# Visual Understanding of Complex Table Structures from Document Images

Sachin Raja
IIIT-Hyderabad
sachin.raja@research.iiit.ac.in

Ajoy Mondal
IIIT-Hyderabad
ajoy.mondal@iiit.ac.in

Jawahar C V
IIIT-Hyderabad
jawahar@iiit.ac.in

## 1. Organization

We organize our supplementary paper in the following manner:

- Sections 2 through 4 describe the implementation details along with specific pre-processing and post-processing steps that would help to reproduce our results.

- Section 5 discuss the comparative analysis on the performance of the proposed method over existing techniques on various benchmark datasets for both physical and logical table structure recognition tasks.

- Section 6 shows some qualitative results obtained by the proposed method.

## 2. Implementation Details

Both networks - TOD-Net and TSR-Net have been trained and evaluated with the table images scaled to a fixed size of $1536 \times 1536$ while maintaining the original aspect ratio as the input. We use NVIDIA TITAN X GPU with 12 GB memory for our experiments and a batch-size of 1. We use a dilated convolution with the filter size of $2 \times 2$ and a dilation parameter of 2. Also, we use the ResNet-101 backbone that is pre-trained on MS-COCO dataset. To compute region proposals, we use 0.5, 1 and 2 as the anchor scale and anchor box sizes of 8, 16, 32, 64 and 128. Further, for the generation of the row/column adjacency matrices, we use 2400 as the maximum number of vertices keeping in mind dense tables. Since every input table may contain hundreds of table cells, training can be a time-consuming process. To achieve faster training, we employ a 2-stage training process. During both stages, we use 0.001 as the learning rate, 0.9 as the momentum, and 0.0001 as the weight decay regularization.

To achieve faster training, we employ a 2-stage training process. In the first step, we use 2014 anchors and 512 ROIs. With this setting, the model can learn high and low-level features but resulted in many false negatives. To combat this, network is trained with 3072 anchors and 2048 ROIs. This significantly reduces the number of false negatives. For the first step, we train a total of 30 epochs, for the first 8, we train all FPN and subsequent layers, for the next 15, we train FPN + last 4 layers of ResNet-101 and for the last 7 epochs, we train all the layers of the model. For the second step, we train a total of 10 epochs, for the first 3, we train all FPN and subsequent layers, for the next 4, we train FPN + last 4 layers of ResNet-101 and for the last 3 epochs, we train all the layers of the model. During both stages, we use 0.001 as the learning rate, 0.9 as the momentum, and 0.0001 as the weight decay regularisation.

## 3. Dataset Pre-processing

We use FinTabNet [41] dataset to train TOD-Net for cell, row, and column detection. Since FinTabNet has bounding boxes wrapped around the cell's content, we pre-process the ground truth to obtain cell level coordinates. To do that, we first find out the minimum and maximum x and y coordinates for every row and column based on the content boxes. We then split the y-coordinate difference between adjacent rows by half and extend the maximum y end coordinate of the previous row and reduce the current row's minimum y-coordinate accordingly. Similarly, we split the x-coordinate difference between adjacent columns by half and extend the maximum x end coordinate of the previous column and reduce the minimum x-coordinate of the current column accordingly. The resulting dataset follows all the constraints that we model in the TOD-Net. Please refer figures 1 and 2 for a step-by-step visual illustration of our algorithm.

## 4. Post-processing after Detection of Rows and Columns

After obtaining row and column bounding boxes, we employ additional post-processing steps to get cell bounding boxes along with the tabular structure in terms of row and column spans. The first step is to obtain cells by finding intersections of row and columns. If a row and column overlap with greater than 75% area, it is most likely a cell that spans multiple rows and/or columns. Further, rectilinear adjacencies output also helps in locating multi-row and multi-column spanning cells. Next step includes the removal of objects that have a high overlap with oth-

ers. Removal of some objects might result in some empty regions that are then filled-up based on the average height of rows and columns' average width. If the gap between sorted adjacent objects is greater than the average height of the column's row or average width, we add extra row(s) or column(s) in that empty space. In the other case, we split the gap between the adjacent objects to remove any gap or any small overlapping region that was missed from the first step. Next, we use Tesseract's [28] HOCR (open standard of data representation for formatted text obtained from optical character recognition) output to finetune cells' bounding boxes. This is done by first identifying words belonging to every cell and updating start and end coordinates in a way such that the content gets completely encapsulated within the cell boxes. This particularly helps us in making sure that as much content of a cell gets retained as possible thereby reducing false negatives with respect to content. After doing this, we align the y-coordinates of detected rows and the x-coordinates of the detected columns based on the output of the previous step. Rectilinear ajdacencies help in building associations between cells. Adjacent cells that have high overlap areas and positive adjacencies are said to belong to the same row or column. Overlapping area percentage helps in resolving conflicting cases for some pair of cells where the rectilinear adjacencies output is not completely accurate.

## 5. Quantitative Results

Table 1 presents physical table structure recognition using the proposed method and comparison with the state-of-the-art techniques on UNLV [27] dataset. Tables 2 and 3 illustrate logical table structure recognition using the proposed method and comparison with state-of-the-art techniques on PubTabNet [42] and TableBank [17] datasets, respectively. Table 4 further presents results across varying IoU thresholds.

## 6. Qualitative Results

Please refer figures 3 through 12 for visual results on sample table images from scientific and business domains tested on standard evaluation datasets. Please observe that for every table image sample, we provide 6 images — (i) original test table image, (ii) Tesseract's [28] HOCR output, (iii) cells detection output using TOD-Net[†], (iv) cells detection output using TOD-Net[‡], (v) structure showing rows in alternating green, red and blue colors, and (vi) structure showing columns in alternating green, red and blue colors.

Figure 3 clearly shows the performance of our method on a sample ICDAR-2013 [8] dataset in the presence of horizontal and vertical separators where the table image contains many empty cells. What is interesting to note in this example is that the Tesseract output gets messed up as part

of the horizontal and vertical separators get falsely recognised as legitimate characters. Figure 4 shows the robustness of our method on a sample image from ICDAR-2019 [7] archival tables dataset. Furthermore, the superiority of our method is bolstered for both dense and sparse tables from Figures 5, 9 and 11. A sample image from TableBank [17] dataset, Figure 6 also demonstrates qualitative performance for table containing multi-line cells.

**Effect of Trainable Loss Weights** In our experiments, we used dynamic loss weights for each of the three loss weights based on the region of interest's visual characteristics. Figure 13 shows that with the same weight initialization, the model with dynamic loss weights converges faster and slightly better by 0.4%.

| Method | Training Dataset | Test on UNLV | | | | | |
|---|---|---|---|---|---|---|---|
| | | SEC | | | NEC | | |
| | | P↑ | R↑ | F1↑ | P↑ | R↑ | F1↑ |
| DGCNN[†] [22,24] | FinTabNet | 0.849 | 0.833 | 0.841 | 0.827 | 0.810 | 0.818 |
| DGCNN[‡] [22,24] | FinTabNet | 0.873 | 0.867 | 0.870 | 0.858 | 0.851 | 0.854 |
| SPLIT [30] | Private | 0.795 | 0.776 | 0.785 | 0.748 | 0.715 | 0.731 |
| TS-Net [24] | FinTabNet | 0.852 | 0.834 | 0.843 | 0.833 | 0.819 | 0.826 |
| Ours[†] | FinTabNet | 0.867 | 0.849 | 0.858 | 0.852 | 0.838 | 0.845 |
| Ours[‡] | FinTabNet | 0.890 | 0.884 | 0.887 | 0.882 | 0.876 | 0.879 |

Table 1. Shows the performances of the proposed networks and its comparison against the existing techniques on physical table structure recognition in ICDAR-2019 dataset. DGCNN[†] indicates TOD-Net[†]+DGCNN+PP, DGCNN[‡] indicates TOD-Net[‡]+DGCNN+PP TS-Net indicates TabStruct-Net, Ours[†] indicates TOD-Net[†]+TSR+PP, and Ours[‡] indicates TOD-Net[‡]+TSR+PP. IoU threshold is set to 0.6. **P:** indicates precision, **R:** indicates recall, and **F1:** indicates F1 score.

| Method | Training | | Average Over Test Set (PubTabNet [40]) |
|---|---|---|---|
| | Dataset | #Image | TEDS↑ |
| EDD [40] | PubTabNet [40] | 420K | 88.30 |
| TabStruct-Net [24] | SciTSR [4] | 012K | 90.10 |
| TOD-Net[†]+TSR+PP | FinTabNet [39] | 91K | 89.50 |
| TOD-Net[‡]+TSR+PP | FinTabNet [39] | 91K | 90.70 |

Table 2. Shows the performances of the proposed networks and its comparison against the existing techniques on logical table structure recognition in PubTabNet [40] dataset. IoU threshold is set to 0.6. **TEDS:** indicates Tree-Edit-Distance-based Similarity.

| Method | Training | | Average Over Test Set | | |
|---|---|---|---|---|---|
| | Dataset | #Image | Document Type | | |
| | | | Text | Latex | Text+Latex |
| | | | BLEU↑ | BLEU↑ | BLEU↑ |
| Image-to-Text (Word) [17] | TableBank [17] | 145K | 0.751 | 0.673 | 0.714 |
| Image-to-Text (Latex) [17] | TableBank [17] | 145K | 0.405 | 0.765 | 0.582 |
| Image-to-Text (Word+Latex)[17] | TableBank [17] | 145K | 0.712 | 0.765 | 0.738 |
| TabStruct-Net [24] | SciTSR [4] | 012K | 0.914 | 0.937 | 0.916 |
| TOD-Net[†]+TSR+PP | FinTabNet [39] | 091K | 0.920 | 0.944 | 0.927 |
| TOD-Net[‡]+TSR+PP | FinTabNet [39] | 091K | 0.924 | 0.955 | 0.937 |

Table 3. Shows the performances of the proposed networks and its comparison against the existing techniques on logical table structure recognition in TableBank [17] dataset.

| Method | IoU | FinTabNet | | ICDAR-13 | | Sci-TSR | | TUCD | |
|---|---|---|---|---|---|---|---|---|---|
| | | CD-F1↑ | TSR-F1↑ | CD-F1↑ | TSR-F1↑ | CD-F1↑ | TSR-F1↑ | CD-F1↑ | TSR-F1↑ |
| TabStruct-Net [24] | | 0.911 | 0.898 | 0.922 | 0.904 | 0.899 | 0.876 | 0.916 | 0.90 |
| TOD-Net$^\dagger$ +TSR-Net+PP | 0.5 | 0.930 | 0.906 | 0.92 | 0.903 | 0.907 | 0.88 | 0.911 | 0.889 |
| TOD-Net$^\ddagger$ +TSR-Net+PP | | **0.951** | **0.944** | **0.924** | **0.904** | **0.92** | **0.894** | **0.932** | **0.918** |
| TabStruct-Net [24] | | 0.874 | 0.848 | 0.91 | 0.886 | 0.891 | 0.864 | 0.896 | 0.871 |
| TOD-Net$^\dagger$ +TSR-Net+PP | 0.6 | 0.915 | 0.892 | 0.918 | 0.903 | 0.90 | 0.878 | 0.911 | 0.889 |
| TOD-Net$^\ddagger$ +TSR-Net+PP | | **0.944** | **0.920** | **0.926** | **0.904** | **0.92** | **0.894** | **0.932** | **0.918** |
| TabStruct-Net [24] | | 0.732 | 0.704 | 0.746 | 0.72 | 0.705 | 0.682 | 0.74 | 0.722 |
| TOD-Net$^\dagger$ +TSR-Net+PP | 0.7 | 0.830 | 0.802 | 0.846 | 0.82 | 0.785 | 0.746 | 0.822 | 0.797 |
| TOD-Net$^\ddagger$ +TSR-Net+PP | | **0.884** | **0.868** | **0.88** | **0.852** | **0.858** | **0.823** | **0.866** | **0.839** |
| TabStruct-Net [24] | | 0.547 | 0.496 | 0.644 | 0.597 | 0.611 | 0.565 | 0.638 | 0.582 |
| TOD-Net$^\dagger$ +TSR-Net+PP | 0.8 | 0.603 | 0.561 | 0.726 | 0.675 | 0.683 | 0.637 | 0.70 | 0.659 |
| TOD-Net$^\ddagger$ +TSR-Net+PP | | **0.715** | **0.680** | **0.80** | **0.748** | **0.76** | **0.714** | **0.792** | **0.735** |
| TabStruct-Net [24] | | 0.195 | 0.120 | 0.375 | 0.292 | 0.326 | 0.255 | 0.347 | 0.289 |
| TOD-Net$^\dagger$ +TSR-Net+PP | 0.9 | 0.368 | 0.325 | 0.396 | 0.307 | 0.347 | 0.296 | 0.37 | 0.301 |
| TOD-Net$^\ddagger$ +TSR-Net+PP | | **0.452** | **0.404** | **0.508** | **0.454** | **0.425** | **0.368** | **0.464** | **0.408** |

Table 4. Shows the comparison between the performances of the proposed network and TabStruct-Net [?] on cell detection and table structure recognition of dataset over various IoU thresholds. CD: indicates cell detection and TSR: indicates table structure recognition. TOD-Net$^\dagger$: indicates TOD-Net for direct cell detection and TOD-Net$^\ddagger$: indicates cell detection through intersection of row and column predictions using TOD-Net. We use FinTabNet [?] dataset for training.

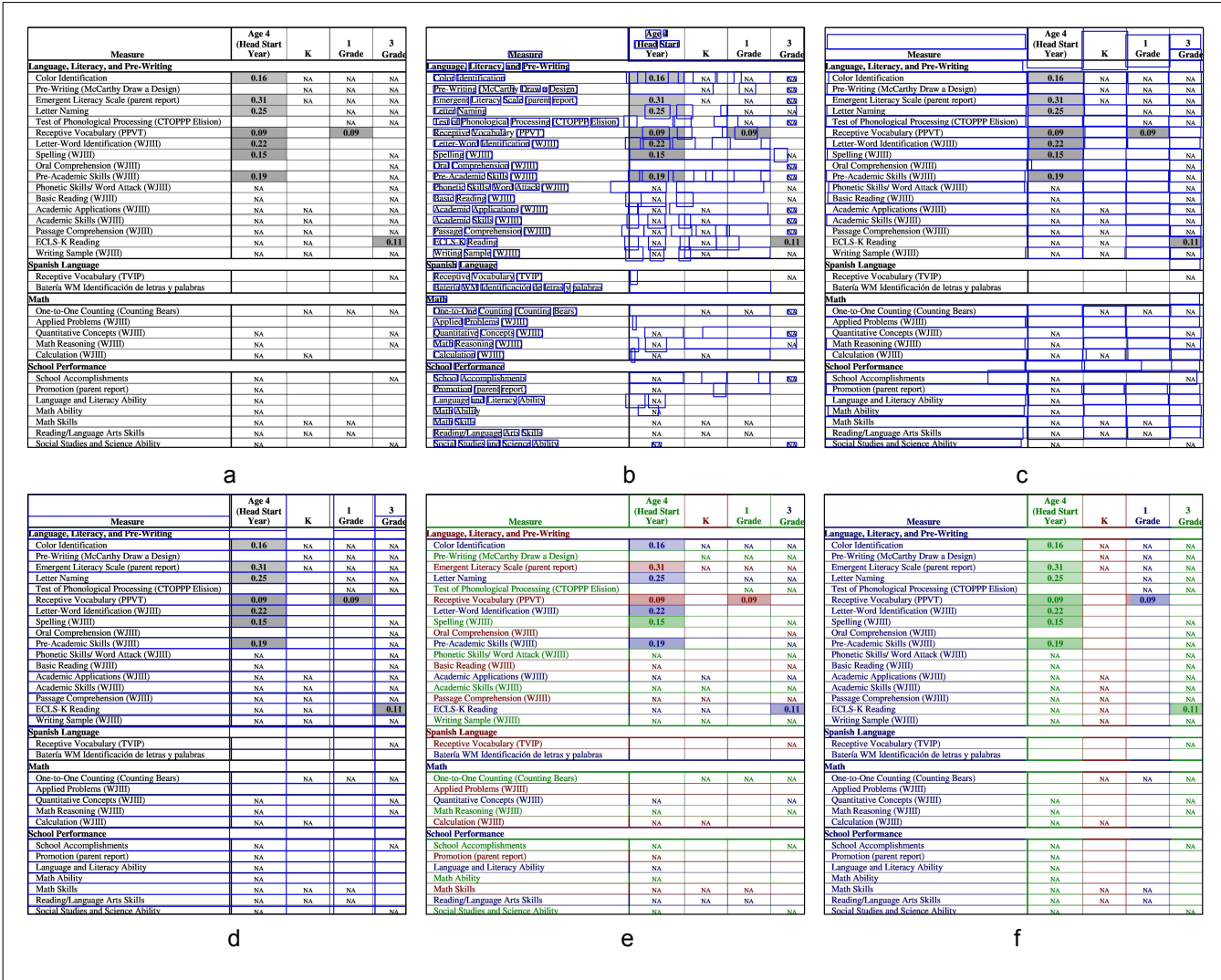Figure 1. Present modification of content level bounding box annotation to cell bounding box annotation using the proposed pre-processing algorithm. (a) shows original ground truth content level bounding box annotations, (b) intermediate annotation of cell level bonding box, and (c) final annotation of cell level bounding box.

Figure 2. Present modification of content level bounding box annotation to cell bounding box annotation using the proposed pre-processing algorithm. (a) shows original ground truth content level bounding box annotations, (b) intermediate annotation of cell level bonding box, and (c) final annotation of cell level bounding box.

a

b

c

d

e

f

Figure 3. Illustrates visual results of ICDAR-2013. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.
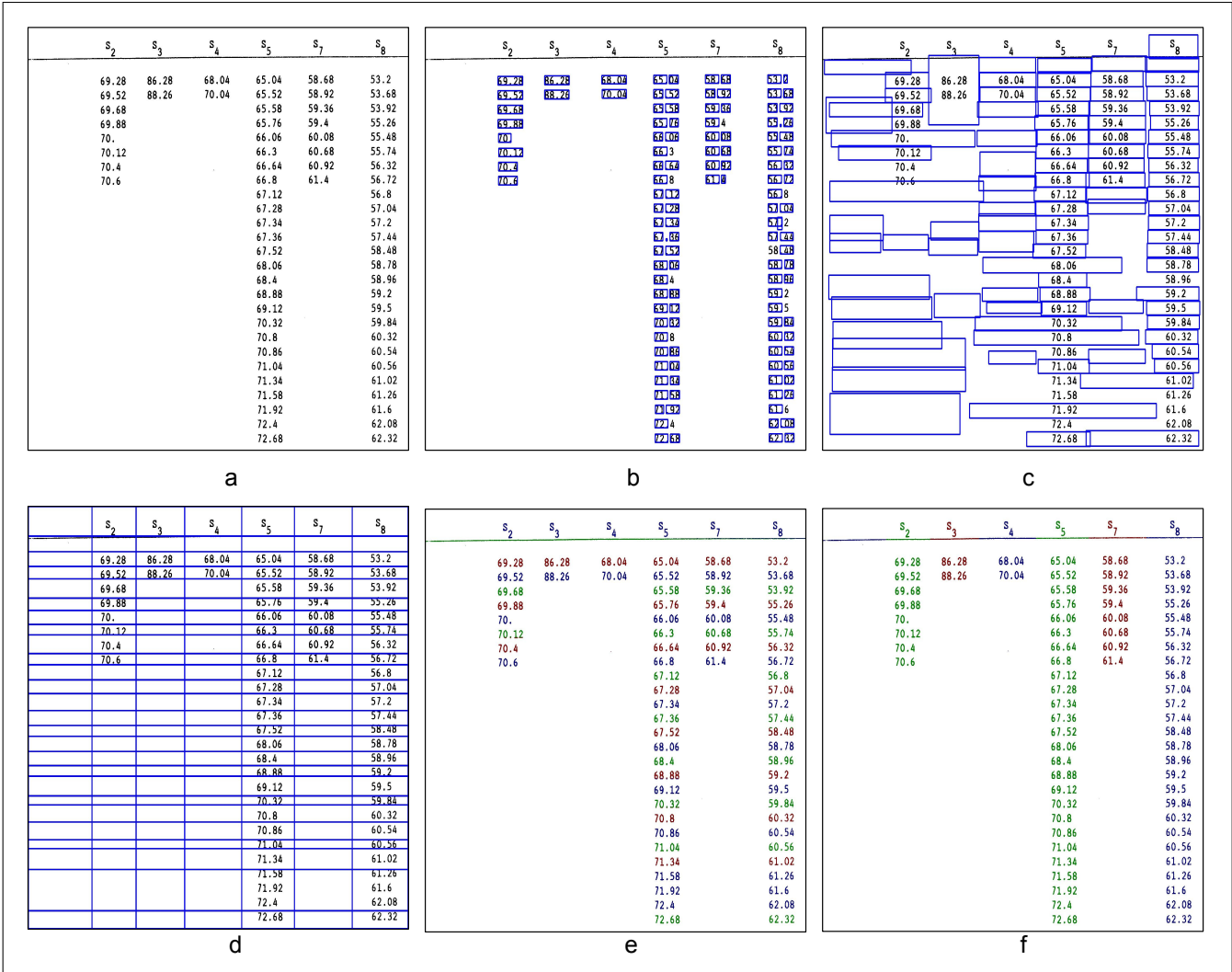
Figure 4. Illustrates visual results of ICDAR-2019. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.

Figure 5. Illustrates visual results of UNLV. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.
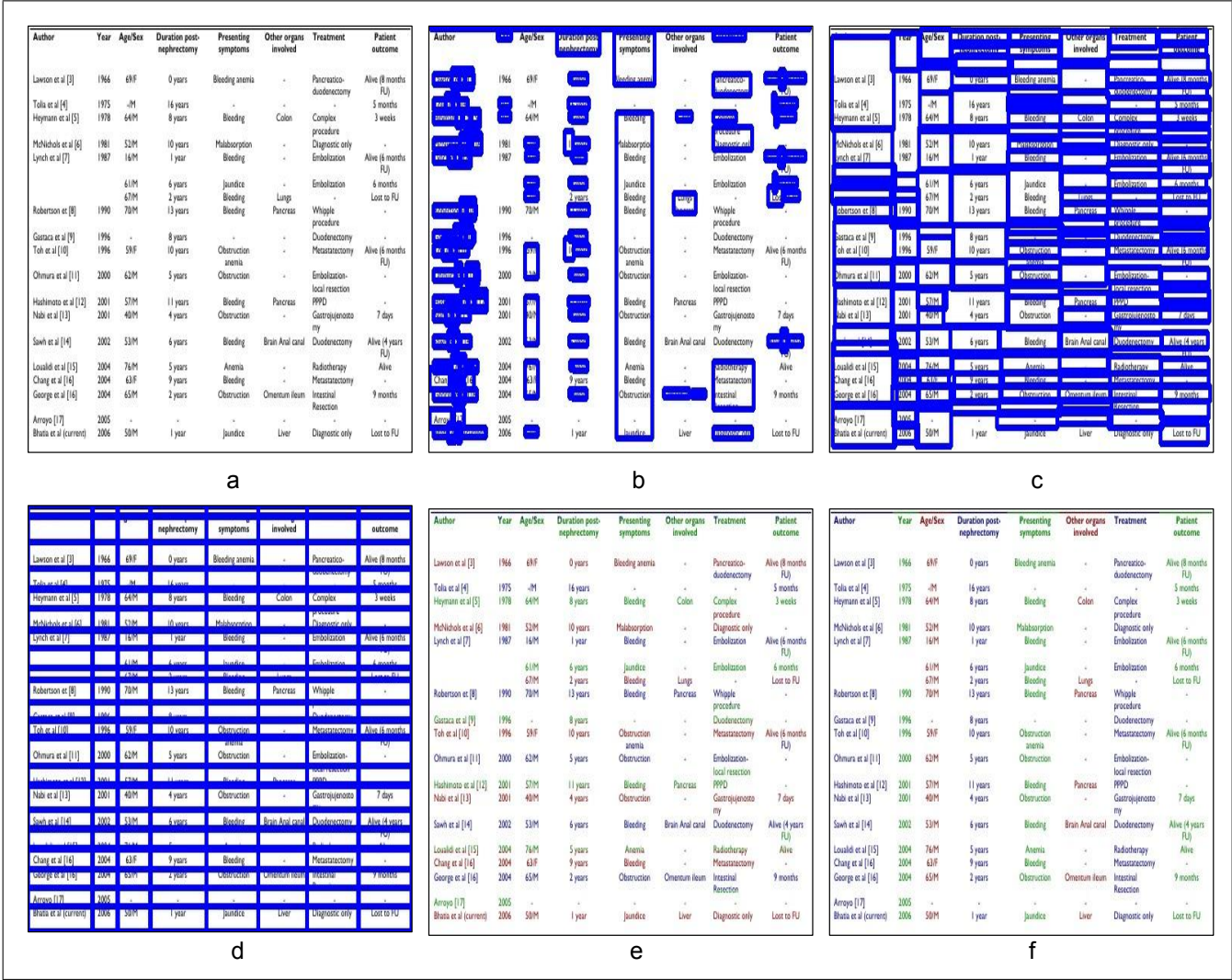
Figure 6. Illustrates visual results of TableBank. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.
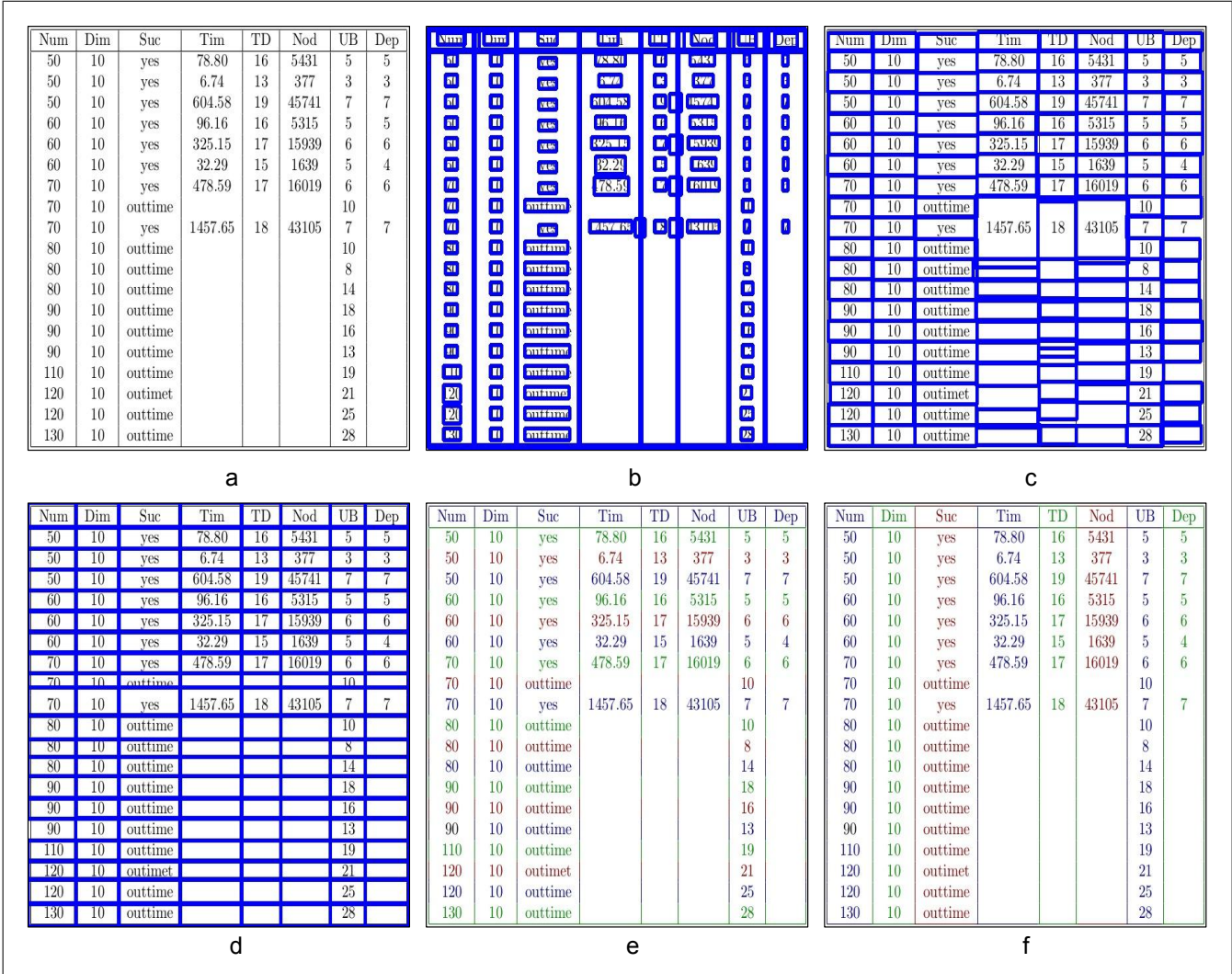
Figure 7. Illustrates visual results of PubTabNet. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.

Figure 8. Illustrates visual results of SciTSR. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.

Figure 9. Illustrates visual results of TUCD. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.

Figure 10. Illustrates visual results of FinTabNet. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.

Figure 11. Illustrates visual results of FinTabNet. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.

Figure 12. Illustrates visual results of FinTabNet. (a) original table image, (b) word detection using Tesseract, (c) cell detection uding TOD-Net, (d) cell detection using intersection of detected row and column using TOD-Net, (e) cells which are row adjacent using TSR-Net, and (f) cells which are column adjacent using TSR-Net.
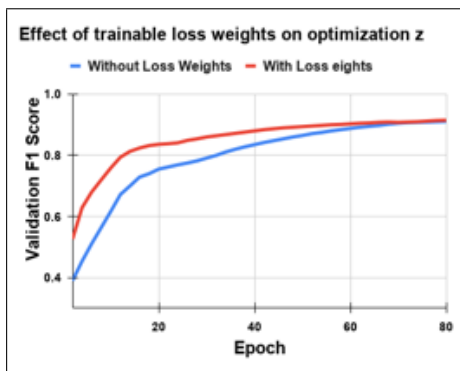
Figure 13. Shows the plot of validation accuracy against epochs. Blue line indicates scores without the use of loss weights, while red line indicates scores when trainable loss weights were used. In both the cases, same weight initialization was used.