

# Recognition and Extraction of Tabular Data from Scanned Document Images

Information Technologies Institute  
Centre of Research & Technology - Hellas  
6th km Xarilaou - Thermi, 57001, Thessaloniki, Greece

**Abstract**— In this paper we propose a heuristics-based method for automatic table detection in document images. Our method works on the output of Tesseract, Google’s open-source OCR engine, which provides the bounding box and the text of each word in the document. These elements are grouped on a bottom-up way, in order to identify tabular arrangements of data. The approach aims at improving the results of OCR, when used on documents with tables.

**Index Terms**—table recognition, OCR, scanned document images.

## I. INTRODUCTION

Tables are visually oriented arrangements of text, widely used to present complex information and data to human readers. While for regular text it is sufficient to acquire just the textual content of each word, in the case of tables it is critical to identify their basic building blocks (table cells) and group them together in rows and columns. Thus, in order to extract and understand the information stored in a table, we have to process the raw data generated by an OCR engine and recreate the table’s layout and structure.

Kieninger & Dengel[1] propose a bottom-up approach for table recognition, by grouping words in segments and looking for horizontally overlapping structures in the text. Oro & Ruffolo[2] further build on this approach, by adding line-by-line analysis and selecting initial table areas, before moving to the bottom-up grouping of the data. Our method is partially based on the above approaches, however we propose a novel method for the grouping of the word segments in table columns and multiple-line table rows.

Section 2 describes in detail our approach, section 3 presents some experimental results in comparison with the state-of-the-art commercial OCR engine ABBYY FineReader 11 and section 4 proposes some future improvements.

## II. PROPOSED METHOD

The goal of our approach is to export in HTML format the image of a document, while reconstructing any tables present in the document.

The steps of the algorithm are the following:

### A. Image Processing

Initially, some image processing is required, for the OCR engine to generate better results. The input image is binarized

and resized (300 dpi, average letter height 40-60 px). Since Tesseract does not support multiple-column text, we check for continuous vertical white areas in the image and segment it in blocks, each one corresponding to a text column. Finally, we remove any horizontal and vertical grid lines, because the OCR engine has trouble distinguishing text when it is surrounded by thick dense grid lines.

### B. Optical Character Recognition

The processed image is given as input to the OCR engine, which in return exports for each recognized word:

- A string containing the word characters.
- The Cartesian coordinates of the bounding box of the word.
- The recognition confidence of the word
- Font information for the word (font size, bold, italic, underlined)
- Dictionary information (whether the word is found in Tesseract’s dictionary or not).

### C. Text Lines, Word Segments and Line Types

Using the coordinates for each word, we create text lines from words that overlap vertically. For each line, we create word segments, by merging together words whose horizontal distance is below a threshold  $T_{seg}$  (dense text). Finally, we assign each line to a specific type:

- Type 1 / Text – Lines with a single long segment (longer than half of the text width).
- Type 2 / Table – Lines with more than one segments.
- Type 3 / Unknown – Lines with a single short segment.

In case of multiple continuous images, the top and bottom areas of each image are checked for similar repetitive word segments. When found, these segments are considered footers/headers and are removed, except from the headers on the first image. Next, the images are processed as a single continuous document.

### D. Initial Table Areas

We create initial table areas, by grouping together adjacent type 2 and type 3 lines. However, for a type 3 line to be assigned to a table area, it must have at least one type 2 line above it (a table cannot begin with a type 3 line at the top).

### E. Table Column Generation

Next, we select word segments as column generators, for each table area. The selection algorithm is the following:

1. Select the left-most word segment that is not assigned to a column
2. Find all the word segments that horizontally almost align with it (the horizontal distance between the left limit of the word segments is below a threshold  $T_{col}$ )
3. Find the average length of these word segments.
4. Select as a column generator, the word segment that is closest to the average length found on step 4.
5. Assign to a new column all the word segments that horizontally overlap with the column generator (some longer word segments can be assigned to more than one column).
6. Go back to step 1, until all the segments in a table area are assigned to a column.

Next, the newly generated column are customized, based on the following criteria:

- If a column has only one segment, which is on the 1<sup>st</sup> table row (possibly misaligned table header), while the column on its left does not have a segment in the same row, these two columns are merged.
- If a column has more multiple-column segments than single-column segments, it is merged to the column on its left
- Tables that end up with a single column are discarded and treated as simple text.
- Type 3 lines at the bottom of a table, with a multiple-column segment, are removed from the table
- Scrambled tables, i.e. tables that have a very inconsistent format, usually generated by random word segments with big white spaces between them (justified text alignment), are discarded and treated as simple text
- If all the columns of a table have more empty cells than cells with data, the table is discarded (simple formatted text).

### F. Multiple-Line Table Rows

Some of the rows in the table areas are merged to multiple-line rows, based on the following criteria:

- If a table row does not have a word segment in the first column, and there is one-to-one correspondence with the word segments of the table row above it, these two rows are merged together.
- If a type 3 table row has its single segment assigned to the first column, and the row below it has a segment in the first column which is indented to the right, these two rows are merged together.

- Tables that end up with a single row or with two single-line rows are discarded and treated as simple text.

### G. Data Extraction

Finally the text and tabular data are extracted in HTML format, using a custom made HTML writer.

## III. EXPERIMENTAL RESULTS

Below are displayed some experimental results, covering multiple table layouts. The figures displayed are the input image, the HTML generated by our and the HTML generated by state-of-the-art commercial OCR engine ABBYY FineReader11 is displayed on the right.

The algorithm successfully recognizes and reconstructs the table in most cases, and in some cases it even outperforms the commercial OCR engine.

Spot	mean.co	mean.pd	SD.co	SD.pd	log2.fold	p.val.Wilcox	LDA.Coeff	Delta.norm
1	86	4.48E-03	1.56E-03	2.88E-03	9.72E-04	-1.521	3.32E-04	-425.75
2	87	9.34E-03	4.16E-03	7.19E-03	3.91E-03	-1.166	1.26E-02	-236.86
3	335	9.41E-03	2.73E-03	9.70E-03	1.83E-03	-1.786	1.96E-03	178.98
4	362	9.56E-03	4.38E-03	6.24E-03	3.48E-03	-1.126	2.72E-03	382.08
5	365	3.95E-02	1.40E-02	2.77E-02	1.01E-02	-1.493	8.10E-04	39.32
6	368	1.44E-02	6.63E-03	9.32E-03	4.55E-03	-1.122	7.82E-03	177.71
7	369	4.31E-02	1.47E-02	2.29E-02	1.02E-02	-1.557	6.06E-05	-165.94
8	382	1.86E-02	7.94E-03	1.64E-02	3.27E-03	-1.224	3.47E-04	-100.27
9	657	5.92E-03	2.10E-03	7.31E-03	1.37E-03	-1.497	9.97E-03	-299.77

  

Spot	mean.co	mean.pd	SD.co	SD.pd	log2.fold	p.val.Wilcox	LDA.Coeff	Delta.norm
1	86	4.48E-03	1.56E-03	2.88E-03	9.72E-04	-1.521	3.32E-04	-425.75
2	87	9.34E-03	4.16E-03	7.19E-03	3.91E-03	-1.166	1.26E-02	-236.86
3	335	9.41E-03	2.73E-03	9.70E-03	1.83E-03	-1.786	1.96E-03	178.98
4	362	9.56E-03	4.38E-03	6.24E-03	3.48E-03	-1.126	2.72E-03	382.08
5	365	3.95E-02	1.40E-02	2.77E-02	1.01E-02	-1.493	8.10E-04	39.32
6	368	1.44E-02	6.63E-03	9.32E-03	4.55E-03	-1.122	7.82E-03	177.71
7	369	4.31E-02	1.47E-02	2.29E-02	1.02E-02	-1.557	6.06E-05	-165.94
8	382	1.86E-02	7.94E-03	1.64E-02	3.27E-03	-1.224	3.47E-04	-100.27
9	657	5.92E-03	2.10E-03	7.31E-03	1.37E-03	-1.497	9.97E-03	-299.77

  

Spot	mean.co	mean.pd	SD.co	SD.pd	log2.fold	p.val.Wilcox	LDA.Coeff	Delta.norm
1	86	4.48E-03	1.56E-03	2.88E-03	9.72E-04	-1.521	3.32E-04	-425.75
2	87	9.34E-03	4.16E-03	7.19E-03	3.91E-03	-1.166	1.26E-02	-236.86
3	335	9.41E-03	2.73E-03	9.70E-03	1.83E-03	-1.786	1.96E-03	178.98
4	362	9.56E-03	4.38E-03	6.24E-03	3.48E-03	-1.126	2.72E-03	382.08
5	365	3.95E-02	1.40E-02	2.77E-02	1.01E-02	-1.493	8.10E-04	39.32
6	368	1.44E-02	6.63E-03	9.32E-03	4.55E-03	-1.122	7.82E-03	177.71
7	369	4.31E-02	1.47E-02	2.29E-02	1.02E-02	-1.557	6.06E-05	-165.94
8	382	1.86E-02	7.94E-03	1.64E-02	3.27E-03	-1.224	3.47E-04	-100.27
9	657	5.92E-03	2.10E-03	7.31E-03	1.37E-03	-1.497	9.97E-03	-299.77

Fig. 1. Table with misaligned headers, no gridlines, no text. Top – input image, centre - HTML generated by our algorithm, bottom – HTML generated by ABBYY FineReader11

Search Operator	Function	Value
EQUAL	Search for words in text equal to the one in value.	Words specified according to search patterns.
DIFFERENT	Search for words in text not equal to the one in value.	
CONTAINS	Search for words in text containing the pattern specified in value (respects order, adjacency, the missing word operator(#) is supported).	
NCONTAINS	Search for words in text not containing the pattern specified in value (respects order, adjacency, the missing word operator(#) is supported).	
PHRASE_LIKE	Search for words in text containing the pattern specified in value (respect adjacency, the missing word operator(#) is not supported).	
NPHRASE_LIKE	Search for words in text not containing the pattern specified in value (respect adjacency, the missing word operator(#) is not supported).	
SEARCH_OPERATOR	Search for words in text equal to the one in value.	Words specified according to search patterns.
EQUAL	Search for words in text equal to the one in value.	
DIFFERENT	Search for words in text not equal to the one in value.	
CONTAINS	Search for words in text containing the pattern specified in value (respects order, adjacency, the missing word operator(#) is supported).	
NCONTAINS	Search for words in text not containing the pattern specified in value (respects order, adjacency, the missing word operator(#) is supported).	
PHRASE_LIKE	Search for words in text containing the pattern specified in value (respect adjacency, the missing word operator(#) is not supported).	
NPHRASE_LIKE	Search for words in text not containing the pattern specified in value (respect adjacency, the missing word operator(#) is not supported).	
SEARCH_OPERATOR	Search for words in text containing the pattern specified in value (respects order, adjacency, the missing word operator(#) is supported).	Words specified according to search patterns.
EQUAL	Search for words in text equal to the one in value.	
DIFFERENT	Search for words in text not equal to the one in value.	
CONTAINS	Search for words in text containing the pattern specified in value (respects order, adjacency, the missing word operator(#) is supported).	
NCONTAINS	Search for words in text not containing the pattern specified in value (respects order, adjacency, the missing word operator(#) is supported).	
PHRASE_LIKE	Search for words in text containing the pattern specified in value (respect adjacency, the missing word operator(#) is not supported).	
NPHRASE_LIKE	Search for words in text not containing the pattern specified in value (respect adjacency, the missing word operator(#) is not supported).	
SEARCH_OPERATOR	Search for words in text containing the pattern specified in value (respects order, adjacency, the missing word operator(#) is supported).	Value

Fig. 2. Table with multiple-line rows, with gridlines, no text. Left – input image, centre - HTML generated by our algorithm, right – HTML generated by ABBYY FineReader11

FROST TABLE FOR VICTORIA AREA				
STATION	CHANCE OF FROST ON OR AFTER THIS DATE			
	50%	33%	25%	10%
Comox	April 19	April 25	April 28	May 8
Duncan	May 1	May 5	May 10	May 23
East Sooke, Anderson Cove	April 29	May 9	May 13	May 25
Gonzales Heights	March 3	March 11	March 13	March 25
Quesnel Avenue	April 2	April 19	April 21	May 19
Saanichton	April 13	April 21	April 30	May 5
Salt Spring Island	April 3	April 10	April 16	April 27
Sidney	May 27	March 30	April 6	April 27
Tillicum Mall area	May 3	May 11	May 12	May 25
Victoria International	April 21	April 25	April 28	April 30
Victoria Marine (WiffinSpin)	March 30	April 3	April 7	May 12
William Head	March 13	March 30	April 2	April 2

  

FROST TABLE FOR VICTORIA AREA				
STATION	CHANCE OF FROST ON OR AFTER THIS DATE			
	50%	33%	25%	10%
Comox	April 19	April 25	April 28	May 8
Duncan	May 1	May 5	May 10	May 23
East Sooke, Anderson Cove	April 29	May 9	May 13	May 25
Gonzales Heights	March 3	March 11	March 13	March 25
Quesnel Avenue	April 2	April 19	April 21	May 19
Saanichton	April 13	April 21	April 30	May 5
Salt Spring Island	April 3	April 10	April 16	April 27
Sidney	May 27	March 30	April 6	April 27
Tillicum Mall area	May 3	May 11	May 12	May 25
Victoria International	April 21	April 25	April 28	April 30
Victoria Marine (WiffinSpin)	March 30	April 3	April 17	May 12
William Head	March 13	March 30	April 2	April 2

  

FROST TABLE FOR VICTORIA AREA				
STATION	CHANCE OF FROST ON OR AFTER THIS DATE			
	50%	33%	25%	10%
Comox	April 19	April 25	April 28	May 8
Duncan	May 1	May 5	May 10	May 23
East Sooke, Anderson Cove	April 29	May 9	May 13	May 25
Gonzales Heights	March 3	March 11	March 13	March 25
Quesnel Avenue	April 2	April 19	April 21	May 19
Saanichton	April 13	April 21	April 30	May 5
Salt Spring Island	April 3	April 10	April 16	April 27
Sidney	May 27	March 30	April 6	April 27
Tillicum Mall area	May 3	May 11	May 12	May 25
Victoria International	April 21	April 25	April 28	April 30
Victoria Marine (WiffinSpin)	March 30	April 3	April 17	May 12
William Head	March 13	March 30	April 2	April 2

Fig. 3. Table with a multiple-column header, no gridlines, with text. Left – input image, centre - HTML generated by our algorithm, right – HTML generated by ABBYY FineReader11

Fig. 4. Table with formatted text, partial gridlines, with text. Left – input image, centre - HTML generated by our algorithm, right – HTML generated by ABBYY FineReader11

Important Tables In SAP FI		
<b>Financial Accounting</b>		
<b>Table Name</b> <b>Description</b> <b>Important Fields</b>		
<b>Financial Accounting</b>		
FIAS	Financial Accounting "Basic"	BURKS / BELNR / GJAHB
FIAT	Accounting Document Header	BURKS / BELNR / GJAHB / BUZETI
FSVBP	Index for Vendor Validation of Double Document	BURKS / LIPRN / WAKES / BLDAT / Y
FIOPF	Inter Company Posting Procedure	FIOPG / BURKS / GJAHB / BELNR / BUZETI
FIOPP	Accounting Document Header (docs from External System)	SLBKR / BELNR / GJAHB / CLEBK
FRDPA	Run Date of a Program	PRGID
KPFA	Customer / Vendor Linking	WUCL / BURKS / VGRAS / FWTYP
RDE4	Customer Payment History	KUNNR / BURKS
RDE5	Customer Master Dunning Data	KUNNR / BURKS / MAKER
RDE6	Customer Master Bank Details	KUNNR / BURKS / BANKS / BANEN
RDE7	Customer Master Transaction Figures	KUNNR / BURKS / GJAHB / GJAHB
RDE8	Customer Master Transaction Transactions	KUNNR / BURKS / GJAHB / SHBKZ
LITB	Vendor Master Dunning Data	LITRN / BURKS / MAKER
LPBK	Vendor Master Bank Details	LITRN / BURKS / BANKS / BANEN
LF01	Vendor Master Transaction Figures	LITRN / BURKS / GJAHB / BANEN
LF02	Vendor Master Special GL Transactions	LITRN / BURKS / GJAHB / SHBKZ
YBOPF	Document Header for Standard Parking	AUDIE / BURKS / BELNR / GJAHB
FBASCORE	Financial Accounting General Services "Basic"	AUSKE / BURKS / BELNR / GJAHB
FBAL	Customer Master (Country Code)	KUNNR / BURKS
FKAL	Vendor Master (General Section)	LITRN / BURKS
FKAT	G/L Account Master (Chart of Accounts) - SPARS	KT0PL / SANMR
FKAT	G/L Account Master (Chart of Accounts) - SAKMR	KT0PL / SANMR
MANN	Accounts Blocked by Dunning Selection	KOART / LAUTP / KUNNR / BURKS / GJAHB / SHBKZ / SANMR / BURAS
MANN	Dunning Data (Account Entries)	KUNNR / BURKS
FI-GL-GL (FBS)	<b>General Ledger Accounting: Basic Functions- G/L Accounts</b>	
FIAS	G/L Account Master (Chart of Accounts) - SPARS / KT0PL / SANMR / SCHLW	
FIAS1	G/L Account Master (Company Code)	BURKS / SANMR
FI-GL-GL (FBS)	Functions - R/3 Customizing for G/L Accounts	
FIGLREP	Settings for G/L Posting Reports	MANDT
FIGLREP	General Ledger Accounting: Basic Functions - Fast Data Entry	BURKS / SANMR
FIGLREP	General Ledger Accounting: Basic Functions - Fast Data Entry	BURKS / SANMR
KOMU	Account Assignment Templates for G/L	KOMMAM / KMZEE
KOMU	Account Item	
FI-AR-AR (FBD)	<b>Accounts Receivable: Basic Functions - Customers</b>	
FI-AR-AR (FBD)	Accounts Receivable: Basic Functions - Customers	

Important Tables In SAP FI		
<b>Financial Accounting</b>		
<b>Table Name</b> <b>Description</b> <b>Important Fields</b>		
<b>Financial Accounting</b>		
FIAS	Financial Accounting "Basic"	BURKS / BELNR / GJAHB
FIAT	Accounting Document Header	BURKS / BELNR / GJAHB / BUZETI
FSVBP	Index for Vendor Validation of Double Document	BURKS / LIPRN / WAKES / BLDAT / Y
FIOPF	Inter Company Posting Procedure	FIOPG / BURKS / GJAHB / BELNR / BUZETI
FIOPP	Accounting Document Header (docs from External System)	SLBKR / BELNR / GJAHB / CLEBK
FRDPA	Run Date of a Program	PRGID
KPFA	Customer / Vendor Linking	WUCL / VGRAS / FWTYP
RDE4	Customer Payment History	KUNNR / BURKS
RDE5	Customer Master Dunning Data	KUNNR / BURKS / MAKER
RDE6	Customer Master Bank Details	KUNNR / BURKS / BANKS / BANEN
RDE7	Customer Master Transaction Figures	KUNNR / BURKS / GJAHB / BANEN
RDE8	Customer Master Special GL Transactions	KUNNR / BURKS / GJAHB / SHBKZ
LITB	Vendor Master Dunning Data	LITRN / BURKS / MAKER
LPBK	Vendor Master Bank Details	LITRN / BURKS / BANKS / BANEN
LF01	Vendor Master Transaction Figures	LITRN / BURKS / GJAHB / BANEN
LF02	Vendor Master Special GL Transactions	LITRN / BURKS / GJAHB / SHBKZ
YBOPF	Document Header for Standard Parking	AUDIE / BURKS / BELNR / GJAHB
FBASCORE	Financial Accounting General Services "Basic"	AUSKE / BURKS / BELNR / GJAHB
FBAL	Customer Master (Country Code)	KUNNR / BURKS
FKAL	Vendor Master (General Section)	LITRN / BURKS
FKAT	G/L Account Master (Chart of Accounts) - SPARS	KT0PL / SANMR
FKAT	G/L Account Master (Chart of Accounts) - SAKMR	KT0PL / SANMR
MANN	Accounts Blocked by Dunning Selection	KOART / LAUTP / KUNNR / BURKS / GJAHB / SHBKZ / SANMR / BURAS
MANN	Dunning Data (Account Entries)	KUNNR / BURKS
FI-GL-GL (FBS)	<b>General Ledger Accounting: Basic Functions- G/L Accounts</b>	
FIAS	G/L Account Master (Chart of Accounts) - SPARS / KT0PL / SANMR / SCHLW	
FIAS1	G/L Account Master (Company Code)	BURKS / SANMR
FI-GL-GL (FBS)	Functions - R/3 Customizing for G/L Accounts	
FIGLREP	Settings for G/L Posting Reports	MANDT
FIGLREP	General Ledger Accounting: Basic Functions - Fast Data Entry	BURKS / SANMR
FIGLREP	General Ledger Accounting: Basic Functions - Fast Data Entry	BURKS / SANMR
KOMU	Account Assignment Templates for G/L	KOMMAM / KMZEE
KOMU	Account Item	
FI-AR-AR (FBD)	<b>Accounts Receivable: Basic Functions - Customers</b>	
FI-AR-AR (FBD)	Accounts Receivable: Basic Functions - Customers	

Fig. 5. Table with inconsistent layout, no gridlines, no text. Left – input image, centre - HTML generated by our algorithm, right – HTML generated by ABBYY FineReader11

V4 DATABASE RELATIONSHIPS														
This document outlines the primary tables in the V4 database. For each table, the intended purpose, the significant key columns and suggested relationship to other tables is provided. The format of the content will mirror the following:														
TABLE NAME (Primary Key Columns)														
Purpose and data items of interest														
Recommended relationship criteria														
PARCEL (SWIS_CO_SWIS_TOWN_SWIS_VG_PARCEL_ID)														
Data items relate to parcel specific data that spans assessment roll years, i.e. print key, grid coordinates, parcel address														
Parcel table is the primary "target" table in most situations and will be referred to in that manner.														
For a list of all print keys within a municipality, parcel table may be queried by SWIS code.														
ASSESSMENT (SWIS_CO_SWIS_TOWN_SWIS_VG_PARCEL_ID, ROLL_YR)														
Data items relate to specific assessment roll year settings, i.e. roll section, property class, assessed values for land and total, taxable values for county, municipality, village and school, active code, school code, active code														
Assessment table is related to parcel via SWIS_CO_SWIS_TOWN_SWIS_VG and PARCEL_ID. Data for a given assessment roll can be gathered by including ROLL_YR in selection criteria. Many of the reporting functions in V4 utilize active code to distinguish parcels as useful, i.e. ACTIVE_CODE where 'A' active, 'D' inactive (deleted), 'R' reactivated, 'H' historical.														
Further, several data items on this table have related reference tables that provide decode information and are listed here in part:														
<table border="1"> <thead> <tr><th>PROP_CLASS</th><th>PROPCCLASS_REF</th></tr> </thead> <tbody> <tr><td>ROLL_SECTION</td><td>ROLLSECTION_STR</td></tr> <tr><td>OWN_CODE</td><td>OWNERCODE_STR</td></tr> <tr><td>TAX_CODE</td><td>TAXCODE_STR</td></tr> <tr><td>SCH_CODE</td><td>SCHOOLREF</td></tr> <tr><td>BANK_CODE</td><td>BANK</td></tr> </tbody> </table>			PROP_CLASS	PROPCCLASS_REF	ROLL_SECTION	ROLLSECTION_STR	OWN_CODE	OWNERCODE_STR	TAX_CODE	TAXCODE_STR	SCH_CODE	SCHOOLREF	BANK_CODE	BANK
PROP_CLASS	PROPCCLASS_REF													
ROLL_SECTION	ROLLSECTION_STR													
OWN_CODE	OWNERCODE_STR													
TAX_CODE	TAXCODE_STR													
SCH_CODE	SCHOOLREF													
BANK_CODE	BANK													
Principle ownership of the parcel for an assessment roll year may be determined by referring to the PRIMARY_OWNER data item and following the link to the OWNER table via OWNER_ID (details to follow in discussion of owner-related data).														
EXEMPT (SWIS_CO_SWIS_TOWN_SWIS_VG_PARCEL_ID, ROLL_YR_EX_CODE, EX_NUM when parcel has multiple of same exemption code, UNTL_NBR when exemption is														
V4 DATABASE RELATIONSHIPS														
This document outlines the primary tables in the V4 database. For each table, the intended purpose, the significant key columns and suggested relationship to other tables is provided. The format of the content will mirror the following:														
TABLE NAME (Primary Key Columns)														
PARCEL (SWIS_CO_SWIS_TOWN_SWIS_VG_PARCEL_ID)														
Data items relate to parcel specific data that spans assessment roll years, i.e. print key, grid coordinates, parcel address														
Parcel table is the primary "target" table in most situations and will be referred to in that manner. For a list of all print keys within a municipality, parcel table may be queried by SWIS code.														
ASSESSMENT (SWIS_CO_SWIS_TOWN_SWIS_VG_PARCEL_ID, ROLL_YR)														
Data items relate to specific assessment roll year settings, i.e. roll section, property class, assessed values for land and total, taxable values for county, municipality, village and school, active code, school code, active code														
Assessment table is related to parcel via SWIS_CO_SWIS_TOWN_SWIS_VG and PARCEL_ID. Data for a given assessment roll can be gathered by including ROLL_YR in selection criteria. Many of the reporting functions in V4 utilize active code to distinguish parcels as useful, i.e. ACTIVE_CODE where 'A' active, 'D' inactive (deleted), 'R' reactivated, 'H' historical.														
Further, several data items on this table have related reference tables that provide decode information and are listed here in part:														
<table border="1"> <thead> <tr><th>PROP_CLASS</th><th>PROPCCLASS_REF</th></tr> <tr><th>ROLL_SECTION</th><th>ROLLSECTIONSTR</th></tr> <tr><th>OWN_CODE</th><th>OWNERCODESTR</th></tr> <tr><th>TAX_CODE</th><th>TAXCODESTR</th></tr> <tr><th>SCH_CODE</th><th>SCHOOLREF</th></tr> <tr><th>BANK_CODE</th><th>BANK</th></tr> </thead> </table>			PROP_CLASS	PROPCCLASS_REF	ROLL_SECTION	ROLLSECTIONSTR	OWN_CODE	OWNERCODESTR	TAX_CODE	TAXCODESTR	SCH_CODE	SCHOOLREF	BANK_CODE	BANK
PROP_CLASS	PROPCCLASS_REF													
ROLL_SECTION	ROLLSECTIONSTR													
OWN_CODE	OWNERCODESTR													
TAX_CODE	TAXCODESTR													
SCH_CODE	SCHOOLREF													
BANK_CODE	BANK													
Principle ownership of the parcel for an assessment roll year may be determined by referring to the PRIMARY_OWNER data item and following the link to the OWNER table via OWNER_ID (details to follow in discussion of owner-related data).														
EXEMPT (SWIS_CO_SWIS_TOWN_SWIS_VG_PARCEL_ID, ROLL_YR_EX_CODE, EX_NUM when parcel has multiple of same exemption code, UNTL_NBR when exemption is														

Fig. 6. Table with formatted text, no gridlines, with text. Left – input image, centre - HTML generated by our algorithm, right – HTML generated by ABBYY FineReader11

#### IV. FUTURE IMPROVEMENTS

The success of the algorithm relies heavily on the correct recognition of the words by the Tesseract OCR engine. Specifically, some of the problems that lead to failure of the algorithm because of limitations of the OCR engine are:

- Documents with both multiple-column and single-column text. As Tesseract cannot distinguish multiple-column text, these parts are reconstructed as big 2-column tables.
- Documents with non-Manhattan Layout give wrong results as Tesseract is unable to segment them
- Documents with vertical text. Tesseract cannot correctly recognize text that has both vertical and horizontal alignment, leading to wrong results.
- Documents with images. Tesseract tries to recognize images as text, leading to “imaginary” words.

Some other minor problems:

- The algorithm sometimes fails to merge consecutive table lines to a single row, as they do not always meet all the line-merging criteria stated in section 2.
- Some words are merged to a single word segment, despite the fact that they are part of different columns, because the horizontal distance between them is very small (below the  $T_{seg}$  threshold).
- Formatted text sometimes tends to be recognized as a table (e.g. source code aligned to the left and comments aligned to the right)

#### REFERENCES

- [1] T. Kieninger and A. Dengel, “The T-Recs table recognition and analysis system,” DAS ’98, LNCS 1655, pp. 255-270, 1999
- [2] E. Oro and M. Ruffolo, “PDF-TREX: An approach for recognizing and extracting tables from pdf documents,” ICDAR’09, pp. 906-910, 2009