

TDATR: Improving End-to-End Table Recognition via Table Detail-Aware Learning and Cell-Level Visual Alignment

Supplementary Material

A. Document Data Processing

For data from different sources, we employed distinct processing workflows due to their varying formats [3, 28].

Chinese and English webpages: We render HTML file to image using `khtmltopdf`¹. Then we utilized a commercial OCR² engine to recognize text lines on the webpages, extracting both the textual content and their corresponding coordinates.

Chinese and English papers: For papers with LaTeX source code, we first compile the LaTeX code into a PDF, and then use the `PyMuPDF`³ parser to extract text lines and their coordinates from the compiled PDF. For papers available only in PDF format, we utilize a commercial engine to extract text lines and coordinates. Specifically, for mathematical formulas in papers, we employ `LatexOCR`⁴ tool.

README files: We downloaded README files and their referenced content from various GitHub projects. First, we filter out invisible elements from the README files, such as web links, jump markers, and comments, to ensure consistency between the text and rendered images. We then used `Pandoc`⁵ to convert the filtered README files into HTML. Finally, we utilized `wkhtmltopdf` to convert the HTML content into images. To limit the image size, we segmented the images and extracted the corresponding mark-down content as labels.

WuKong dataset [12] and in-house data: We utilized a commercial OCR engine to recognize text lines on images.

B. Table Data Processing

Real-world table refers to images captured through photographing or scanning. Such images often contain geometric distortions, background noise, and low resolution, making recognition considerably more challenging. Digitally-born table refers to images rendered directly from code or digital documents. These images have clean characters and well-aligned layouts.

B.1. Unified Multi-source Table Data Processing

To obtain labels for the table auxiliary tasks from various datasets, we designed a unified processing pipeline.

In the first step, we unify the table label from different sources into a consistent format. In document images, we

represent a table using table box, table cells, and table grids. The table box indicates the position of the table area within the document image. Table cells contain cell coordinates, logical coordinates, the text content within each cell and cell ID. Table grids represent the fine-grained structure of a table, showing the results after splitting merged cells. Each table grid includes the ID of the corresponding cell and its coordinates.

In the second step, we conduct data cleaning to eliminate inconsistently labeled table data, ensuring high data quality. First, we remove table data with overlapping logical coordinates for cells. Next, we exclude entries with incomplete table grids, specifically those where grids have not been assigned to their corresponding cells. Finally, we eliminate redundant table grids, which occur when adjacent grid rows and grid columns are identical.

The last step is training data generation. We extract table images from document images by cropping based on the table box. Table cell information is used for label generation in the table HTML parsing task, table cell detection task, and table cell spotting task. Table grid information is utilized for label generation in the table span cell detection task and the table row and column detection task.

B.2. Table Data Augmentation

High-quality labeled table data for photographic scenes is limited [57, 64]. We expanded the `iFLYTAB` [64] dataset inspired by an identity matrix-based augmentation method [6], resulting in the `iFLYTAB-Aug` dataset with 82.5k samples.

We made the following modifications to the identity matrix-based augmentation to ensure the generation of complex and realistic table data.

- We restrict the selected table regions to have more than 4 rows and columns.
- We ensure that the selected table sub-region always contains at least one span cell, and all rows and columns containing the span cell are retained. This ensures that the table has a complex structure.
- For wireless tables, the selected region always starts from the first row and the first column. Because the row and column headers provide essential information for distinguishing between rows and columns.

C. Implementation Details

In this section, we provide the detailed input–output designs of the table detail-aware learning tasks, as illustrated in the

¹<https://wkhtmltopdf.org/>

²<https://www.xfyun.cn/services/common-ocr>

³<https://github.com/pymupdf/PyMuPDF>

⁴<https://github.com/lukas-blecher/LaTeX-OCR>

⁵<https://pandoc.org/>

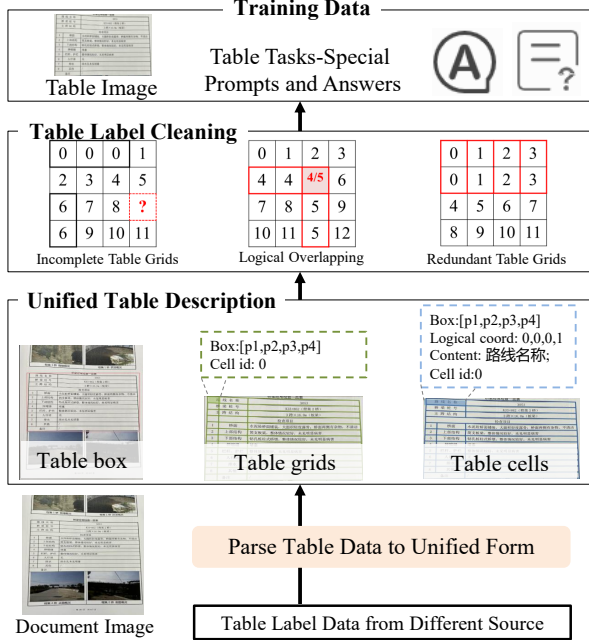


Figure 4. The pipeline of unified multi-source table data processing. The pipeline normalizes heterogeneous table annotations from various sources into a unified representation for model training.

Fig. 5 and 6.

D. Baseline Protocol

Thanks for pointing out the ambiguity in Tables 3 and 4. Our compared baselines can be grouped into: (1) Dataset-specific setting: methods fine-tuned on each target dataset. (2) Unified setting (marked with “†”): a single checkpoint evaluated across multiple datasets. Table 9 presents the training data configurations of all baseline methods used in this paper.

E. Structure-guided Cell Localization Module

We leverage logical relationships between cells to perform bidirectional structure-guided enhancement. We take the row-based enhancement as an example to illustrate the computation process. We first project the cell representation C' into a row feature space using a linear layer to obtain row-level similarity features C^{row} , as shown in Eq. 2. We then compute pairwise similarity scores via inner product to estimate whether two cells belong to the same row. After thresholding, we obtain the row similarity matrix, i.e., a binary relationship matrix indicating which cell pairs share the same row, as defined in Eq. 3. As illustrated in Fig. 7, Cell 2 and Cell 3 are in the same row, thus $M_{2,3}^{row} = 1$. This matrix is subsequently used as a mask in self-attention to reinforce feature interactions among cells within the same

row.

F. Evaluation Benchmarks

iFLYTAB-full obtains 5,419 test samples. The samples come from diverse sources—including screenshots, scans, and camera-captured images—covering a wide range of image qualities that allow evaluation of model robustness. The dataset exhibits large variations in image resolution, testing the model’s capability to handle multi-resolution inputs. It also contains grid tables, bordered three-line tables, and borderless tables. The absence of visible cell boundaries in borderless tables introduces significant challenges for TR.

TabRecSet contains 7,548 validation samples, all captured in real-world scenarios with strong perspective distortion and low image quality. Borderless and three-line tables are generated by erasing the ruling lines of grid tables, creating a domain gap between these synthetic styles and real-world data. The dataset includes both Chinese and English tables.

PubTabNet consists of 9,015 validation samples and 9,064 test samples, with the validation set commonly used for benchmarking. Its annotations are produced by an automated pipeline, resulting in low-resolution images and inconsistent visual-HTML alignment (e.g., cell over-segmentation). Such inconsistencies lead to contradictory training signals and may underestimate performance during evaluation. Models often require dataset-specific fine-tuning to adapt to these inconsistencies.

PubTables-1M contains 93,834 test samples and is sourced from the same corpus as PubTabNet. It applies automated consistency checks to correct the annotation inconsistencies present in PubTabNet, resulting in significantly improved label reliability.

OmniDocBench v1.5. Following PaddleOCR-VL, we crop 512 table samples from the benchmark. The dataset covers a wide spectrum of table types, including challenging note-style tables where continuous content and background ruling lines visually disrupt cell boundaries, often causing over-segmentation. Successful recognition requires semantic understanding of cell content beyond visual boundary cues.

CC-OCR. The 300 table test samples in CC-OCR cover both Chinese and English, spanning real-world and digital-document scenarios. The dataset includes long tables, dense tables, and heavily rotated cases, posing significant challenges for structure parsing and spatial reasoning.

OCRBench includes 700 table-related samples in both Chinese and English. Using the provided table boxes and our internal table detector, we crop table regions for recognition. Many samples come from financial reports, whose formatting introduces unique difficulties, e.g., large spacing between “\$” and numbers is easily mistaken for column separators.

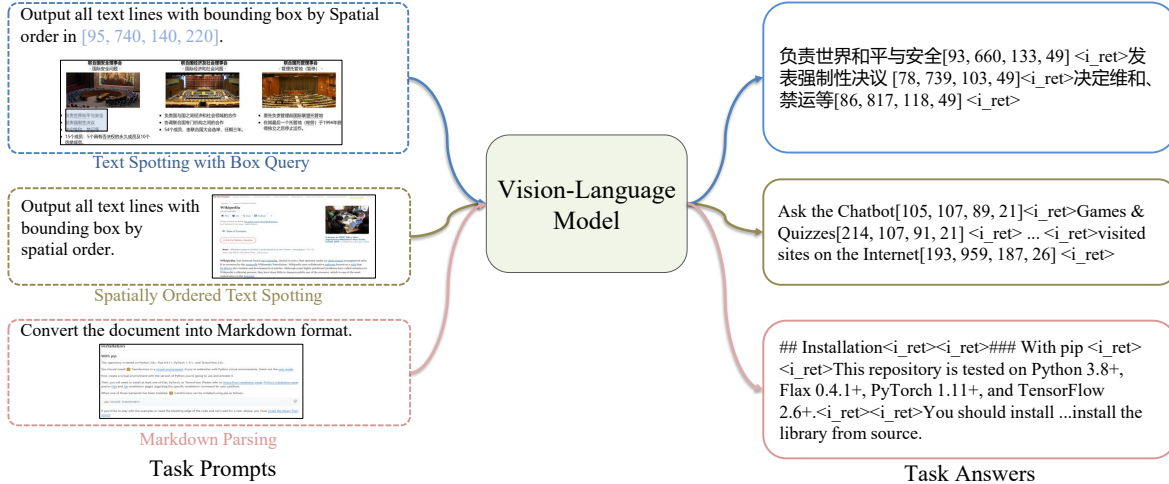


Figure 5. Illustration of table content recognition tasks. These tasks leverage diverse document data to enable text recognition, text localization, and reading-order understanding.

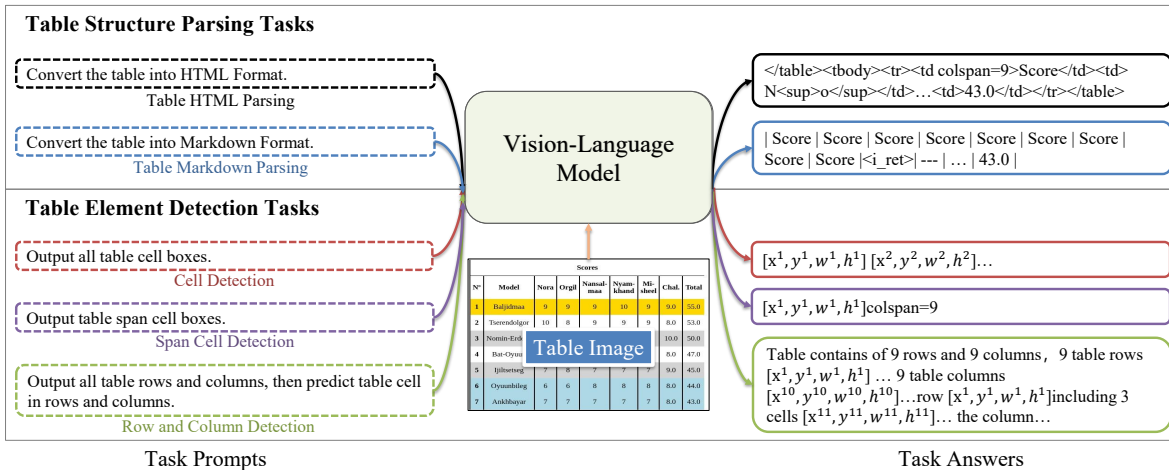


Figure 6. Illustration of table structure understanding tasks. These tasks equip the model with structure-awareness from both the cell level and the row/column level.

G. Single-dataset Training Variant

We conduct a single-dataset comparison by performing only-PubTabNet table detail fusion fine-tuning starting from our table detail-aware pretrained model. On PubTabNet-val, we achieve TEDS-S 96.78 / TEDS 96.10, outperforming the second-best TR dataset-specific baseline, TableFormer (96.75 / 93.60). This supports the effectiveness of our “perceive-then-fuse” paradigm in a single-dataset setting.

H. SGCL Inference Efficiency

TDATR leverages SGCL to localize cells in parallel, conditioned on the generated cell tokens. Since TR baselines such as Dolphin or EDD do not output cell boxes, a direct

efficiency comparison is not applicable. For a fair comparison, we implement a matched baseline, “ED Loc Gen” in Table. 6, that autoregressively generates discretized coordinates after the cell tokens. We evaluate 40 randomly sampled PubTabNet images (max side length 1024), with an average of 26.75 cells and 190.38 TR tokens. Measured on an NPU with batch size 1, TDATR achieves 9.7s end-to-end latency, compared to 15.7s for the baseline (1.6× faster). Importantly, SGCL contributes only 0.28s to the end-to-end latency, confirming that parallel refinement keeps localization overhead low. TDATR and the baselines have comparable max reserved memory (15.77 GiB vs. 15.36 GiB).

Table 9. Summary of the training data configurations of the baseline methods. For each method, we report the paradigm, table training data, auxiliary data, whether table-specific fine-tuning is applied, and additional notes.

Method	Paradigm	Table training data	Extra data	Dataset specific	Notes
TableMaster	TSR	PubTabNet	–	Yes	–
LORE	TSR	PubTabNet, TabRecSet, and iFLYTAB	–	Yes	20k images were randomly sampled from PubTabNet for training. TabRecSet and iFLYTAB were reproduced by us based on the released code.
BGTR (PT)	TSR	TabRecSet, iFLYTAB, PubTabNet, FinTabNet and SynthTabNet	–	Yes	–
UniTabNet	TSR	iFLYTAB, PubTables-1M, and PubTabNet.	Pre-training: a synthetic dataset comprising 1.4 million Chinese and English samples from SynthDog, and PubTables-1M	Yes	–
EDD	E2E-TR	PubTabNet	–	Yes	–
SEMr3 + PPOCR	M-TR	PubTabNet and iFLYTAB	PPOCR relies on general text recognition data	Yes	“+PPOCR” indicates that the cell content is obtained from the PPOCR model.
GTE	TSR	PubTabNet and FinTabNet	–	Yes	The model is pre-trained on PubTabNet and fine-tuned on multiple datasets.
Davar-Lab	TSR	PubTabNet	–	Yes	–
LGPMA + R2AM	M-TR	PubTabNet	Additional data required by R2AM	Yes	“+R2AM” indicates that the cell content is obtained from the R2AM model.
TableFormer + GT	M-TR	PubTabNet, FinTabNet, and SynthTabNet	–	Yes	“+GT” indicates that the ground-truth cell content.
RapidTable	M-TR	–	–	–	–
OmniParser	OCR-VLM	PubTabNet and FinTabNet	Large-scale document parsing data	Yes	–
DocOwl1.5	OCR-VLM	TURL and PubTabNet	Unified structure-learning data from documents, webpages, charts, and natural images	No	–
Dolphin	OCR-VLM	PubTabNet and PubTab1M	Large-scale document parsing data	Yes	–
MinerU2.5	OCR-VLM	In-house	Large-scale document parsing data	No	–
DeepSeek-OCR	OCR-VLM	In-house	Large-scale document parsing data	No	–
PaddleOCR-VL	OCR-VLM	In-house	Large-scale document parsing data	No	–
dots.ocr	OCR-VLM	In-house	Large-scale document parsing data	No	–
GOT	OCR-VLM	In-house	Large-scale document parsing data	No	–
TabPedia	TSR	PubTabNet and PubTab1M	–	No	–
DETR + PDF	M-TR	PubTables-1M	–	Yes	“+PDF” indicates that the cell content is obtained from the source PDF files.

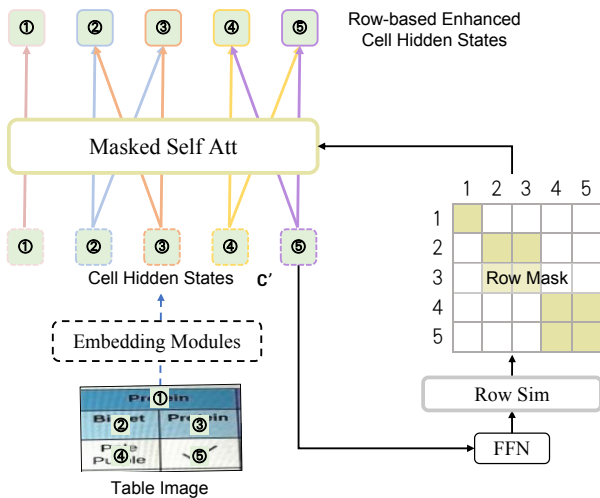


Figure 7. The architecture of the structure-guided cell localization module, illustrated with the row-based cell enhancement example.

I. Visualization of Table Recognition

We visualize several challenging table samples. Real-world tables (Fig. 8) contain background noise, perspective distortion, and uneven illumination. Long tables (Fig. 9) feature lengthy sequences, numerous cells, and long text contents. Complex-structure tables (Fig. 10) include extensive row or column spanning. Our method performs robustly across all these cases, demonstrating strong generalization and effectiveness.

序号	项目	考核内容	评分标准	检查情况	得分
路面工程 (4分)	涵洞 (1分)	涵洞进出口、洞身及帽石、八字墙、一字墙顺直、好透,洞内无杂物、淤泥及积水现象,管身、涵底铺砌、拱圈、盖板无裂缝,涵洞处路面平顺、无跳车现象等	现场检查,每发现一处不合格扣0.5分,扣完为止	/	/
	基层 (1分)	表面平整、无坑洼;施工接茬平整;芯样完整、密实;材料质量、干净等	现场检查,每发现一处不合格扣0.5分,扣完为止	/	/
	水泥混凝土面层 (3分)	混凝土板表面无蜂窝、印痕、裂纹、缺边掉角;表面抗滑构造符合要求,表面无严重泛砂,芯样级配良好;无断板裂板;板缝无明显缺陷	现场检查,每发现一处不合格扣0.5分,扣完为止	混凝土边缝缺边掉角,表面有裂纹	/
	路肩 (1分)	路肩培土平整、密实	现场检查,每发现一处不合格扣0.5分,扣完为止	/	/
交通安全设施 (3分)	交通安全标志线及护栏 (3分)	标志面不得有划痕、气泡、翘翘、变形、开裂,着色不均匀等现象,标线应清晰、无污迹、无气泡及无网状裂缝等现象,玻璃珠洒布应均匀,附着牢固,反光均匀,护栏波形梁线形应顺适,色泽一致,立柱顶部无明显磕边、变形、开裂	现场检查,每发现1处扣0.5分,扣完为止,工程交(竣)工没有及时设置标志标线,此项不得分	/	/

序号	项目	考核内容	评分标准	检查情况	得分
路面工程 (4分)	涵洞 (1分)	涵洞进出口、洞身及帽石、八字墙、一字墙顺直、好透,洞内无杂物、淤泥及积水现象,管身、涵底铺砌、拱圈、盖板无裂缝,涵洞处路面平顺、无跳车现象等	现场检查,每发现一处不合格扣0.5分,扣完为止	/	/
	基层 (1分)	表面平整、无坑洼;施工接茬平整;芯样完整、密实;材料质量、干净等	现场检查,每发现一处不合格扣0.5分,扣完为止	/	/
	水泥混凝土面层 (3分)	混凝土板表面无蜂窝、印痕、裂纹、缺边掉角;表面抗滑构造符合要求,表面无严重泛砂,芯样级配良好;无断板裂板;板缝无明显缺陷	现场检查,每发现一处不合格扣0.5分,扣完为止	混凝土边缝缺边掉角,表面有裂纹	/
	路肩 (1分)	路肩培土平整、密实	现场检查,每发现一处不合格扣0.5分,扣完为止	/	/
交通安全设施 (3分)	交通安全标志线及护栏 (3分)	标志面不得有划痕、气泡、翘翘、变形、开裂,着色不均匀等现象,标线应清晰、无污迹、无气泡及无网状裂缝等现象,玻璃珠洒布应均匀,附着牢固,反光均匀,护栏波形梁线形应顺适,色泽一致,立柱顶部无明显磕边、变形、开裂	现场检查,每发现1处扣0.5分,扣完为止,工程交(竣)工没有及时设置标志标线,此项不得分	/	/

序号	项目	摘要	预算费用	备注	支出明细		支出明细		结余金额	备注	
					摘要	支出金额	支出方式	摘要			支出金额
1	材料费	机器折旧	850	明细见《表一》	项目1	600.00	建行	项目1	600.00	建行	6,900.00
		易耗材料	20810		项目2	500.00	农行	项目2	500.00	农行	8,000.00
		机器易耗配件	850		项目3	900.00	现金	项目3	900.00	现金	9,600.00
2	人工成本	工资	3200	明细见《表一》	项目4	1,200.00	建行	项目4	1,200.00	建行	11,600.00
		员工保险费	90		项目5	1,500.00	农行	项目5	1,500.00	农行	11,600.00
		员工房租	1600		项目6	2,100.00	现金	项目6	2,100.00	现金	13,700.00
3	累计成本	33130	3	项目7	2,500.00	建行	项目7	2,500.00	建行	13,700.00	
4	利润	10%	3313	3313	项目8	3,000.00	农行	项目8	3,000.00	农行	13,900.00
5	税金	6%	2326	累计成本+利润×6%	项目9	3,500.00	现金	项目9	3,500.00	现金	13,000.00
6	总计		38769		项目10	4,000.00	农行	项目10	4,000.00	建行	10,800.00
					项目11	4,500.00	农行	项目11	4,500.00	农行	9,800.00
					项目12	5,000.00	现金	项目12	5,000.00	现金	9,000.00

Figure 10. Visualization of table recognition results complex-structure tables. In each subfigure, the left shows the input original table image, and the right presents the HTML-rendered visualization of the corresponding recognition result.

J. Visualization of Cell Localization

We qualitatively compare the cell localization results of several SOTA models, as shown in Fig. 11. SEMv3, which follows a “split-and-merge” strategy by detecting row/column separators to form table grids, is prone to confusing separators with inter-word gaps. LORE employs CornerNet for cell localization and relies solely on visual cues, making it unreliable for empty cells. UniTabNet predicts cell boxes through a single cell token, but compressing spatial information into one token limits its performance on dense tables. “ED Loc Gen” generates cell coordinates sequentially, resulting in excessively long answer sequences that are easily truncated on long tables.

K. Failure Cases Analysis

We observed that these hard cases on iFLYTAB-full are mainly fall into three main error types: (1) Boundary confusion: In borderless tables containing multi-line text the model struggles to distinguish text line spacing from cell delimiters. (2) Span number errors: For cells with large row/column spans (more than 15), the model occasionally predicts the error number. (3) Localization instability: Dense and empty borderless cells lack explicit visual cues causing instability visual-based regression in SGCL. We plan to enhance the decoder’s semantic reasoning to re-

solve these visual ambiguities.

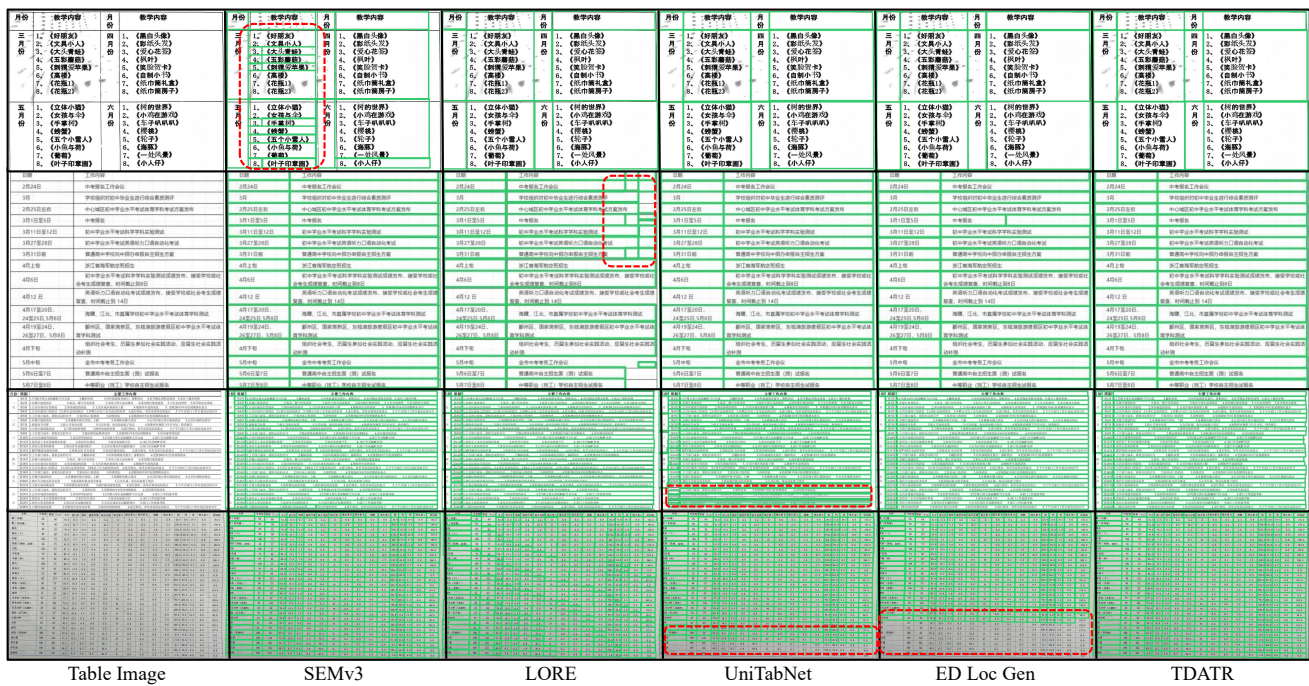


Figure 11. Qualitatively comparison of the cell localization.