
	DeepDeSRT [6]	91.44†	87.36†	95.93†		✓

† result not directly comparable due to evaluation on a random subset of 34 tables.

We also experimented with giving the models additional input channels, though the model does not require this. If the input table images are rendered from PDFs, then by examining drawing command metadata we can render additional images that contain only text, only vector art (e.g., lines), or only raster images. These additional images can be concatenated with the rendered PDF and slightly improve network performance.

#### A. ICDAR 2013

Our results on the ICDAR 2013 dataset are from models trained on the private collection. We attempted to validate that our improved performance results from a better deep model and not just from a larger training set. We did so by reimplementing the DeepDeSRT model [6] and training on the same data as our proposed model.

Table I shows the results on the ICDAR 2013 Table Competition dataset (task 2). Our methods with the *-PDF* suffix indicate that the additional PDF rendered input channels were used. The evaluation metric for this dataset is the F-measure over detected *adjacency relationships* [15]. Roughly speaking, this measures the percentage of correctly detected pairs of adjacent cells, where both cells are segmented correctly and identified as neighbors.

2005				2010			
Number	Margin of error (±)	Percent	Margin of error (±)	Number	Margin of error (±)	Percent	Margin of error (±)
291,099	*****	100.0	(X)	303,858	*****	100.0	(X)

(a) Blank Cells next to Headers

THRESHOLD FOR RELEASES			Inadequate housing uni	
to air kg/year	to water kg/year	to land kg/year	(%)	Unadjusted OR
100 million	-	-	(4.7)	Ref.
100	-	-	(5.9)	1.1
100 000	-	-		

(b) Separators Through Text

ERDF	Convergence	EURbn	Total	EURbn
3.2			5.5	
0			0.21	
13.4			22.2	
1.9			3.0	
11.2			19.8	
2.4			3.9	
3.4			5.7	

(c) Join Columns

(d) Split Columns

Figure 8. Example heuristic post-processing on ICDAR2013 table segments. Removed separators are shown in red, added separators in green, and unmodified predicted separators in blue.

For this dataset, the Merge model failed to provide adequate post-processing for the output of the Split model. After performing the predicted merges, our post processing combines additional cells to prevent cells from having an

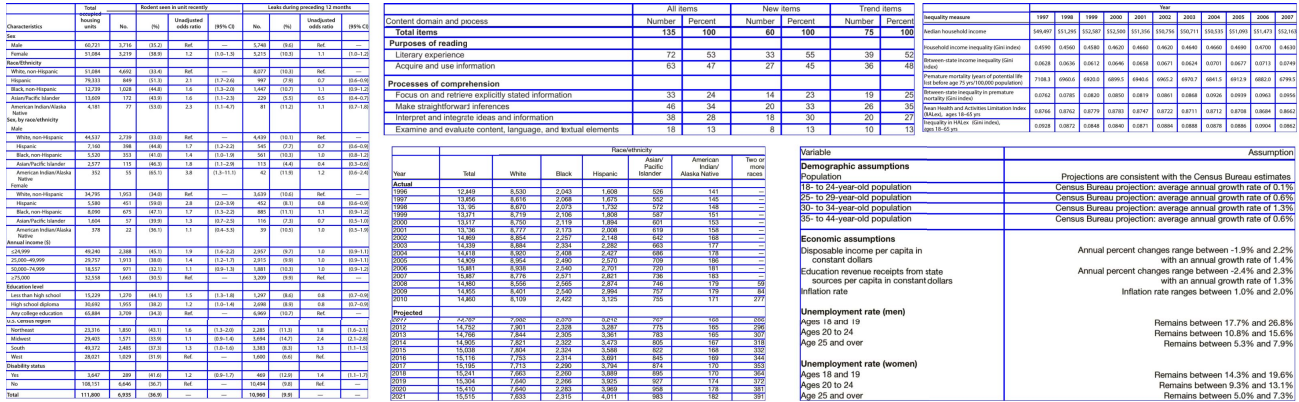


Figure 9. Example predictions on challenging (mostly unlined) ICDAR 2013 tables, including one failure case.

Table II  
PRIVATE COLLECTION RESULTS.

Category	Method	Table Accuracy	F-measure	Recall	Precision	PDF Input	Image Input
Proposed Models	Split	80.79	95.46	<b>95.87</b>	95.06		✓
	SPLERGE	83.09	95.78	95.50	96.07		✓
	Split-PDF	81.07	95.37	95.64	95.10	✓	✓
	SPLERGE-PDF	<b>85.65</b>	<b>95.92</b>	95.59	<b>96.27</b>	✓	✓
Our Implementation	DeepDeSRT	51.76	76.85	75.32	78.44		✓
	DeepDeSRT-PDF	53.48	78.44	77.03	79.91	✓	✓
Commercial	Acrobat	61.11	82.62	82.66	82.59	✓	

L-shape in the final output. In several of the large header regions in the ICDAR 2013 dataset, large groups of individual cells were merged into a single cell due to a few incorrect pairwise merge predictions that created L-shapes. Instead of further refining the heuristic to prevent L-shapes, we instead implemented some simple heuristics that could replace the Merge model. These heuristics include

- Merge cells where predicted separators pass through text.
- Merge adjacent columns when the vast majority of paired cells (after row 3) are either both blank or only one cell per pair is non-blank. This merges a content column with a (mostly) blank column.
- In the first row (likely header row), merge non-blanks cells with adjacent blank cells.
- Split columns that have a consistent whitespace gap between vertically aligned text.

Some example tables that are fixed by the heuristics are shown in Fig. 8. While the Split model performs well in identifying the target grid, it sometimes makes easily correctable mistakes and cannot by itself handle spanning cells. When combined with simple heuristics to handle these cases, it achieves an F-measure of 95.26% compared to the previous best result of 94.60%. The Merge model fails to generalize from the private collection to the ICDAR 2013 dataset, but as shown in Table II, it does improve performance on the private collection. Figure 9 shows some example predictions by *Split-PDF + Heuristics* for unlined

tables, which are harder to recognize than lined tables.

There is a large performance difference from when PDF information (text, path, image channels) is given as input to the Split model and when it is not. Because the difference is not as large on the private collection (Table II), we conclude that the usefulness of PDF input channels depends on the dataset. The ICDAR tables are primarily lined, have larger headers, and may have distinct visual appearances compared to the training dataset. Thus, the additional PDF information may help more in the unfamiliar domain because the text and path elements are explicit inputs instead of needing to be visually inferred by the model.

We reimplemented the DeepDeSRT table structure model [6] and trained it on the same private data as our proposed model. However, we were unable to obtain reasonable performance, even after exploring a variety of values for post processing thresholds and training hyperparameters (values not specified in [6]). In [6], they report an FM of 91.44%, but over a random subset of 34 tables, so a direct comparison cannot be done. We felt that this gap indicated that we could not faithfully reproduce their model for a fair comparison. However, the training set we used is very different and there are significant differences between the private collection and the ICDAR 2013 dataset. These reasons could explain the performance gap, but the gap is large enough that we are unsure that our implementation is a faithful reproduction of Deep DeSRT, so we omit exact performance numbers to avoid direct comparison.

## B. Private Collection

On this dataset, we evaluate methods using precision and recall over correctly detected cells. We also report accuracy as the percentage of tables with perfect precision and recall. A predicted bounding box (BB) is a correct prediction if it entirely contains only a single ground truth cell content BB. In particular, predicted BBs that intersect multiple GT BBs or that do not completely contain any GT BB are marked as false positives. Unmatched GT BBs are marked as false negatives. Because blank cells were not manually annotated, we exclude predicted BBs that do not intersect any GT BBs. This way, methods are not penalized if they correctly predict unmarked blank cells.

Table II shows the results on the test split of 5000 tables. Following [2], reported precision and recall are computed per-table and then averaged. We were unable to find any official implementations of prior works, so for comparison, we use the commercial software system Acrobat Pro DC and our reimplementations of the DeepDeSRT model [6].

All variants of our proposed models significantly outperform the two baselines on all metrics. We also see that the Merge model significantly improves table accuracy, as the Split model alone cannot handle tables with spanning cells. For tables that require merges, the average number of needed merges per-table is much less for the private collection than for ICDAR 2013, resulting in fewer L-shaped predictions. We also observe that using PDF information as input does bring an improvement, but more marginally than with the ICDAR 2013 data. This may reflect the fact that domain differences between ICDAR 2013 and the private collection are more apparent in the rendered PDF, but the differences are lessened when examining just the text or path image channels. This suggests that this method may be effective

## V. CONCLUSION

We have proposed a new method for table structure extraction, SPLERGE. It is composed of a pair of deep learning models that together split a table image into the basic grid of cells, and then merge cells together to recover cells that span multiple rows and columns. The key insight of the models is to pool information over large regions of the table image such as entire rows/columns of pixels or previously predicted cell regions. When evaluating the split model on the ICDAR 2013 Table Competition dataset, we achieve state-of-the-art performance. We also show that PDF information, such as if page elements are text/path/image, can be encoded as input to the deep network and improve performance. However, if such information is not available (e.g., scanned documents), the model can use only the grayscale image as input. Lastly, we have shown that the Merge model is effective on a private collection of tables extracted from the web.

for scanned documents, where only the captured image is available for input, if appropriate training data is provided.

## REFERENCES

- [1] J. Hu, R. Kashi, D. Lopresti, G. Nagy, and G. Wilfong, "Why table ground-truthing is hard," in *International Conference on Document Analysis and Recognition*. IEEE, 2001, pp. 129–133.
- [2] M. Göbel, T. Hassan, E. Oro, and G. Orsi, "ICDAR 2013 table competition," in *International Conference on Document Analysis and Recognition*. IEEE, 2013, pp. 1449–1453.
- [3] A. Nurminen, "Algorithmic extraction of data in tables in PDF documents," Master's thesis, Tampere University of Technology, 2013.
- [4] R. Rastan, H.-Y. Paik, and J. Shepherd, "TEXUS: A task-based approach for table extraction and understanding," in *Symposium on Document Engineering*. ACM, 2015, pp. 25–34.
- [5] A. Shigarov, A. Mikhailov, and A. Altaev, "Configurable table structure recognition in untagged PDF documents," in *Symposium on Document Engineering*. ACM, 2016, pp. 119–122.
- [6] S. Schreiber, S. Agne, I. Wolf, A. Dengel, and S. Ahmed, "DeepDeSRT: Deep learning for detection and structure recognition of tables in document images," in *International Conference on Document Analysis and Recognition*, vol. 1. IEEE, 2017, pp. 1162–1167.
- [7] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [8] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The Pascal visual object classes (VOC) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [9] S. Clinchant, H. Déjean, J.-L. Meunier, E. M. Lang, and F. Kleber, "Comparing machine learning approaches for table recognition in historical register books," in *International Workshop on Document Analysis Systems*. IEEE, 2018, pp. 133–138.
- [10] Q. Xu, Q. Wang, C. Xu, and L. Qu, "Collective vertex classification using recursive neural network," *arXiv preprint arXiv:1701.06751*, 2017.
- [11] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [12] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," in *Computer Vision and Pattern Recognition*, 2016, pp. 4724–4732.
- [13] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [14] D. M. Greig, B. T. Porteous, and A. H. Seheult, "Exact maximum a posteriori estimation for binary images," *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 271–279, 1989.
- [15] M. Göbel, T. Hassan, E. Oro, and G. Orsi, "A methodology for evaluating algorithms for table understanding in PDF documents," in *Symposium on Document Engineering*. ACM, 2012, pp. 45–48.