



# Table Identification and Information extraction in Spreadsheets

Elvis Koci<sup>1,2</sup>, Maik Thiele<sup>1</sup>, Oscar Romero<sup>2</sup>, and Wolfgang Lehner<sup>1</sup>

<sup>1</sup>Technische Universität Dresden, Germany

<sup>2</sup>Universitat Politècnica de Catalunya, Spain

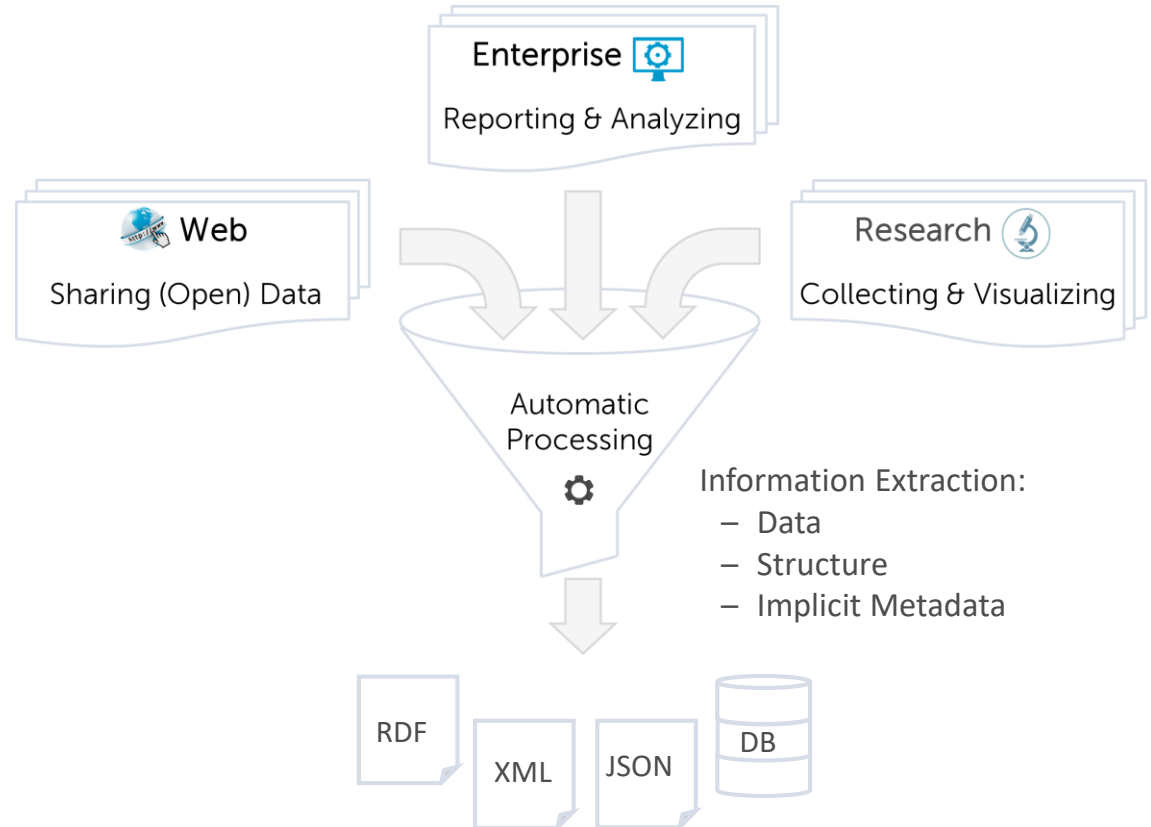
# Spreadsheets in the Realm of Big Data

## IMPORTANT SOURCE OF INFORMATION

- Used in various domains
- Content generated by experts

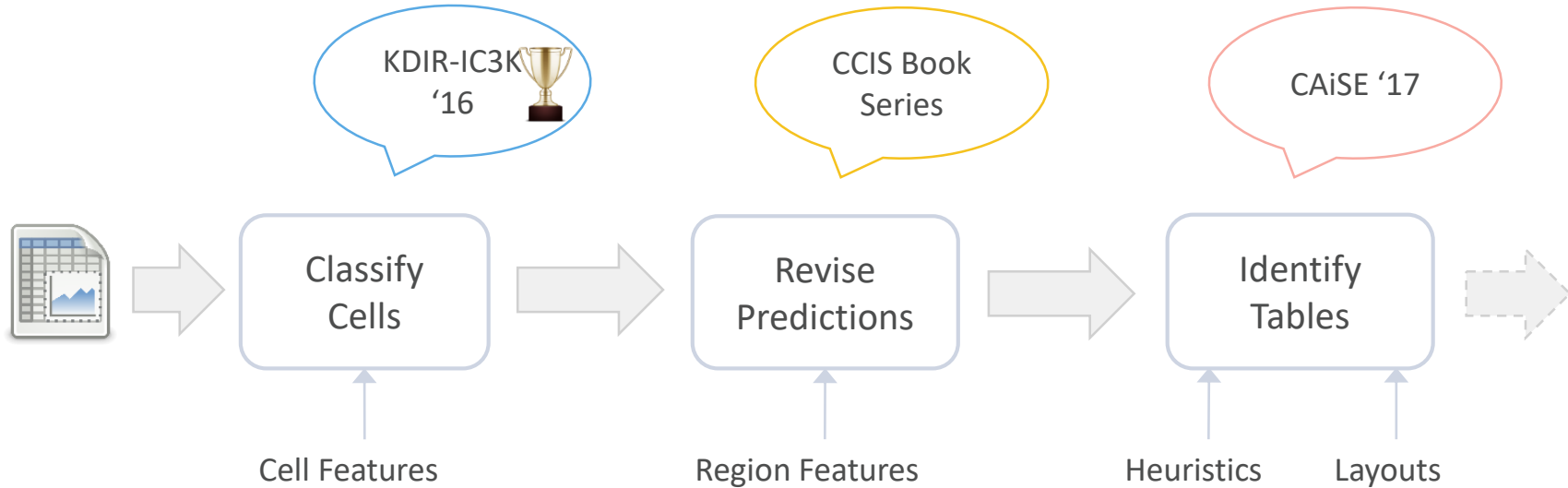
## CHALLENGES

- No description of the structure
- High degree of autonomy
- Implicit Information



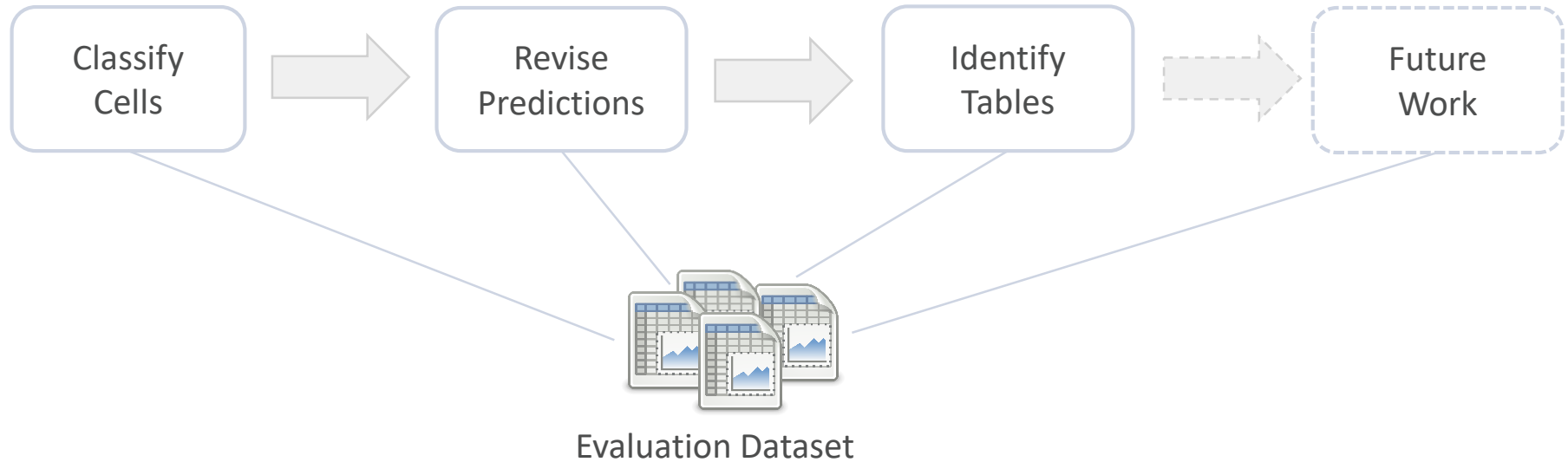
# Proposed Solution

## A PROCESSING PIPELINE



# Proposed Solution

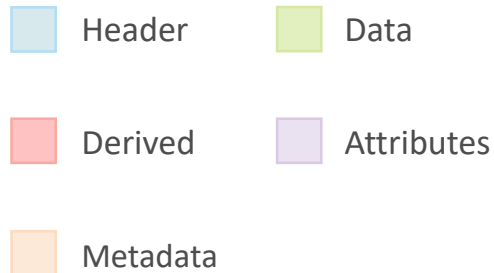
## A PROCESSING PIPELINE



# Cell Classification\*

## UNDERSTANDING THE DATA

- Work at the cell granularity:
  - Omit empty cells
  - Omit floating objects
- Use **5** layout categories (labels):



	A	B	C		D		E	
1		<b>Number of Items (in Units) Sold per Region</b>						
2								
3		Location	Item			Total		
4			Mouse	Monitor	Adapter			
5					VGA		HDMI	
6		<u>Europe</u>						
7		Spain	500	200	85	61	846	
8		France	465	169	80	80	794	
9		<u>Asia</u>						
10		China	422	163	90	44	719	
11		Vietnam	473	182	74	55	784	
12		Item "Keyboard" is omitted. Check next sheet.						

\* Koci et al.: A machine learning approach for layout inference in spreadsheets. KDIR IC3K'16

# Cell Classification\*

## UNDERSTANDING THE DATA

- Work at the cell granularity:
  - Omit empty cells
  - Omit floating objects
- Use **5** layout categories (labels):

 Header

 Data

 Derived

 Attributes

 Metadata

	A	B	C		D		E	
1		<b>Number of Items (in Units) Sold per Region</b>						
2								
3		Location	Item			Total		
4			Mouse	Monitor	Adapter			
5					VGA		HDMI	
6		<u>Europe</u>						
7		Spain	500	200	85	61	846	
8		France	465	169	80	80	794	
9		<u>Asia</u>						
10		China	422	163	90	44	719	
11		Vietnam	473	182	74	55	784	
12		Item "Keyboard" is omitted. Check next sheet.						

\* Koci et al.: A machine learning approach for layout inference in spreadsheets. KDIR IC3K'16

# Cell Classification\*

## UNDERSTANDING THE DATA

- Work at the cell granularity:
  - Omit empty cells
  - Omit floating objects
- Use **5** layout categories (labels):

 Header

 Data

 Derived

 Attributes

 Metadata

	A	B	C		D		E	
1		Number of Items (in Units) Sold per Region						
2								
3		Location	Item			Total		
4			Mouse	Monitor	Adapter			
5					VGA		HDMI	
6		<u>Europe</u>						
7		Spain	500	200	85	61	846	
8		France	465	169	80	80	794	
9		<u>Asia</u>						
10		China	422	163	90	44	719	
11		Vietnam	473	182	74	55	784	
12		Item "Keyboard" is omitted. Check next sheet.						

\* Koci et al.: A machine learning approach for layout inference in spreadsheets. KDIR IC3K'16

# Cell Classification\*

## UNDERSTANDING THE DATA

- Work at the cell granularity:
  - Omit empty cells
  - Omit floating objects
- Use **5** layout categories (labels):

Header

Data

Derived

Attributes

Metadata

	A	B	C	D	E		
1		<b>Number of Items (in Units) Sold per Region</b>					
2							
3		Location	Item			Total	
4			Mouse	Monitor	Adapter		
5					VGA		
6		<u>Europe</u>					
7		Spain	500	200	85	61	846
8		France	465	169	80	80	794
9		<u>Asia</u>					
10		China	422	163	90	44	719
11		Vietnam	473	182	74	55	784
12		Item "Keyboard" is omitted. Check next sheet.					

\* Koci et al.: A machine learning approach for layout inference in spreadsheets. KDIR IC3K'16



# Cell Classification\*

## UNDERSTANDING THE DATA

- Work at the cell granularity:
  - Omit empty cells
  - Omit floating objects
- Use **5** layout categories (labels):

Header

Data

Derived

Attributes

Metadata

	A	B	C		D		E
1		Number of Items (in Units) Sold per Region					
2		Number of Items (in Units) Sold per Region					
3		Location	Item			Total	
4			Mouse	Monitor	Adapter		
5					VGA		HDMI
6		Europe					
7		Spain	500	200	85	61	846
8		France	465	169	80	80	794
9		Asia					
10		China	422	163	90	44	719
11		Vietnam	473	182	74	55	784
12		Item "Keyboard" is omitted. Check next sheet.					

\* Koci et al.: A machine learning approach for layout inference in spreadsheets. KDIR IC3K'16

# Cell Classification\*

## UNDERSTANDING THE DATA

- Work at the cell granularity:
  - Omit empty cells
  - Omit floating objects
- Use **5** layout categories (labels):

Header

Data

Derived

Attributes

Metadata

	A	B	C		D		E
1		Number of Items (in Units) Sold per Region					
2		Number of Items (in Units) Sold per Region					
3		Location	Item			Total	
4			Mouse	Monitor	Adapter		
5					VGA		HDMI
6		Europe					
7		Spain	500	200	85	61	846
8		France	465	169	80	80	794
9		Asia					
10		China	422	163	90	44	719
11		Vietnam	473	182	74	55	784
12		Item "Keyboard" is omitted. Check next sheet.					

\* Koci et al.: A machine learning approach for layout inference in spreadsheets. KDIR IC3K'16

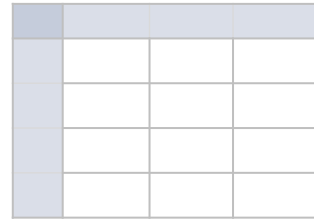
# Supervised Classifications of Cells

## COLLECTING FEATURES

- Binary and numerical >200
- Selected 43 for evaluation

## CLASSIFICATION

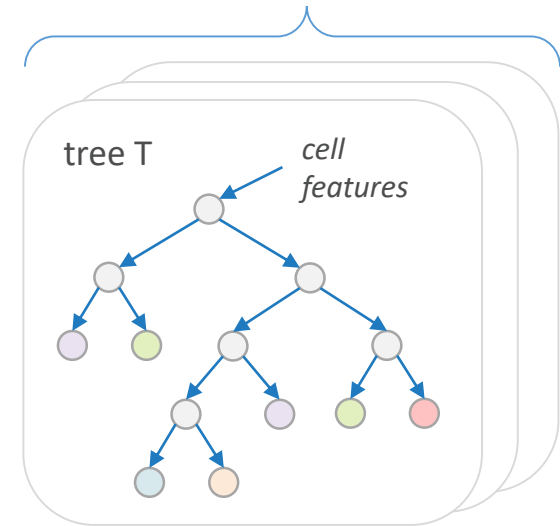
- Evaluated 5 classifiers
- Random Forest most accurate



### Features

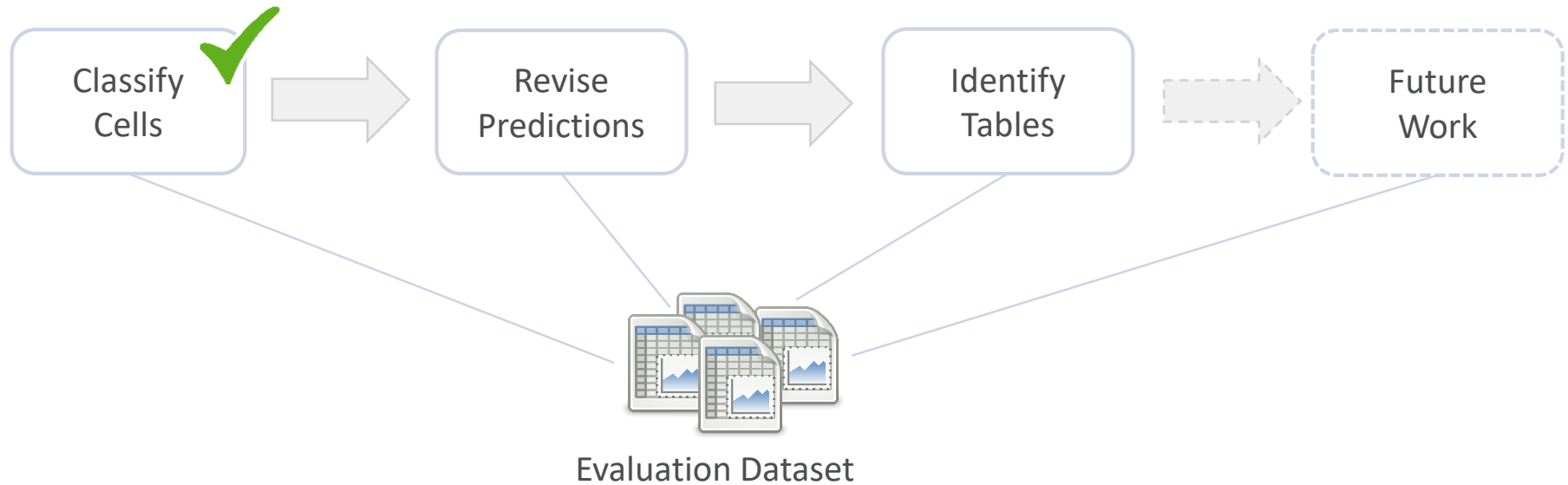
- |                             |                            |
|-----------------------------|----------------------------|
| – <i>Type String</i>        | – <i>Thick Left Border</i> |
| – <i>Is Bold Font</i>       | – <i>Not Aggregation</i>   |
| – <i>No Underline</i>       | – <i>Not Referenced</i>    |
| – <i>Center Aligned</i>     | – <i>Right is Empty</i>    |
| – <i>Default Fill Color</i> | – <i>Top Same Style</i>    |
| ...                         | ...                        |

## Random Forest



Best Results for Data  
and Header

## A PROCESSING PIPELINE

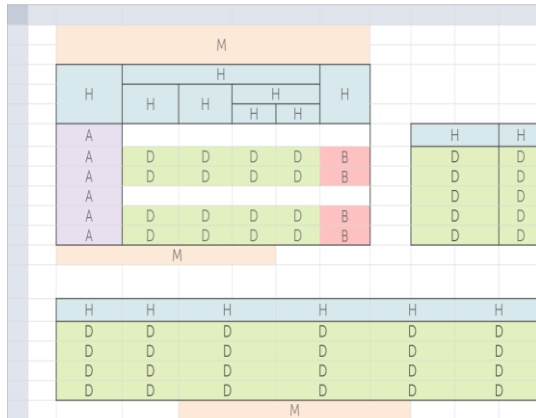


# TIRS Framework\*

## Table identification and Reconstruction

### HEURISTIC PROCEDURE

- Identify tables from the results of cell classification
- A procedure of **8** heuristics tasks



A grid of cells representing classified data. The cells are colored and labeled with letters: A (purple), B (red), D (green), H (blue), and M (orange). The grid is organized into several distinct regions, some of which are enclosed in dashed boxes to represent tables.

Worksheet of Classified Cells



Pre-  
Process



TIRS

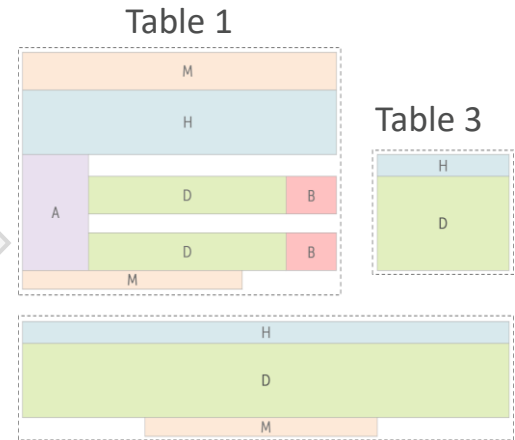


Table 2

\* Koci et al.: Table Identification and Reconstruction in Spreadsheets. CAiSE '17

# Grouping Cells

## NOTIONS

- Work with larger structures
- Decrease complexity

## LABEL REGIONS (LRs)

- Adjacent Cells
- Same Label

M					
H	H				H
	H	H	H		
H			H		
A					
A	D	D	D	D	B
A	D	D	D	D	B
A					
A	D	D	D	D	B
A	D	D	D	D	B
M					

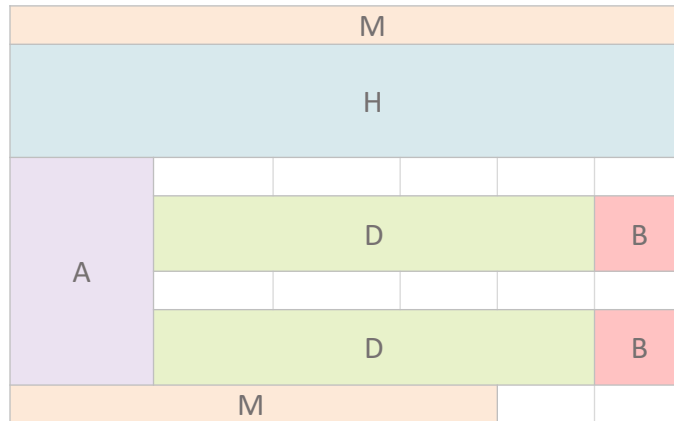
# Grouping Cells

## NOTIONS

- Work with larger structures
- Decrease complexity

## LABEL REGIONS (LRs)

- Adjacent Cells
- Same Label



# Grouping Cells

## NOTIONS

- Work with larger structures
- Decrease complexity

## LABEL REGIONS (LRs)

- Adjacent Cells
- Same Label

M					
H	H				H
	H	H	H		
			H	H	
A					
A	D	D	D	D	B
A	A!	D	D	D	B
A					
D!	D	D	D	D	B
A	D	D	D	D	B
M					



# Grouping Cells

## NOTIONS

- Work with larger structures
- Decrease complexity

## LABEL REGIONS (LRs)

- Adjacent Cells
- Same Label
- Rectilinear polygons

M					
H	H				H
	H	H	H		
			H	H	
A					
A	D	D	D	D	B
A	A!	D	D	D	B
A					
D!	D	D	D	D	B
A	D	D	D	D	B
M					

# Grouping Cells

## NOTIONS

- Work with larger structures
- Decrease complexity

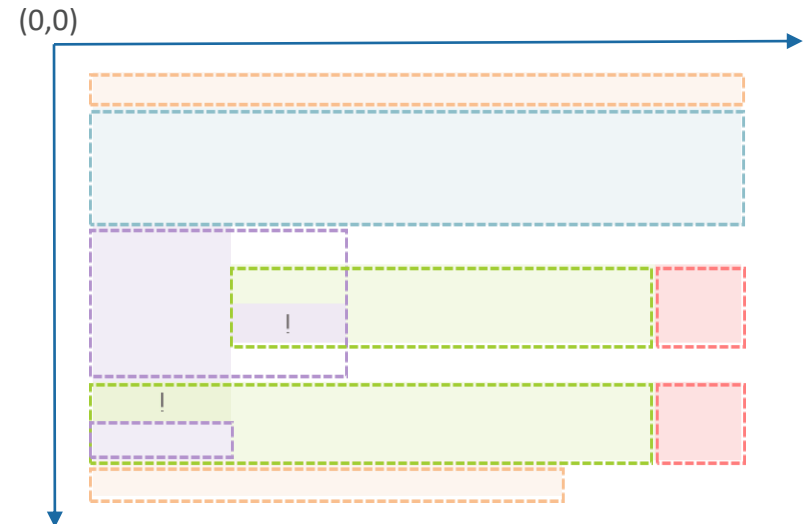
## LABEL REGIONS

- Adjacent Cells
- Same Label
- Rectilinear polygons

## MINIMUM BOUNDING RECTANGLE

- An approximation for LR
- Extensively used for spatial problems

M					
H	H				H
	H	H	H		
A			H	H	
A	D	D	D	D	B
A	A!	D	D	D	B
A					
D!	D	D	D	D	B
A	D	D	D	D	B
M					



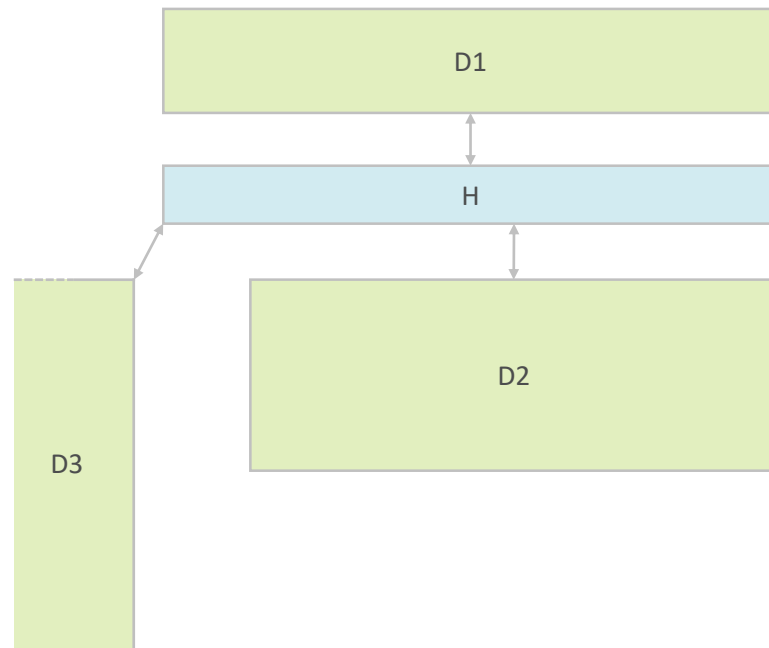
# Constructing Tables - Preliminaries

## CONSTRAINTS FOR PAIRING HEADER WITH DATA

- 1) A Header is on the top of Data
- 2) The distance between them is the smallest
- 3) Their projections overlap significantly:

$$\frac{\text{overlap}(\text{projection}(D), \text{projection}(F))}{\max(\text{projection\_length}(D), \text{projection\_length}(F))} > \theta$$

- Overlap ratio measures this significance
- For the evaluation  $\vartheta$  (empirically) set to **0.5**



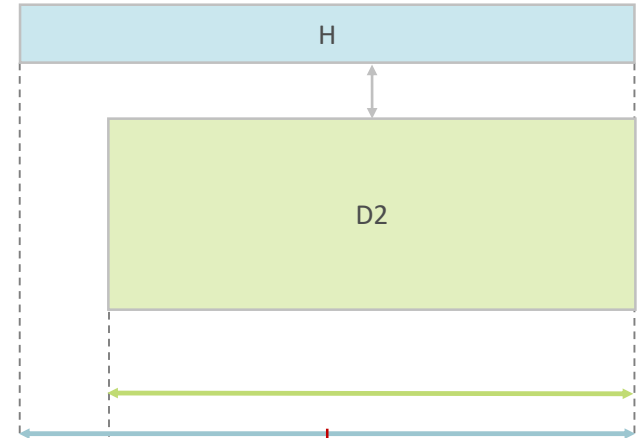
# Constructing Tables - Preliminaries

## CONSTRAINTS FOR PAIRING HEADER WITH DATA

- 1) A Header is on the top of Data
- 2) The distance between them is the smallest
- 3) Their projections overlap significantly:

$$\frac{\text{overlap}(\text{projection}(D), \text{projection}(F))}{\max(\text{projection\_length}(D), \text{projection\_length}(F))} > \theta$$

- Overlap ratio measures this significance
- For the evaluation  $\vartheta$  (empirically) set to **0.5**

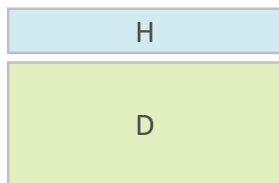


# First 5 Heuristic Steps

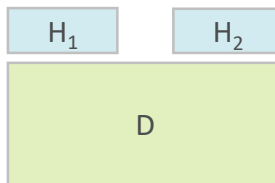
## TABLE PATTERNS

- There are 5 patterns of tables
- Each corresponds to one heuristic task

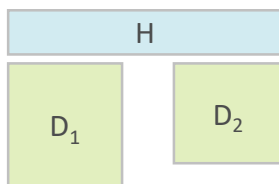
$$\frac{\sum_i^N \sum_j^M \text{overlap}(\text{projection}(H_i), \text{projection}(D_j))}{\max(\sum_i^N \text{projection\_length}(H_i), \sum_j^M \text{projection\_length}(D_j))} > \theta$$



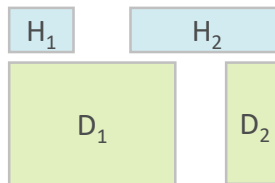
One to One



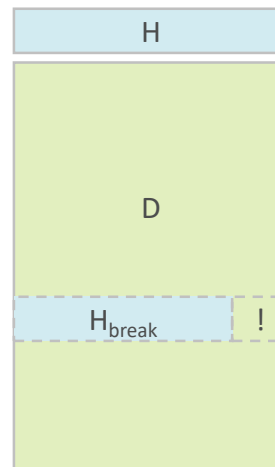
Many to One



One to Many



Many to Many

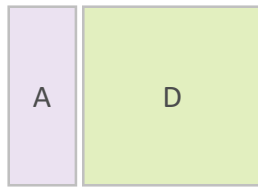


Breaker

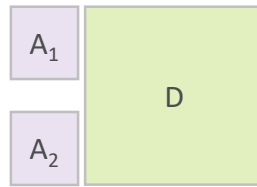
# First 5 Heuristic Steps

## TABLE PATTERNS

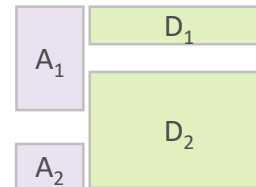
- There are **5** patterns of tables
- Each corresponds to one heuristic task



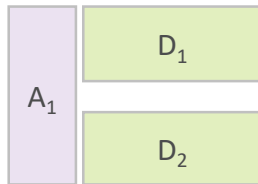
**One to One**



**Many to One**



**Many to Many**



**One to Many**

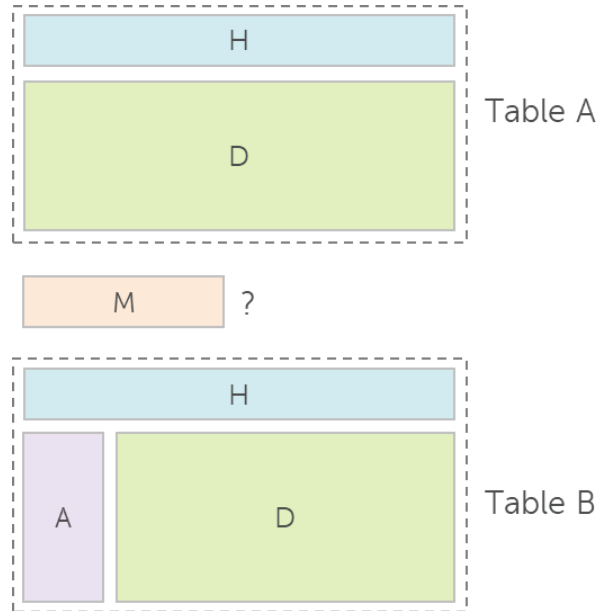


**Breaker**

# Heuristics Steps 6-8

## HANDLE UNPAIRED & REMOVE INCONSISTENCIES

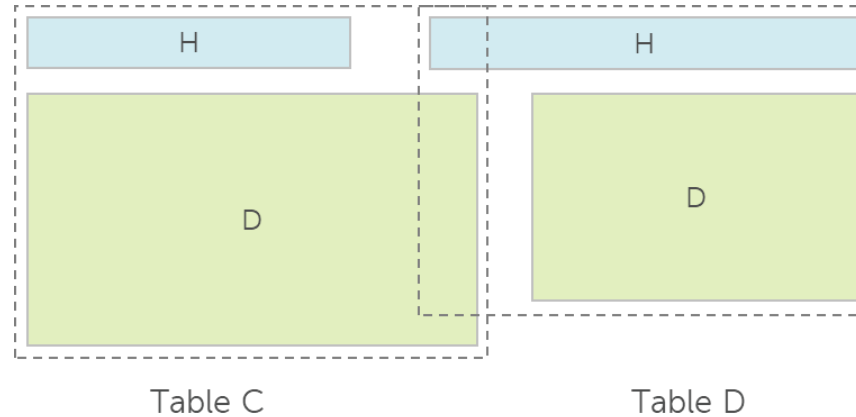
- Assign regions In-between tables



# Heuristics Steps 6-8

## HANDLE UNPAIRED & REMOVE INCONSISTENCIES

- Assign regions In-between tables
- Detect overlapping tables

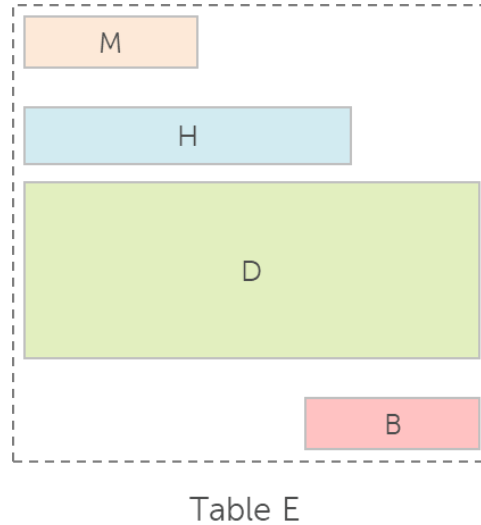




# Heuristics Steps 6-8

## HANDLE UNPAIRED & REMOVE INCONSISTENCIES

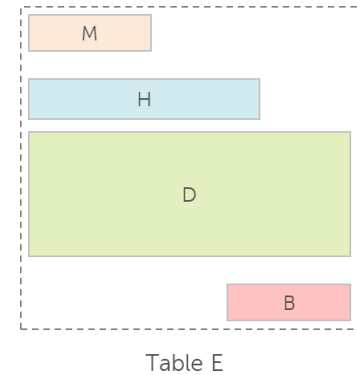
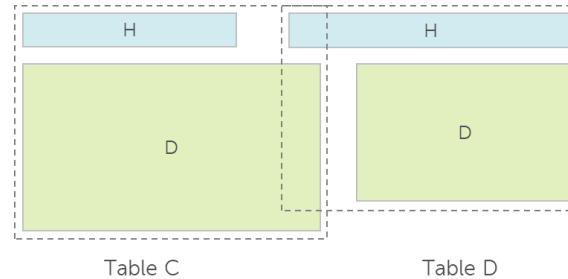
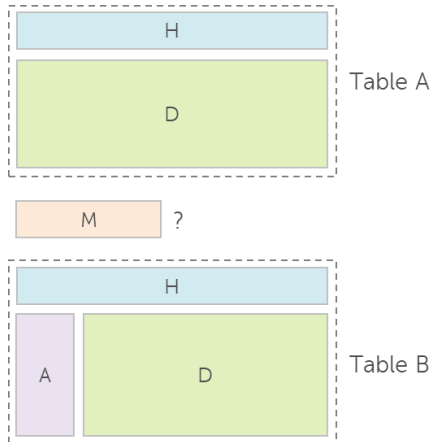
- Assign regions In-between tables
- Detect overlapping tables
- Give regions to nearest table



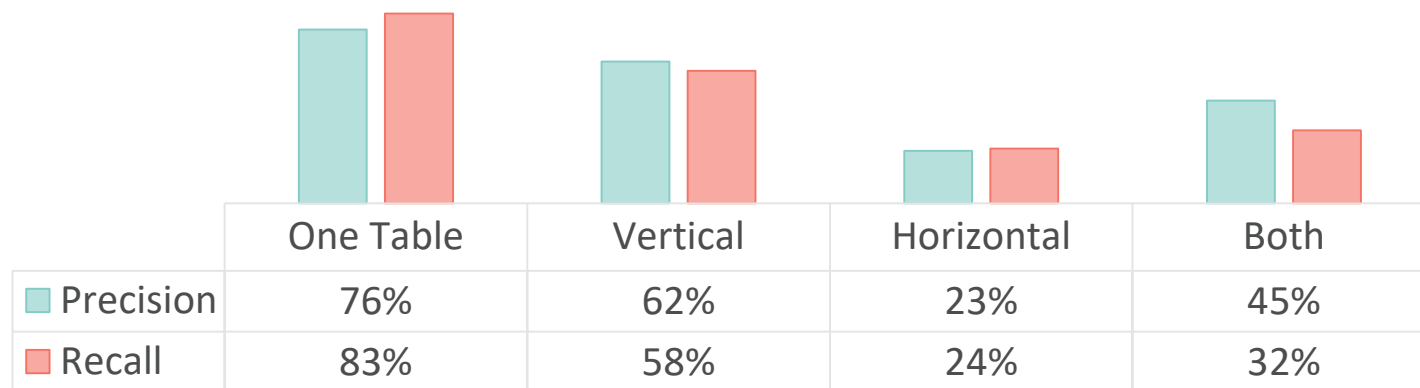
# Heuristics Steps 6-8

## HANDLE UNPAIRED & REMOVE INCONSISTENCIES

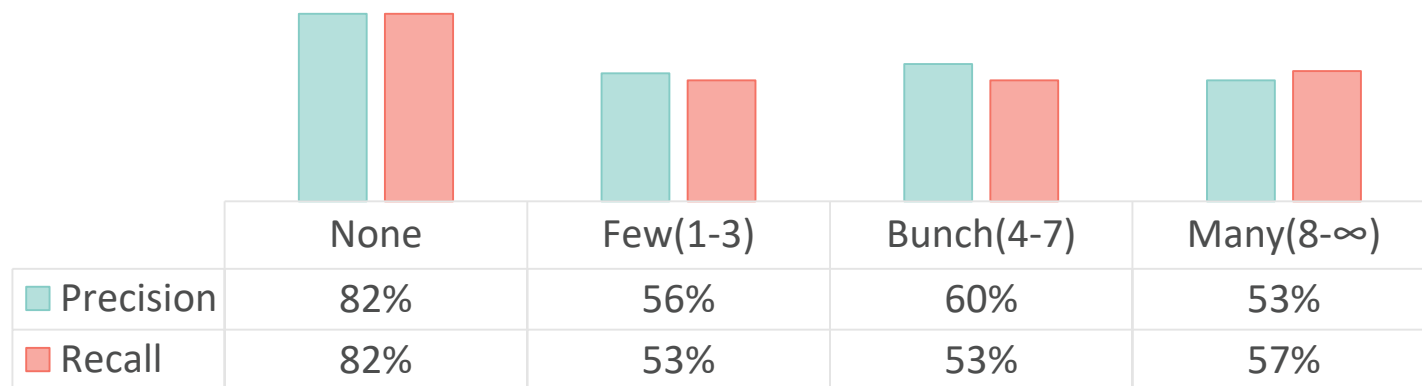
- Assign regions In-between tables
- Detect overlapping tables
- Give regions to nearest table



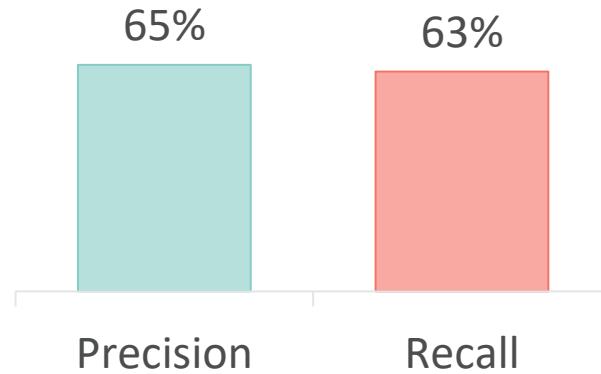
# Evaluation: Table Orientations



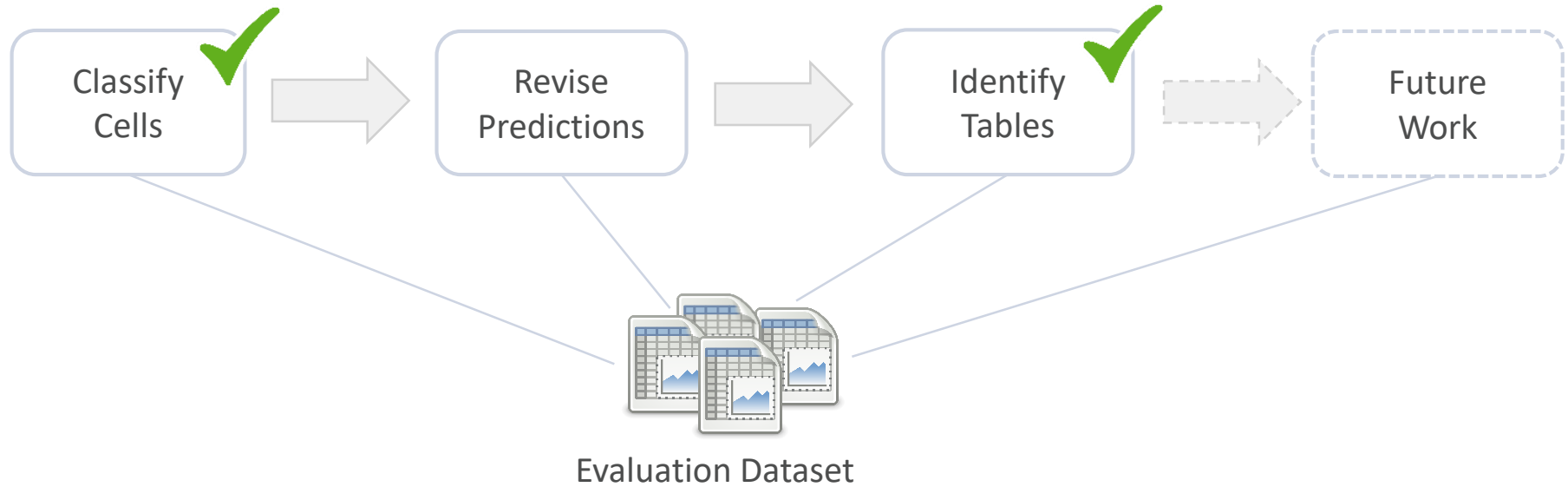
# Evaluation: Number of Misclassified Cells



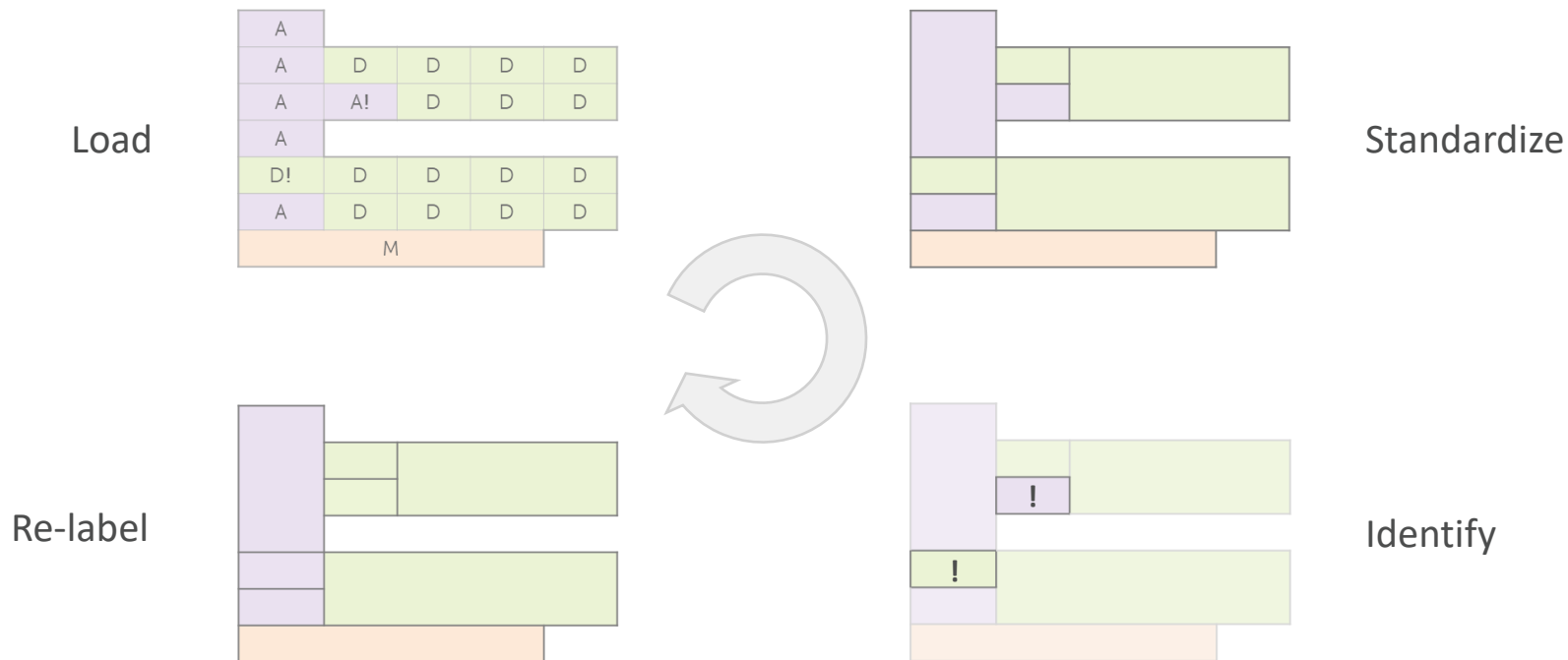
# Overall Results



## A PROCESSING PIPELINE



# Revising Classification Prediction\*



\* Koci et al.: Cell Classification for Layout Recognition in Spreadsheets. Chapter at Communications in Computer and Information Science (CCIS) book series. (Currently In Press)

# Revising Classification Prediction

M					
H	H				H
	H	H	H		
			H	H	
A					
A	D	D	D	D	B
A	A!	D	D	D	B
A					
D!	D	D	D	D	B
A	D	D	D	D	B
M					

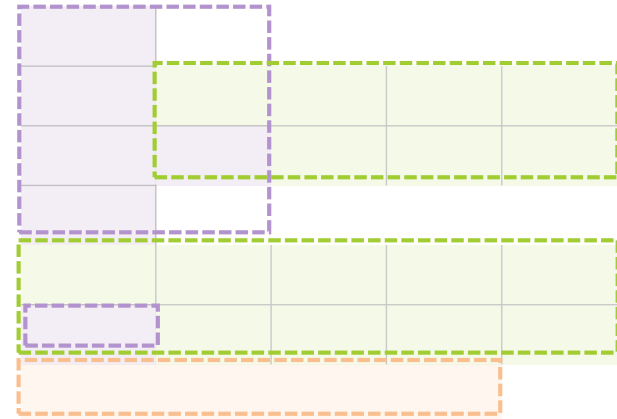
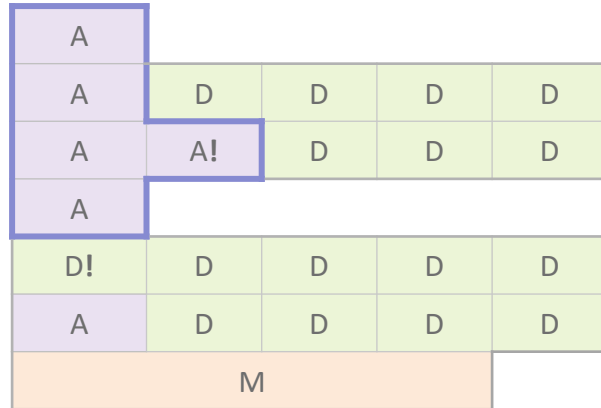


A					
A	D	D	D	D	
A	A!	D	D	D	
A					
D!	D	D	D	D	
A	D	D	D	D	
M					



# Revising Classification Prediction

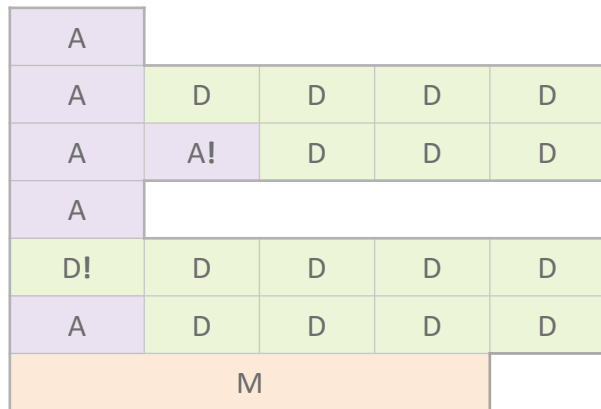
## IRREGULAR REGIONS HINT MISCLASSIFICATIONS



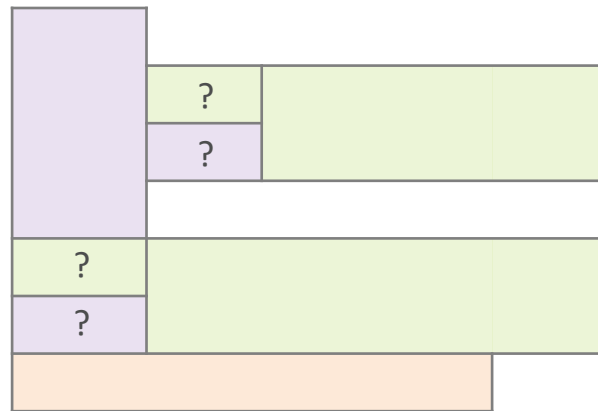
# Revising Classification Prediction

## IRREGULAR REGIONS HINT MISCLASSIFICATIONS

- Make all strictly rectangular
- Small regions are suspicious



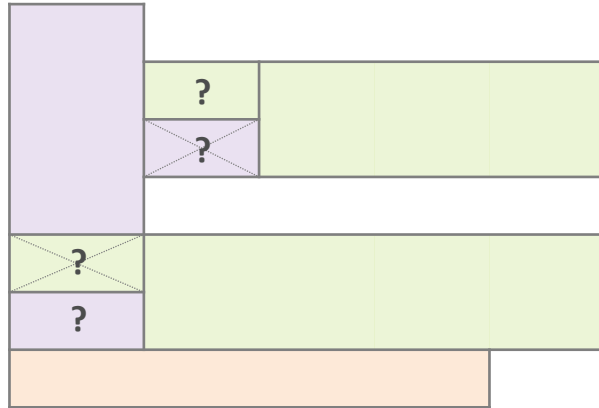
Standardize!



# Revising Classification Prediction

## IDENTIFY TRUE NEGATIVES

- A supervised machine learning task
- Use the context from the neighborhood

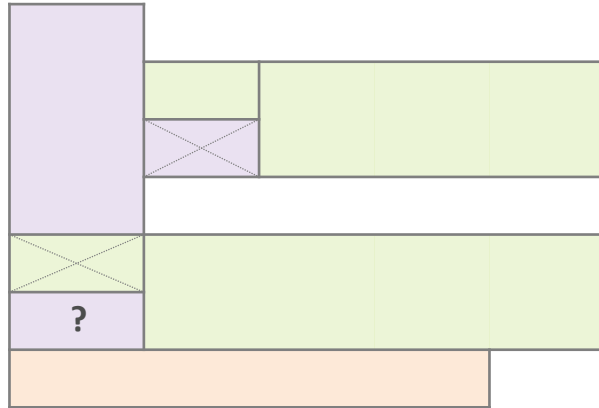


- Similarity
- Dissimilarity

# Revising Classification Prediction

## IDENTIFY TRUE NEGATIVES

- A supervised machine learning task
- Use the context from the neighborhood

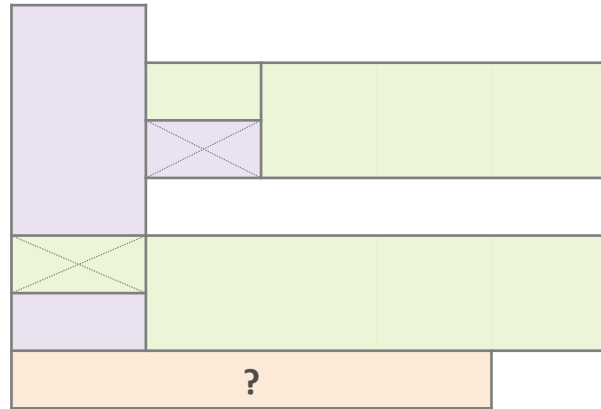


- Similarity
- Dissimilarity
- Influence

# Revising Classification Prediction

## IDENTIFY TRUE NEGATIVES

- A supervised machine learning task
- Use the context from the neighborhood

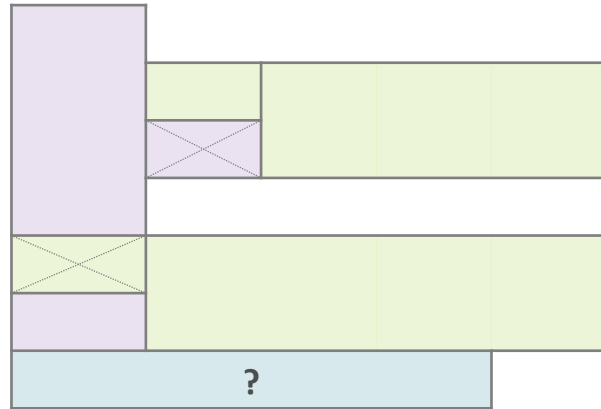


- Similarity
- Dissimilarity
- Influence
- Absence

# Revising Classification Prediction

## IDENTIFY TRUE NEGATIVES

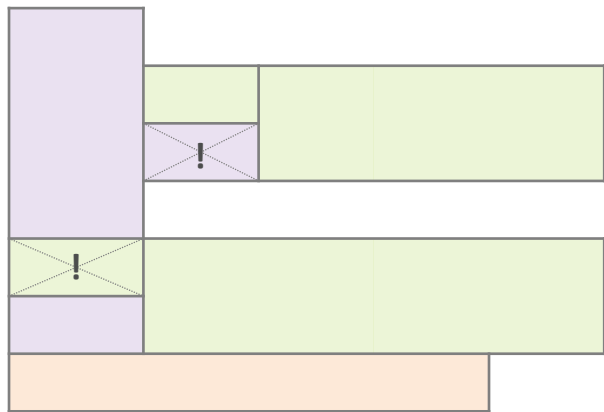
- A supervised machine learning task
- Use the context from the neighborhood



- Similarity
- Dissimilarity
- Influence
- Absence

## IDENTIFY TRUE NEGATIVES

- A supervised machine learning task
- Use the context from the neighborhood



$$weight_i = \text{OverlapRatio}(r, n_i, d) \cdot \frac{1}{1 + \text{Distance}(r, n_i)}$$

$$similarity_i = \begin{cases} 0 & \text{Label}(r) \neq \text{Label}(n_i) \vee n_i \notin \text{Nearest}(r, d) \\ weight_i & \text{otherwise} \end{cases}$$

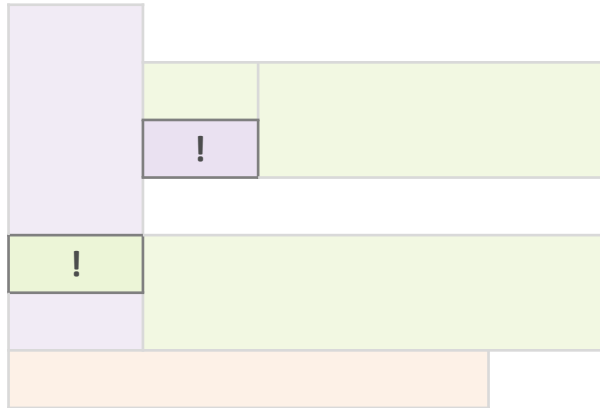
$$dissimilarity_i = \begin{cases} 0 & \text{Label}(r) = \text{Label}(n_i) \vee n_i \notin \text{Nearest}(r, d) \\ weight_i & \text{otherwise} \end{cases}$$

$$influence_i = \begin{cases} 0 & \text{Label}(n_i) \neq l \vee n_i \notin \text{Nearest}(r, d, l) \\ weight_i & \text{otherwise} \end{cases}$$

# Revising Classification Prediction

## PREDICT TRUE LABEL (I.E., RELABELING)

- A supervised machine learning task
- Use the context from the neighborhood



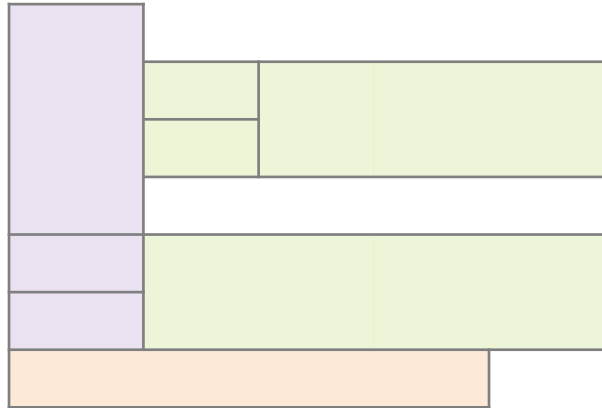
- Similarity
- Dissimilarity
- Influence
- Absence



# Revising Classification Prediction

## PREDICT TRUE LABEL (I.E., RELABELING)

- A supervised machine learning task
- Use the context from the neighborhood

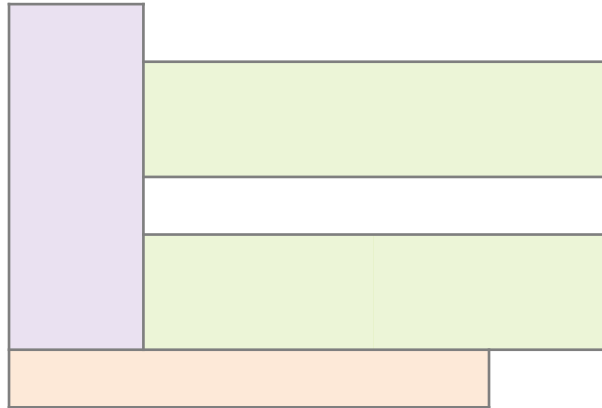


- Similarity
- Dissimilarity
- Influence
- Absence

# Revising Classification Prediction

## PREDICT TRUE LABEL (I.E., RELABELING)

- A supervised machine learning task
- Use the context from the neighborhood



- Similarity
- Dissimilarity
- Influence
- Absence

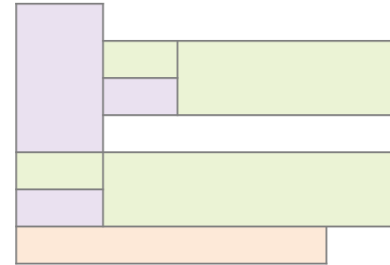
# Revising Classification Prediction

## OVERALL PROCESS

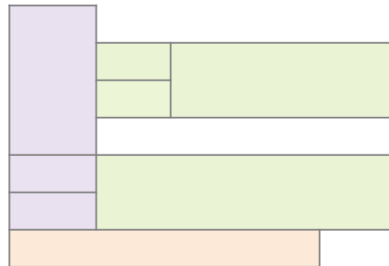
Load

A				
A	D	D	D	D
A	A!	D	D	D
A				
D!	D	D	D	D
A	D	D	D	D
M				

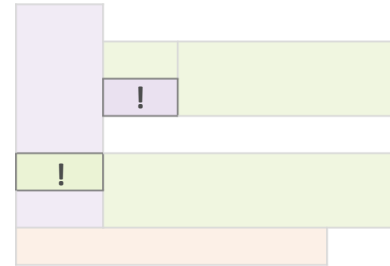
Standardize



Re-label



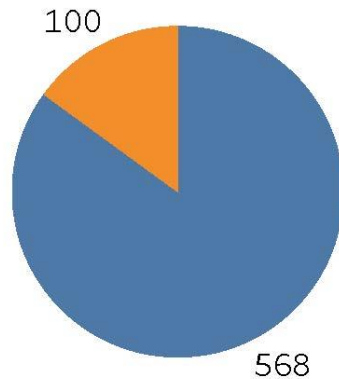
Identify



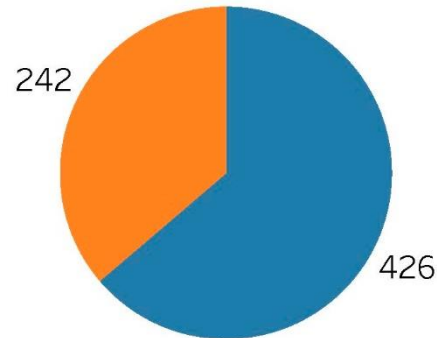
# Revising Classification Prediction

## EVALUATION RESULTS

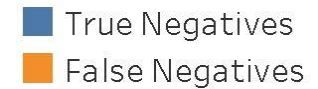
- Successfully identified almost half of misclassified (total **1,237**)
- Relabeling approach needs further improvements



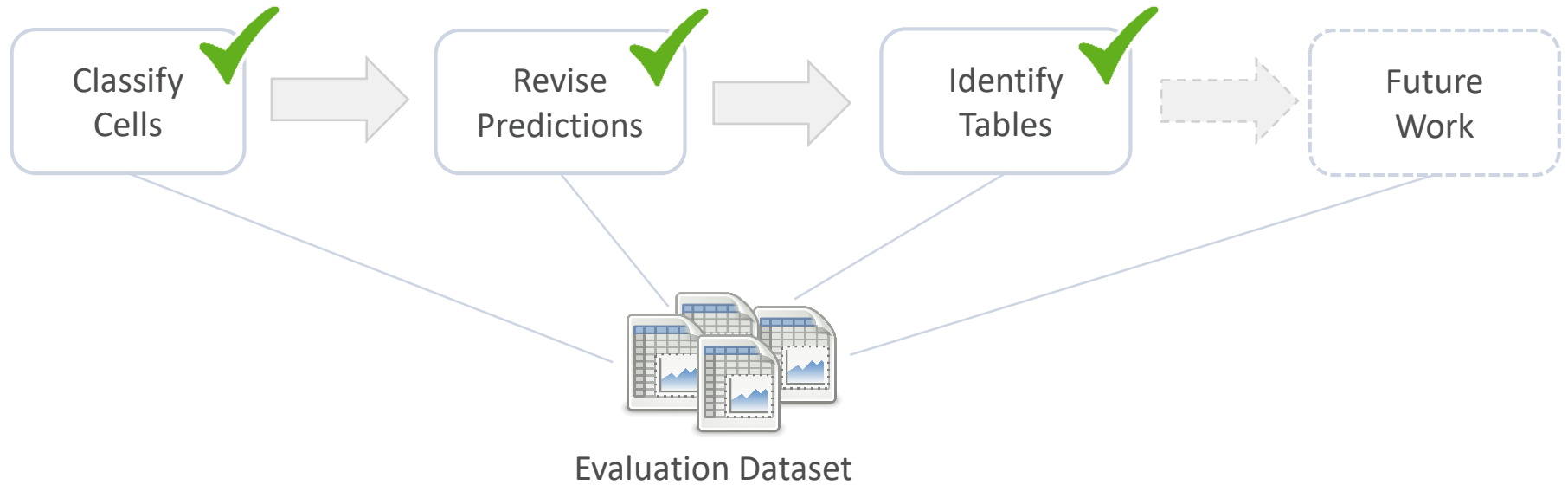
Identifying



Relabeling

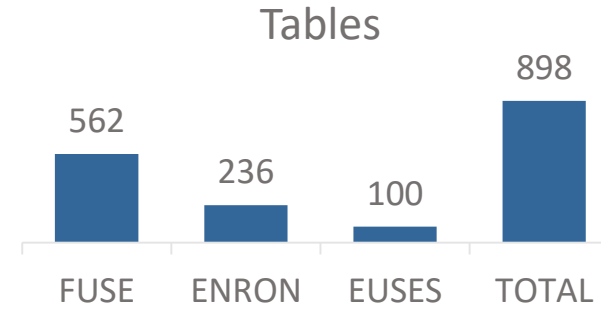
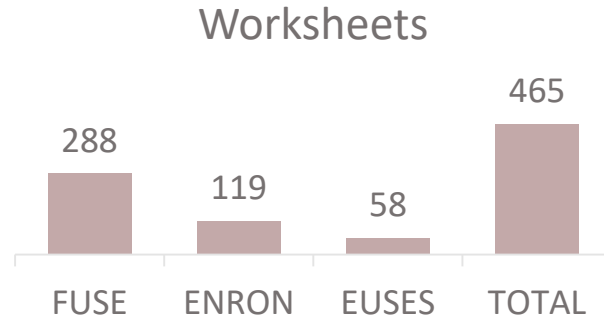


## A PROCESSING PIPELINE



# The Evaluation Dataset

## THREE CORPORA OF SPREADSHEETS



- Barik, T., Lubick, K., Smith, J., Slankas, J., & Murphy-Hill, E. FUSE: A Reproducible, Extendable, Internet-Scale Corpus of Spreadsheets. MSR 2015
- Hermans, F., & Murphy-Hill, E. Enron's Spreadsheets and Related Emails: A Dataset and Analysis. ICSE 2015
- Fisher, M., & Rothermel, G. The EUSES Spreadsheet Corpus: a Shared Resource for Supporting Experimentation With Spreadsheet Dependability Mechanisms. ACM SIGSOFT Software Engineering Notes 2005 (Vol. 30, No. 4, pp. 1-5)

# The Evaluation Dataset

## ANNOTATION TOOL

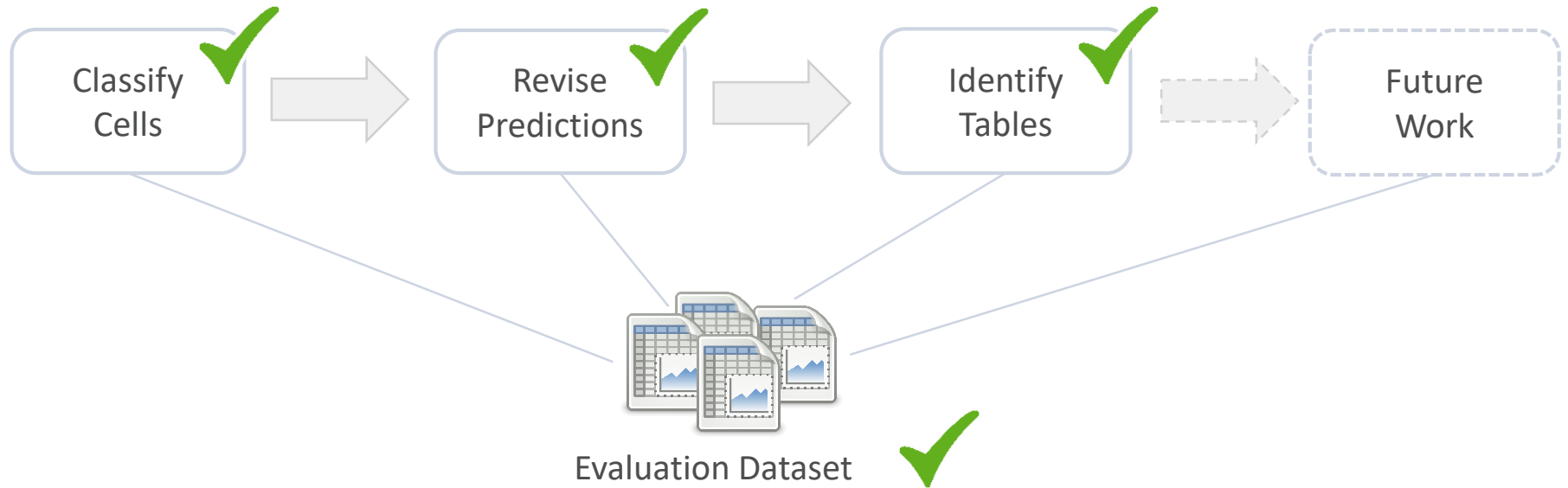
File Annotations Window Preferences

A67

	A	B	C	D	F	G	H	I
	Security/ Ticker	Trade Date	Settlement Date	Instru- ment	Posi- tion	Strike Price	Notional Units	Value
3	<b><u>PUBLICS</u></b>							
4	3TEC Warrants	08-03-00	08-03-03	Swap	Long	\$ 1.18	78,000	\$ 91,937
5	Active Power	08-03-00	08-03-03	Swap	Long	\$ 53.00	1,276,383	\$ 67,648,299
6	Avici Systems	08-03-00	08-03-03	Swap	Long	\$ 162.50	1,093,426	\$ 177,681,725
9	DevX Energy Pref	08-03-00	08-03-03	Swap	Long	\$ 4.07	127,500	\$ 518,400
12	<b><u>PRIVATES</u></b>							
13	Amerada Hess Exposure	08-03-00	08-03-03	Swap	Long			\$ 1,250,000
14	Ameritex	08-03-00	08-03-03	Swap	Long			\$ 4,563,600
51	Quicksilver	08-03-00	08-03-03	Swap	Long	\$ 7.63	484,154	\$ 3,691,676
53	Hughes Rawls Note	08-03-00	08-03-03	Swap	Long			\$ 283,416
54								
56		<b>Totals</b>						\$ 616,559,338
58						Terminations		
59						Public DERIVED		\$ 26,329,435
60						Private		\$ 81,826,449
62						Net Notional		\$ 508,403,454

Amort MPR Raptor Daily Position R ...

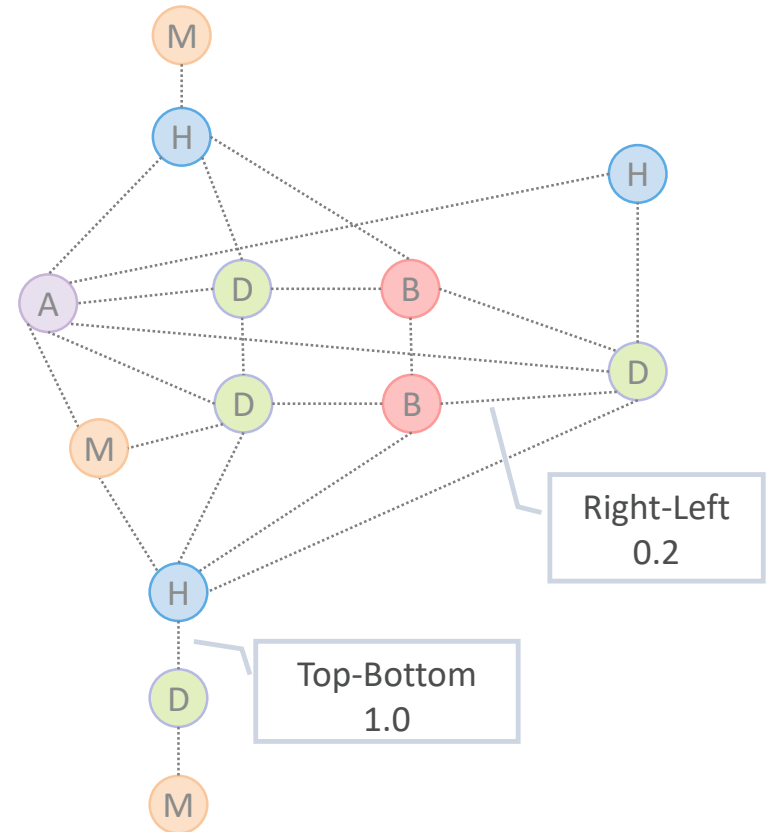
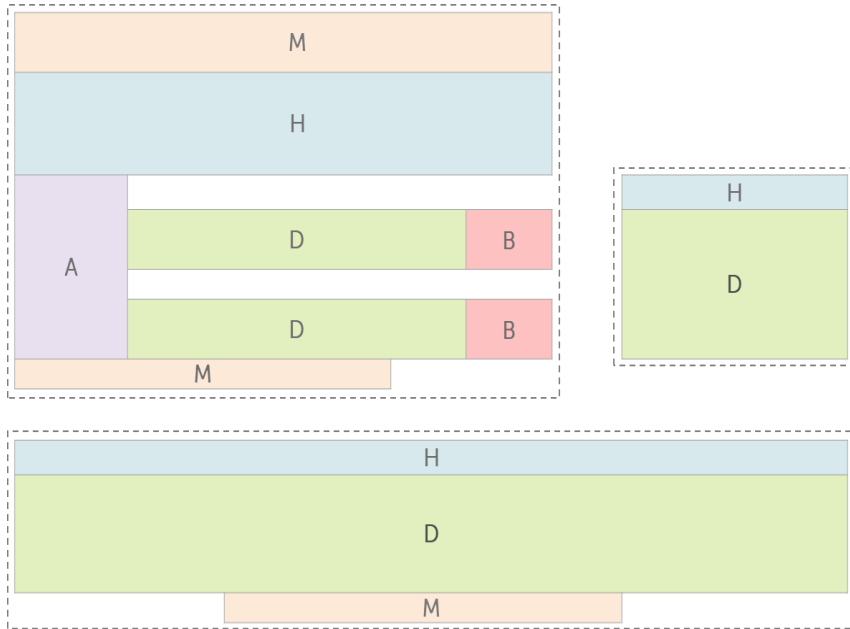
## A PROCESSING PIPELINE





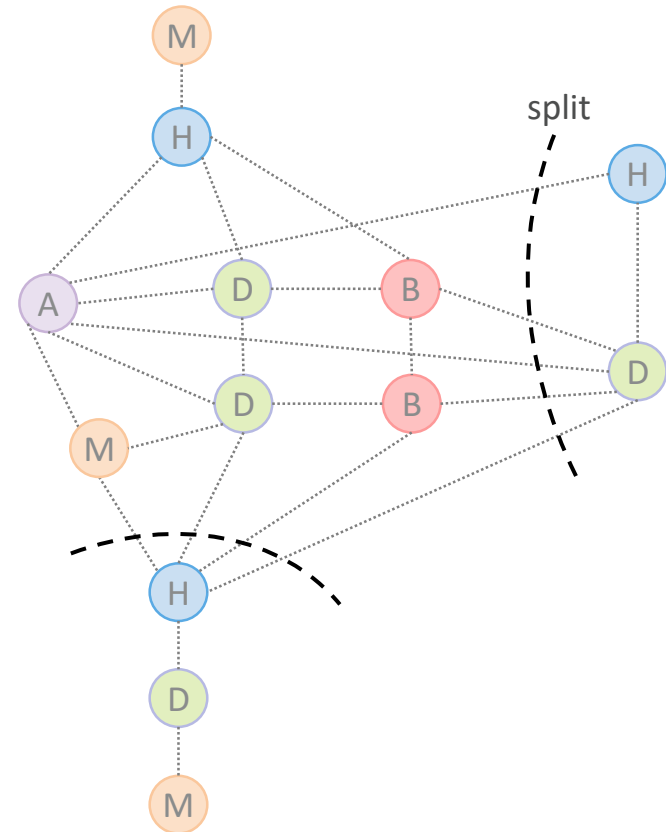
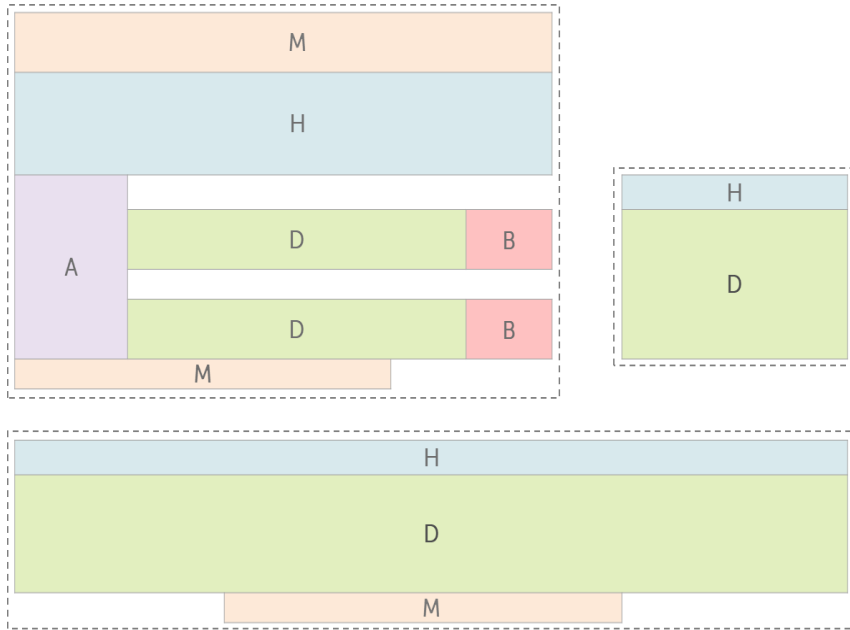
# Future Work

## GRAPH REPRESENTATION



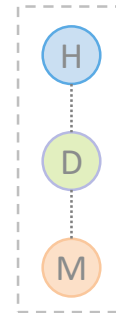
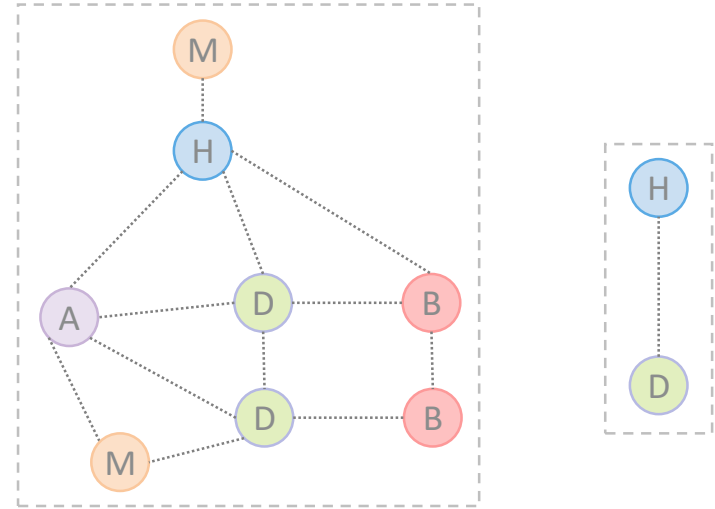
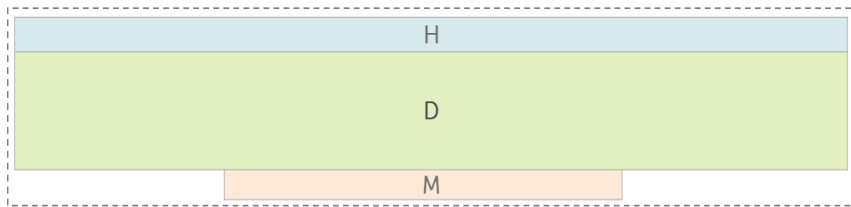
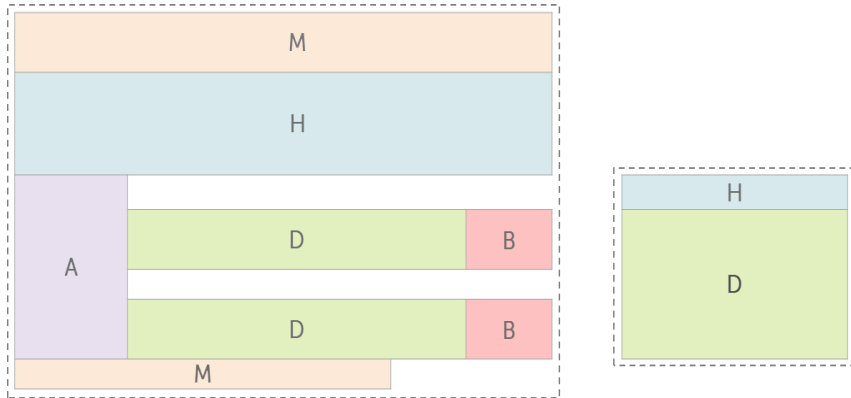
# Future Work

## GRAPH REPRESENTATION



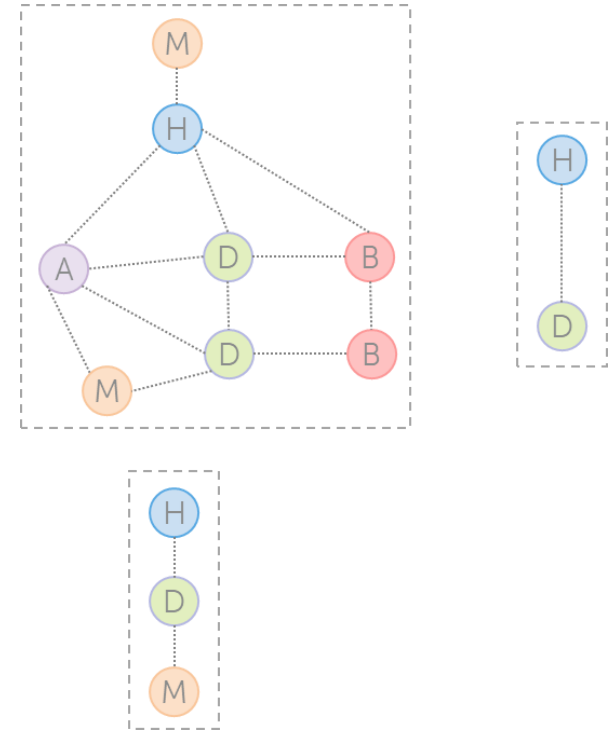
# Future Work

## GRAPH REPRESENTATION



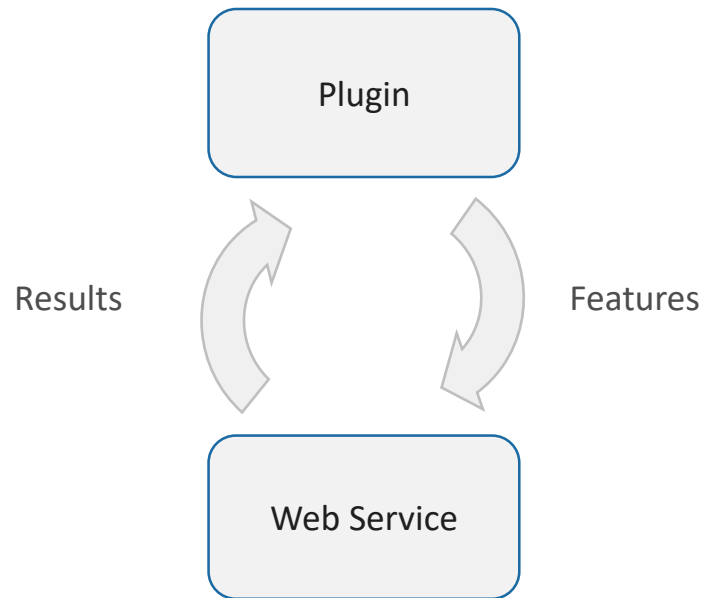
# Future Work

- Graph-based Table Identification (target **ECIR**)
  - Node clustering
  - Guided minimum-cut
  - Fitness function



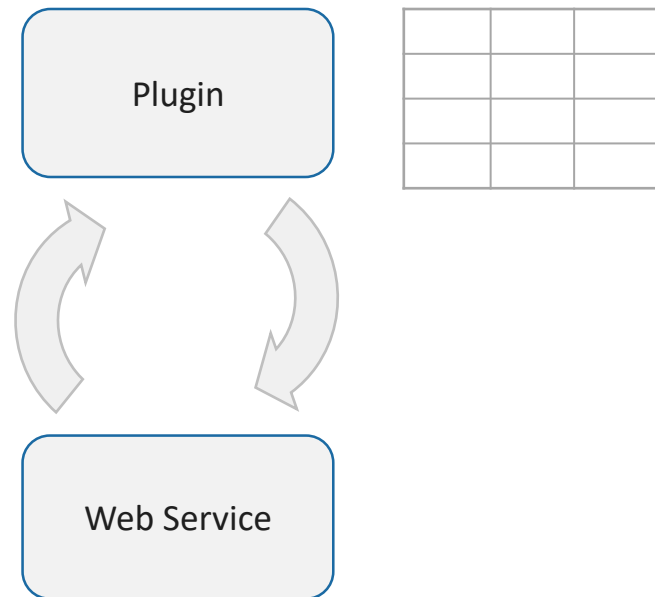
# Future Work

- Graph-based Table Identification (target **ECIR**)
  - Node clustering
  - Guided minimum-cut
  - Fitness function
- Demo: a plugin for Excel to extract information



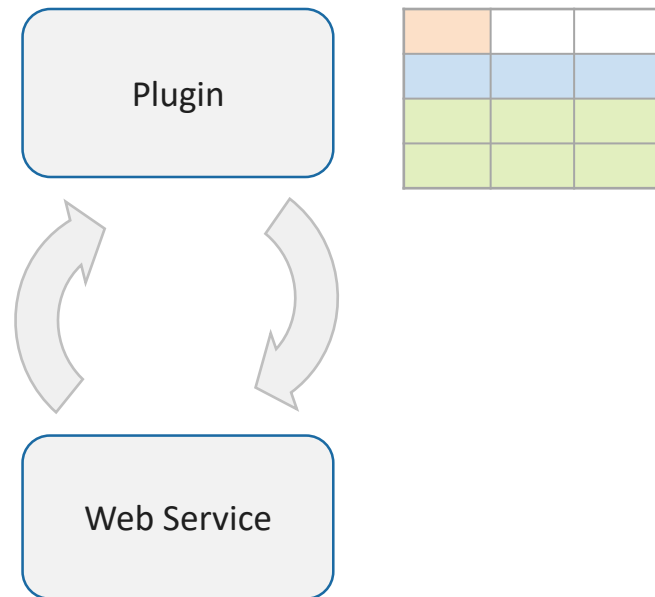
## FUTURE WORK

- Graph-based Table Identification (target **ECIR**)
  - Node clustering
  - Guided minimum-cut
  - Fitness function
- Demo: a plugin for Excel to extract information



## FUTURE WORK

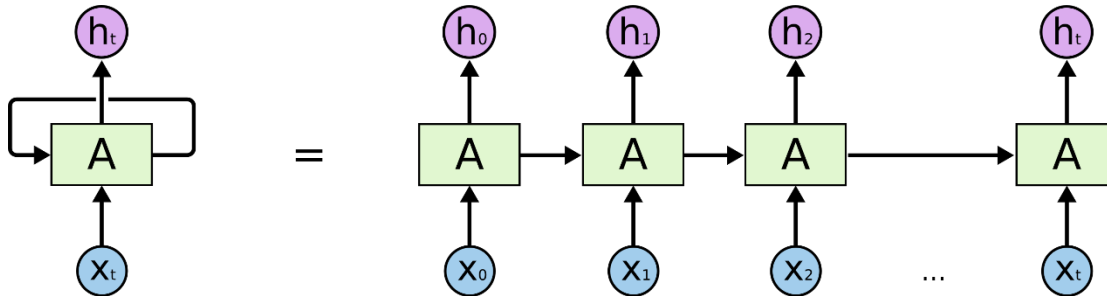
- Graph-based Table Identification (target **ECIR**)
  - Node clustering
  - Guided minimum-cut
  - Fitness function
- Demo: a plugin for Excel to extract information



## FUTURE WORK

- Graph-based Table Identification (target **ECIR**)
  - Node clustering
  - Guided minimum-cut
  - Fitness function
- Demo: a plugin for Excel to extract information
- Cell classification using Recurrent Neural Networks

	A	B	C	D	E		
1		Number of Items (in Units) Sold per Region					
2							
3		Item				Total	
4		Location	Mouse	Monitor	Adapter		
5					VGA		HDMI
6		Europe					
7		Spain	500	200	85	61	846
8		France	465	169	80	80	794
9		Asia					
10		China	422	163	90	44	719
11		Vietnam	473	182	74	55	784
12		Item "Keyboard" is omitted. Check next sheet.					



<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>



# Conclusions & Future Work

## FUTURE WORK

- Graph-based Table Identification (target **ECIR**)
  - Node clustering
  - Guided minimum-cut
  - Fitness function
- Demo: a plugin for Excel to extract information
- Cell classification using Recurrent Neural Networks
- Schema Extraction from identified tables

Table 1

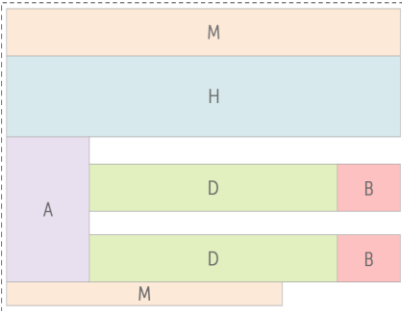


Diagram illustrating Table 1 with cell classification. The table is represented by a grid of colored cells. The top row is orange and labeled 'M'. The second row is light blue and labeled 'H'. The third row consists of a purple cell labeled 'A', a green cell labeled 'D', and a red cell labeled 'B'. The fourth row consists of a green cell labeled 'D' and a red cell labeled 'B'. The bottom row is orange and labeled 'M'. The entire grid is enclosed in a dashed box.

Table 3

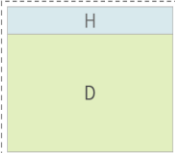


Diagram illustrating Table 3 with cell classification. The table is represented by a grid of colored cells. The top row is light blue and labeled 'H'. The bottom row is green and labeled 'D'. The entire grid is enclosed in a dashed box.

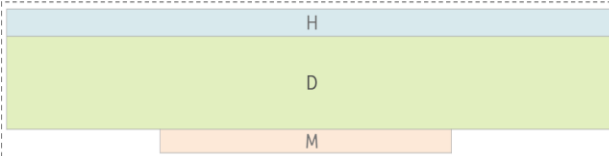


Diagram illustrating Table 2 with cell classification. The table is represented by a grid of colored cells. The top row is light blue and labeled 'H'. The middle row is green and labeled 'D'. The bottom row is orange and labeled 'M'. The entire grid is enclosed in a dashed box.

Table 2

# Conclusions & Future Work

## FUTURE WORK

- Graph-based Table Identification (target **ECIR**)
  - Node clustering
  - Guided minimum-cut
  - Fitness function
- Demo: a plugin for Excel to extract information
- Cell classification using Recurrent Neural Networks
- Schema Extraction from identified tables

Table 1

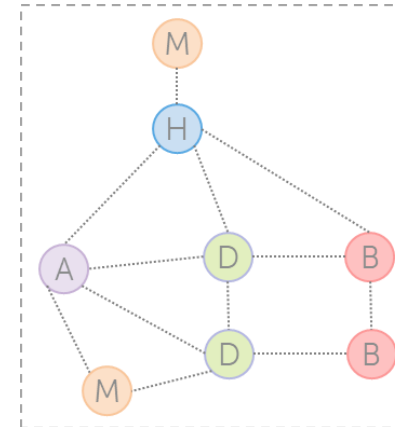
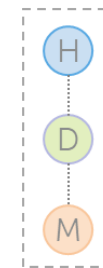


Table 3



Table 2

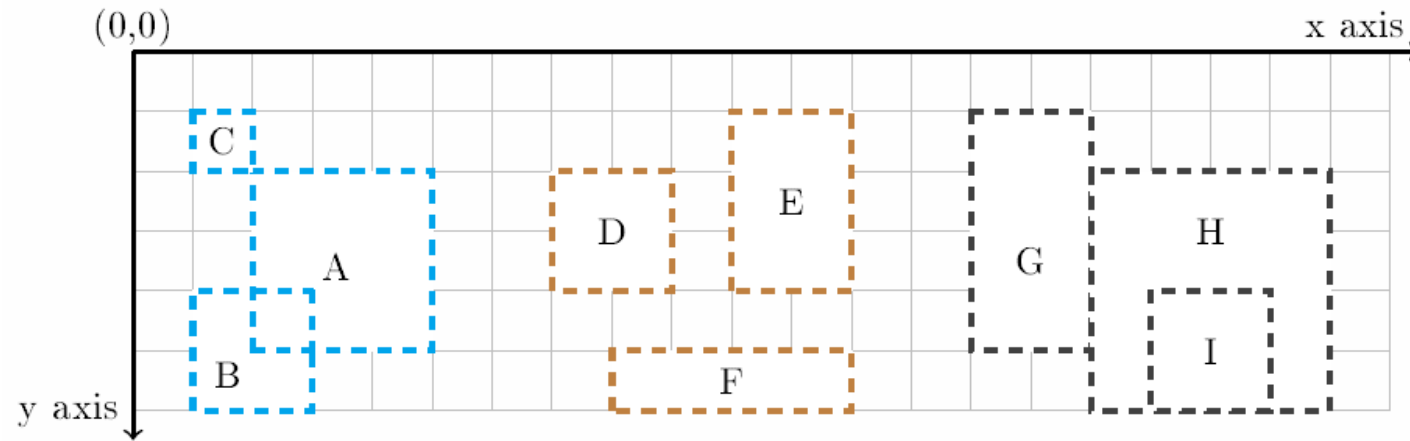


# Related Work

1. Chen, Z., & Cafarella, M. (2013, August). "Automatic Web Spreadsheet Data Extraction". In Proceedings of the 3rd International Workshop on Semantic Search Over the Web (p. 1). ACM.
2. Eberius, J., Werner, C., Thiele, M., Braunschweig, K., Dannecker, L., & Lehner, W. (2013, October). "DeExcelerator: A Framework for Extracting Relational Data from Partially Structured Documents". In Proceedings of the 22nd ACM International Conference on Information & Knowledge Management (pp. 2477-2480). ACM.
3. Adelfio, M. D. and Samet, H. (2013). "Schema Extraction for Tabular Data on the Web". VLDB'13, 6(6): 421–432.
4. Cunha, J., Saraiva, J., & Visser, J. (2009, January). "From Spreadsheets to Relational Databases and Back". In Proceedings of the 2009 ACM SIGPLAN Workshop on Partial Evaluation and Program Manipulation (pp. 179-188). ACM.
5. Abraham, R., & Erwig, M. (2007). "UCheck: A Spreadsheet Type Checker for End Users". Journal of Visual Languages & Computing, 18(1), 71-95.
6. Shigarov, Alexey O. (2015). "Table Understanding Using a Rule Engine". Expert Systems with Applications, 42(2): 929-937.
7. Swidan, A., & Hermans, F. (2017, June). Semi-automatic Extraction of Cross-Table Data from a Set of Spreadsheets. In International Symposium on End User Development (pp. 84-99). Springer, Cham.

Thank You!

# Region Spatial Arrangements

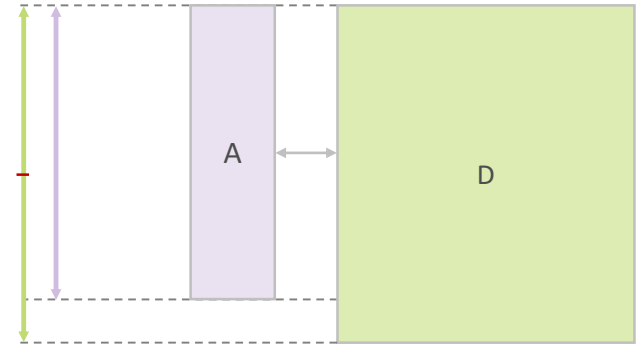


## CONSTRAINS FOR PAIRING ATTRIBUTES WITH DATA

- 1) **Attributes** are on the **left** of Data
- 2) The distance between them is the smallest
- 3) Their projections overlap significantly:

$$\frac{\text{overlap}(\text{projection}(D), \text{projection}(F))}{\max(\text{projection\_length}(D), \text{projection\_length}(F))} > \theta$$

- Overlap ratio measures this significance
- For the evaluation  $\vartheta$  (empirically) set to **0.5**



# Constructing Tables - Preliminaries

## CONSTRAINTS FOR PAIRING ATTRIBUTES WITH DATA

- 1) Attributes are on the left of Data
- 2) The distance between them is the smallest
- 3) Their projections overlap significantly:

$$\frac{\text{overlap}(\text{projection}(D), \text{projection}(F))}{\max(\text{projection\_length}(D), \text{projection\_length}(F))} > \theta$$

- Overlap ratio measures this significance
- For the evaluation  $\vartheta$  (empirically) set to **0.5**

