

# From Partially Structured Documents to Relations

Utilizing and Extending Best Practices in Table Recognition for Automatic Information Extraction in Spreadsheets



Elvis Koci  
elvis.koci@tu-dresden.de

TB2

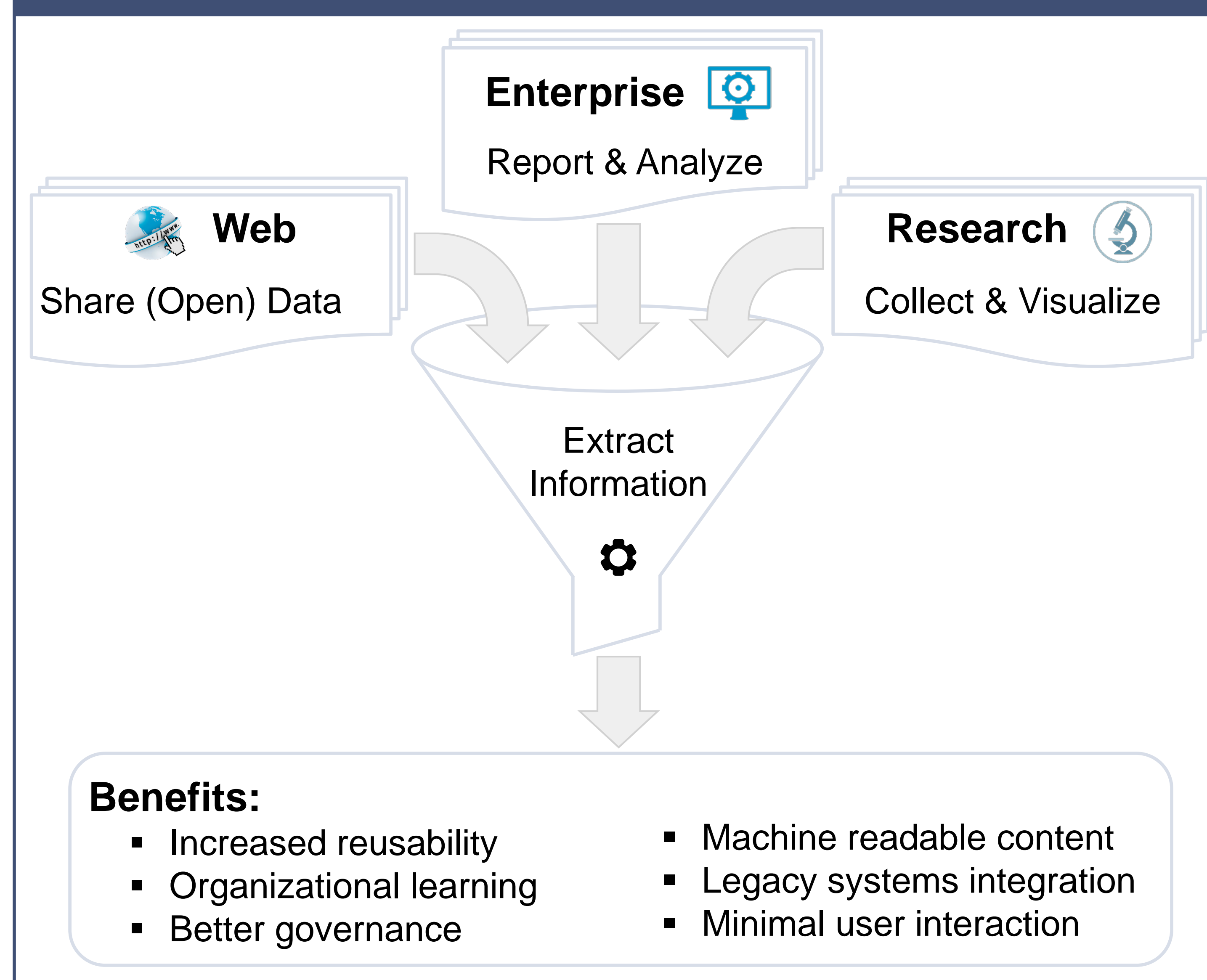
Roles in  
Software Development

2nd  
COHORT

Prof. Lehner  
Prof. Schill

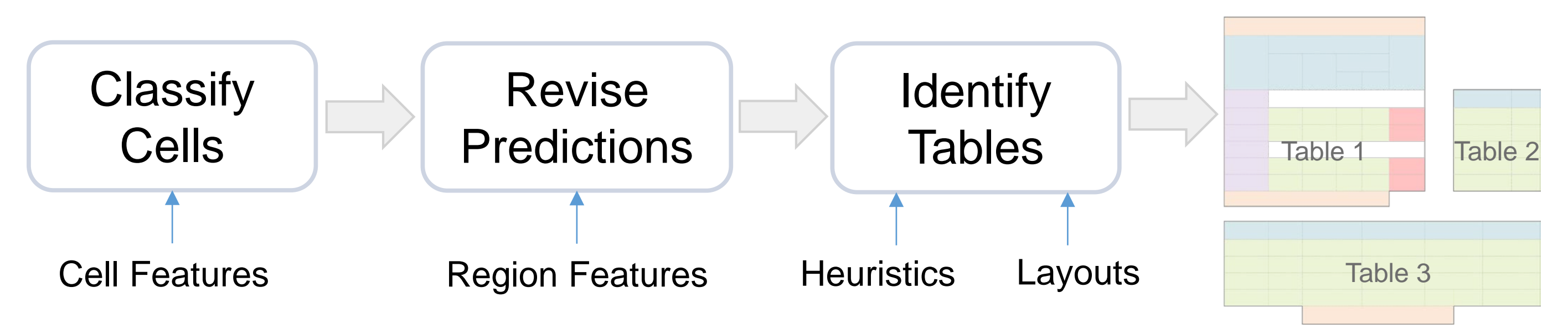
Funded by  
ITBI  
DC

## Motivation

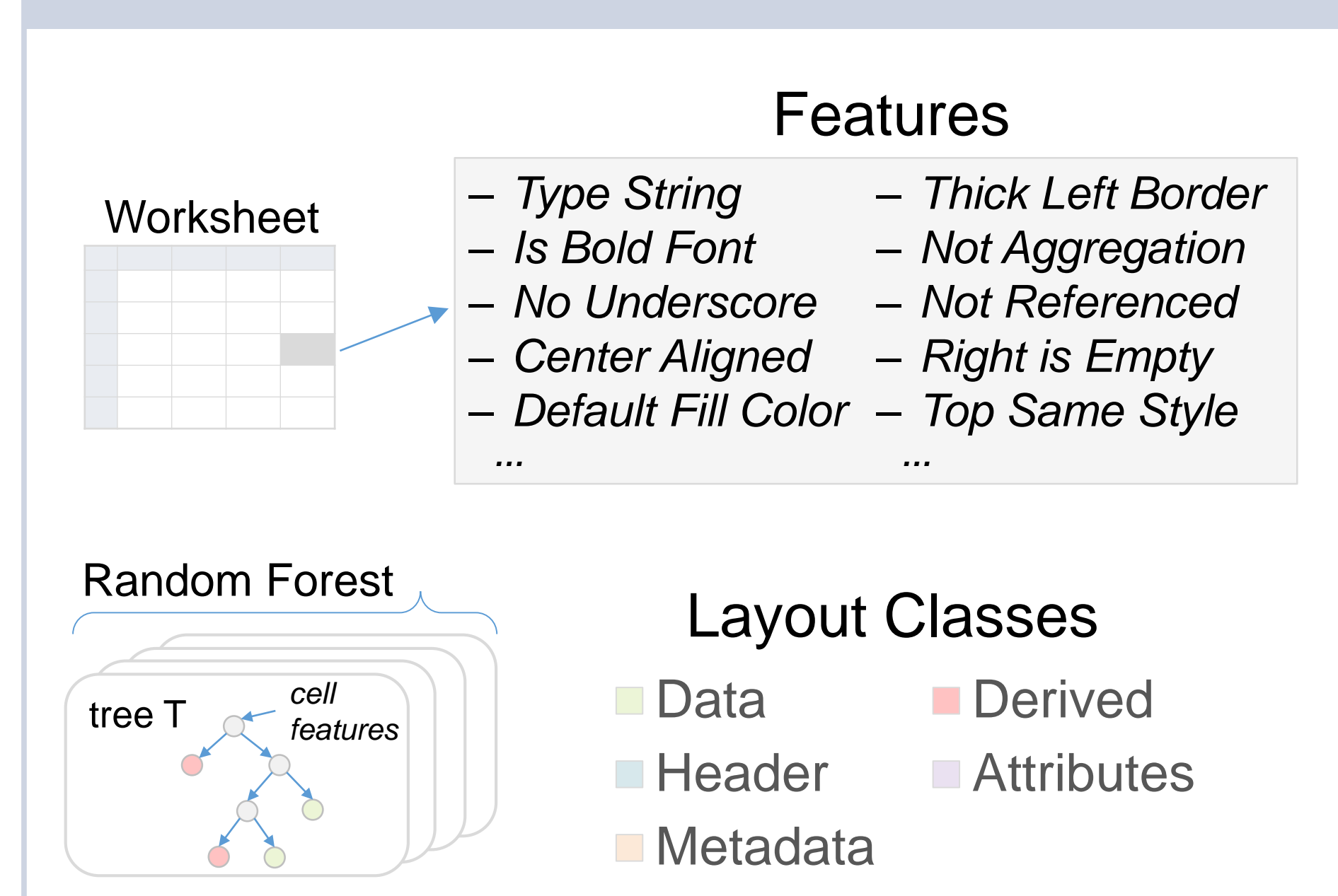


## Proposed Solution

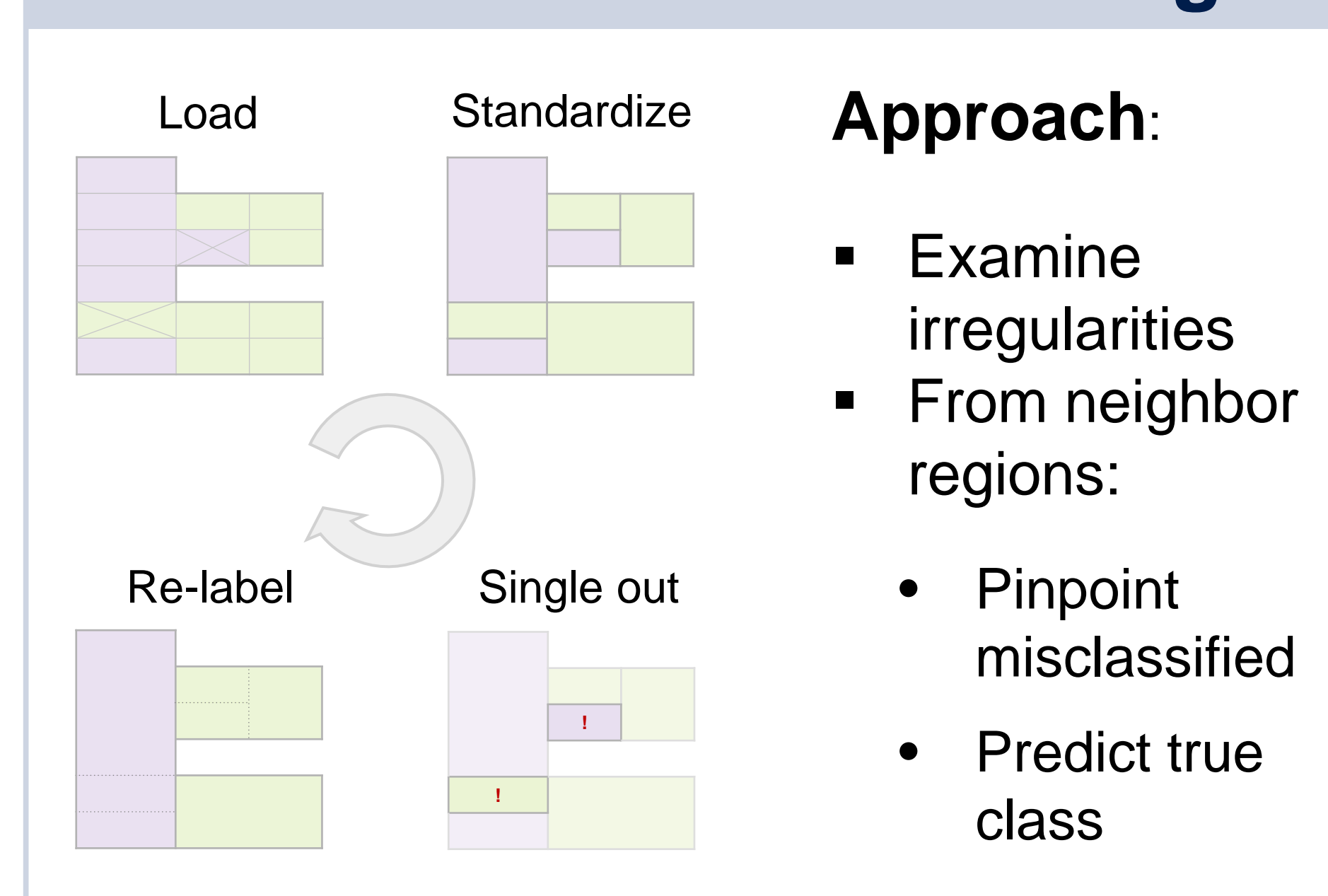
	A	B	C	D	E	F	G	H	I	J
1			Number of items (in units) sold per region							
2										
3										
4										
5										
6										
7										
8										
9										
10										
11										
12										
13										
14										
15										
16										
17										
18										
19										
20										



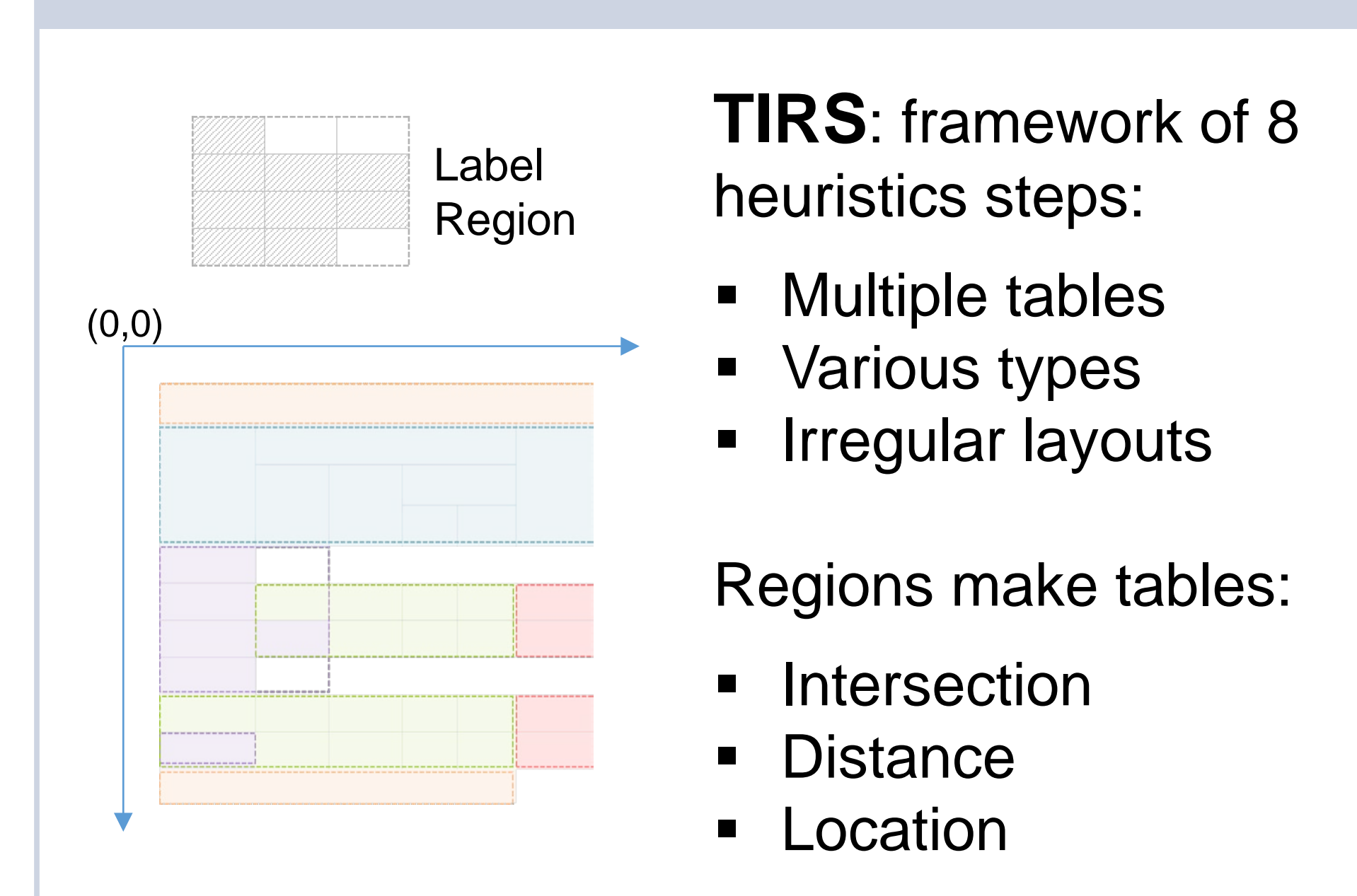
## Cell Classification



## Misclassification Handling

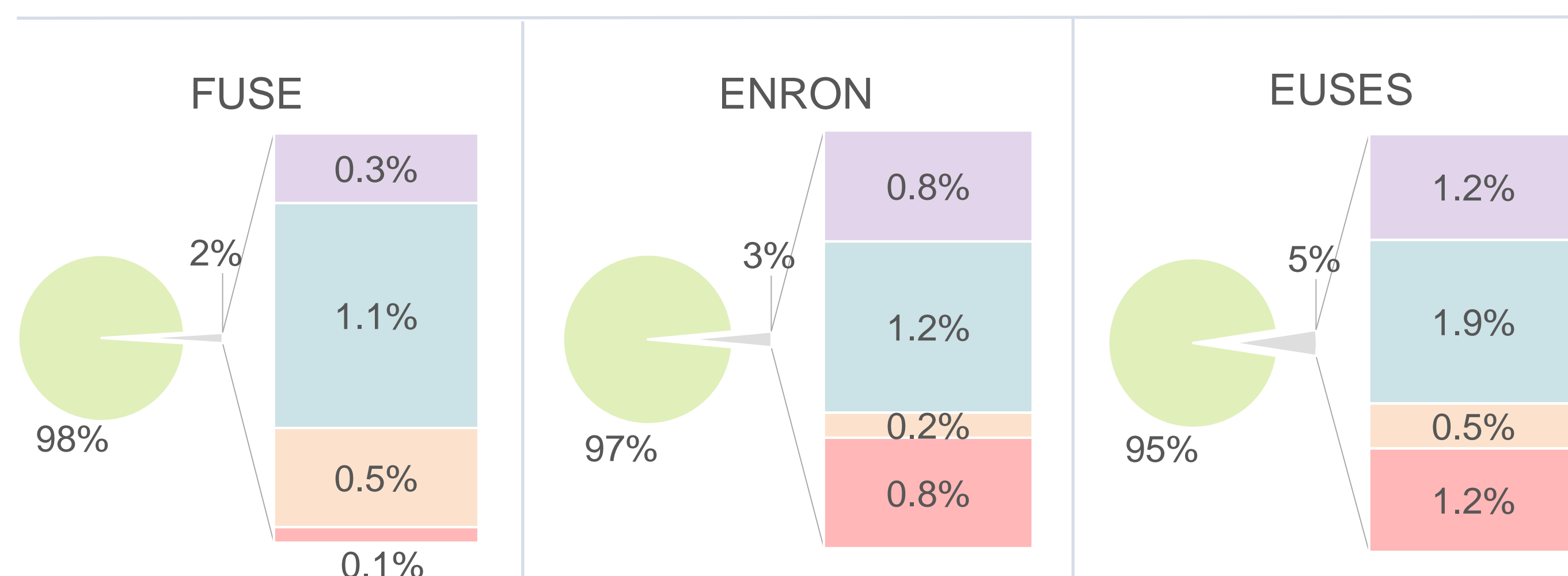


## Table Identification

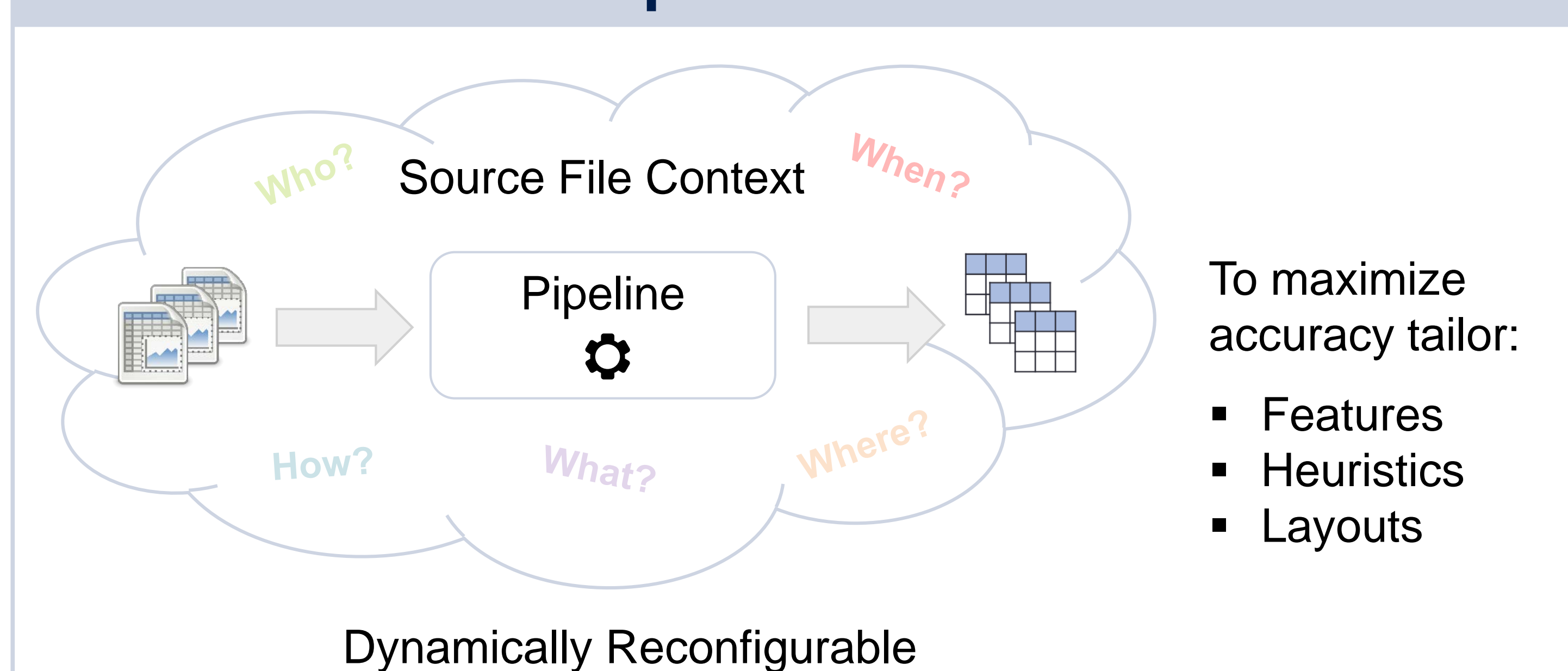


## Annotated Dataset

Contains **465** sheets, **898** tables, **>800,000** cells, from three corpora:



## Adaptive Process



- [1] Elvis Koci, Maik Thiele, Oscar Romero, and Wolfgang Lehner. A machine learning approach for layout inference in spreadsheets. IC3K'16, p. 77–88, SciTePress, 2016.
- [2] Elvis Koci, Maik Thiele, Oscar Romero, and Wolfgang Lehner. Table identification and reconstruction in spreadsheets. CAISE'17. (Accepted Paper)
- [3] Elvis Koci, Maik Thiele, Oscar Romero, and Wolfgang Lehner. Cell Classification for Layout Recognition in Spreadsheets. Chapter in CCIS book series. (In Press)