

# TabbyXL: Software Platform for Rule-Based Spreadsheet Data Extraction and Transformation

A. Shigarov, V. Khristyuk, A. Mikhailov

*Matrosov Institute for System Dynamics and Control Theory, Siberian Branch of the Russian Academy of Sciences, 134 Lermontov St., Irkutsk, Russia, 664033*

---

## Abstract

Spreadsheets are widely used in science, engineering, business, and other activities. Overall, they conceal a large volume of data in a form intended to be interpreted by humans. We present a novel software platform facilitated for liberating such data. It provides rule-based spreadsheet data extraction and transformation to a structured form. Its core consists of a flexible table object model and a domain-specific rule language for table analysis. They serve to represent knowledge of table layout and content features, as well as their interpretation depending on transformation goals. This enables processing arbitrary tables originating from various domains. Our empirical results demonstrate that one ruleset can be applied to process arbitrary tables having the same features of layout, style, or content. The paper also describes two applications using the software platform to develop programs for rule-based converting data from arbitrary spreadsheet tables.

*Keywords:* table understanding, information extraction, unstructured data management, rule-based programming, spreadsheet data, software development

---

## 1. Introduction

2 Many spreadsheet tables are designed to be interpreted by humans. They  
3 lack metadata (explicit semantics) needed for computer programs to interpret  
4 them as intended by their author or as required by an application. Spread-  
5 sheet tools provide a variety of layout structures and formatting styles for  
6 presenting tables. Their human-centeredness leads to the fact that arbitrary  
7 tables often have an incorrect structure (e.g. one logical cell can be im-  
8 properly split into several physical cells) and “messy” text data (e.g. typos,  
9 non-standardized values, extra spaces).

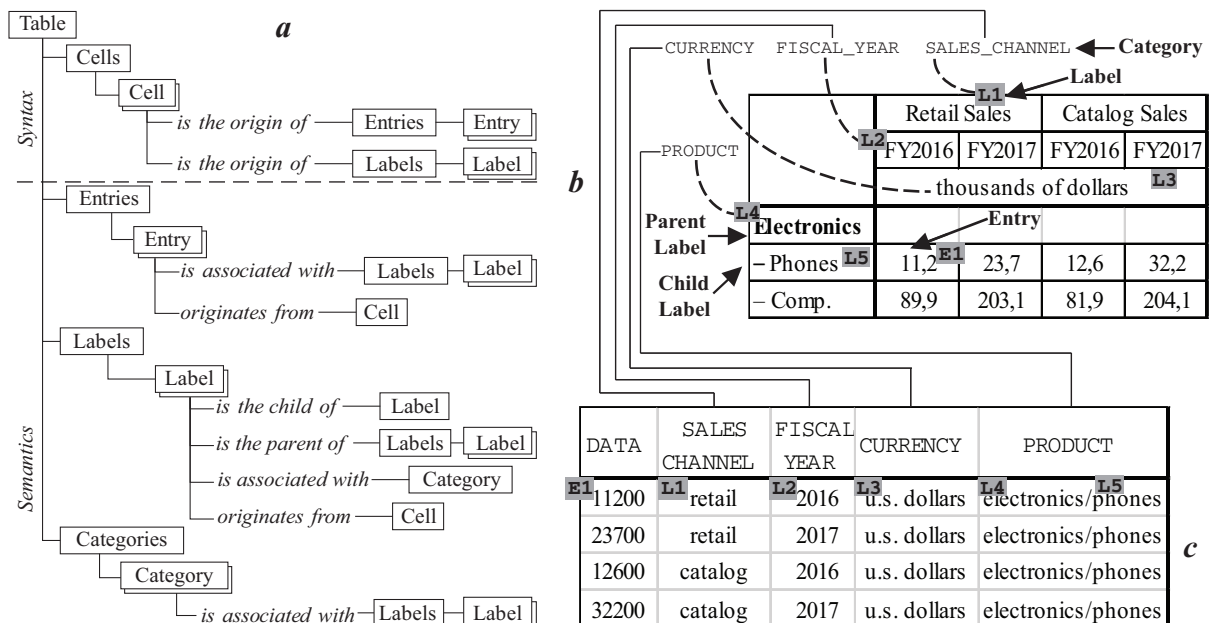


Figure 1: The concept of the table object model — *a*; an instance of an arbitrary table (source) — *b* and in a canonicalized version (target) — *c* (the tag **E1** indicates an entry, the tags **L1**,...,**L5** point out its related labels).

10 Spreadsheet data can be a valuable source in data science and business  
 11 intelligence applications. However, difficulties that inevitably arise with ex-  
 12 traction data from arbitrary tables often hinder the intensive use of them  
 13 in these areas. Typically, experts faced with such tasks resort to general-  
 14 purpose tools. In comparison with the latter, ad hoc software platforms or  
 15 frameworks can allow shortening software development time, hiding inessen-  
 16 tial details and focusing on the issues under consideration.

17 **TabbyXL**<sup>1</sup>, a novel software platform, aims at the development and execu-  
 18 tion of rule-based programs for spreadsheet data extraction and transforma-  
 19 tion from arbitrary (Fig. 1, *b*) to relational tables (Fig. 1, *c*). The platform  
 20 implements our approach to rule-based table understanding (i.e. recovering  
 21 the metadata of arbitrary tables) [1, 2]. It exploits a flexible table object  
 22 model to represent knowledge of layout and content features, as well as user-  
 23 defined rules to analyze and interpret tables, depending on transformation  
 24 goals.

<sup>1</sup><https://github.com/tabbydoc/tabbyxl/releases/tag/v1.0.4>

## 25 2. Background

26 The contemporary data integration solutions (e.g. `Talend`<sup>2</sup>, `OpenRefine`<sup>3</sup>)  
27 provide some transforming operations (e.g. “unpivoting”, “deduplication”,  
28 or “reconciliation”) for flat file databases presented in spreadsheet-like for-  
29 mats. They expect only relational tables as input. Meanwhile we consider  
30 the transformation of arbitrary tables designed by humans for humans.

31 The spreadsheet data extraction and transformation consist in recover-  
32 ing missing metadata describing the structure and content of an arbitrary  
33 table. Typically, spreadsheets do not provide all metadata that are needed to  
34 interpret table structure and content. There are no functional roles, relation-  
35 ships, and external describing concepts of data items placed in an arbitrary  
36 table. Such metadata should be recovered to convert data from arbitrary to  
37 relational tables. We refer to the considered problem as *table analysis* (re-  
38 covering the functional roles and internal relationships) and *interpretation*  
39 (recovering the external relationships).

40 The recent papers are dedicated to the related issues of table analysis and  
41 interpretation, e.g. layout features [3, 4, 5], code smells and formulas [6, 7,  
42 8, 9], programming by examples [10, 11, 12], data models [13, 14, 15], linked  
43 open data [16, 17], domain-specific [18, 19, 20] and rule-based architectures  
44 [2, 21, 22]. There are works [23, 24] with goals similar to ours. They propose  
45 methods for transforming spreadsheet tables with predefined layout features  
46 to a canonical form based on heuristics [23] and machine-learning [24]. In  
47 contrast to them we propose the software platform supporting user-defined  
48 table layout features.

49 Only a few related projects exhibited software they developed. We men-  
50 tion here only those ones that published their software at least partially.  
51 The `SSaaPP`<sup>4</sup> project implemented two frameworks for mapping data from  
52 a spreadsheet to a relational database, `HaExcel` [25] and `MDSheet` [15], as  
53 extensions for `OpenOffice`<sup>5</sup>. The `DeExcelerator`<sup>6</sup> project aimed at develop-  
54 ment of a framework for information extraction from spreadsheet tables [26].  
55 It is published only partially as `XCellAnnotator`<sup>7</sup>, a desktop application for  
56 interactively annotating cell regions in Excel files. The `Senbazuru`<sup>8</sup> project  
57 on development of a spreadsheet database management system [27] pub-

---

<sup>2</sup><https://sourceforge.net/projects/talend-studio>

<sup>3</sup><http://openrefine.org>

<sup>4</sup><http://ssaapp.di.uminho.pt>

<sup>5</sup><https://www.openoffice.org>

<sup>6</sup><https://wwwdb.inf.tu-dresden.de/misc/DeExcelerator>

<sup>7</sup><https://github.com/elviskoci/XCellAnnotator>

<sup>8</sup><http://dbgroup.eecs.umich.edu/project/sheets>

58 lished **Frame Finder**<sup>9</sup>, a software for detecting functional cell regions in a  
59 spreadsheet table having a layout called “data frame” [28].

60 Typically, the related software solutions rely on a predefined table struc-  
61 ture. They support only a few widespread layout types of tables with typical  
62 functional cell regions (e.g. “stub”, “head”, “body”, “derived”). In contrast  
63 to them, we use a domain-independent table model that is not limited by  
64 predefined functional cell regions. It is designed for specifying layout features  
65 of arbitrary tables in user-defined rules.

66 Unlike others, our model associates functional roles with data items but  
67 not cells. This enables supporting a layout where one cell contains two or  
68 more data items (e.g. in some bilingual or statistical tables). Such data items  
69 originating from the same cell can be extracted separately and be linked with  
70 different external concepts.

71 Another distinctive feature of our approach consists in the use of user-  
72 defined rules for mapping a physical structure of cells (layout, formatting  
73 style, and text) to a logical structure (linked functional data items such as  
74 entries, labels, and categories). They can be executed by a rule engine or  
75 be translated to executable programs in a general-purpose language. Our  
76 software platform implements both cases.

### 77 **3. Software Overview**

78 **TabbyXL** is developed as Java application with the command-line user  
79 interface. As an input, it expects a spreadsheet file in Excel (\*.xlsx) for-  
80 mat containing one or more arbitrary tables and a rule-based program with  
81 some user-defined rules for cleansing, analyzing, and interpreting such tab-  
82 ular data. Our software platform uses the rules to transform data from  
83 arbitrary to relational tables. As an output, a spreadsheet file containing  
84 extracted data in the relational form is generated for each source table.

85 End-users can exploit the software platform to develop rule-based pro-  
86 grams for goal-oriented extraction and transformation of data from arbitrary  
87 spreadsheet tables. **TabbyXL** provides two ways to implement and run the  
88 user-defined rules. One of them is to write the rules in a general-purpose  
89 rule-based language (e.g. Drools<sup>10</sup>, Jess<sup>11</sup>) and execute them via an appro-  
90 priate rule engine that is compatible with JSR-94 (Java Rule Engine API)<sup>12</sup>.  
91 The other is to express the rules in CRL [2], our domain-specific language,

---

<sup>9</sup>[https://github.com/chenzheruc/spreadsheet\\_framefinder](https://github.com/chenzheruc/spreadsheet_framefinder)

<sup>10</sup><https://www.drools.org>

<sup>11</sup><https://www.jessrules.com>

<sup>12</sup><https://www.jcp.org/ja/jsr>

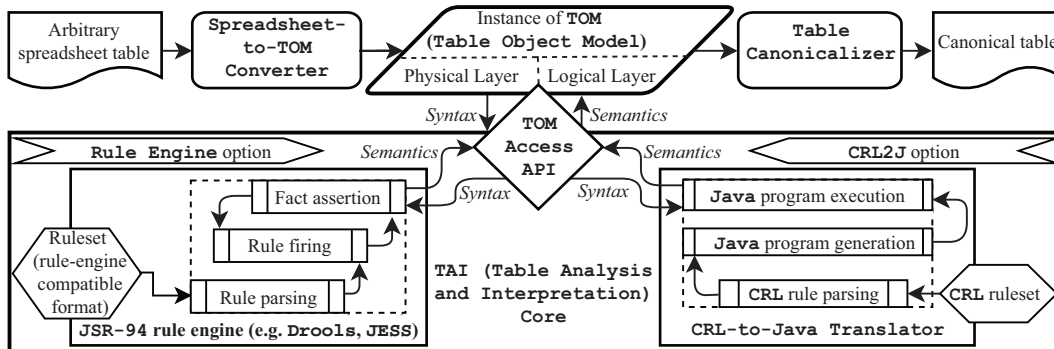


Figure 2: The architecture of the software platform.

92 and translate them to Java source code that can be compiled and executed  
 93 as Java program. The detail user manual on developing the rules is available  
 94 as a part of TabbyXL documentation<sup>13</sup>.

### 95 3.1. Software Architecture

96 The architecture shown in Fig. 2 determines the interaction of the follow-  
 97 ing main components.

98 Table Object Model (TOM) is designed for representing both physical  
 99 *cells* and logical *data items* (Fig. 1). The model includes two interrelated  
 100 layers: *physical (syntactic)* represented by a collection of cells and *logical*  
 101 (*semantic*) that consists of three collections of *entries* (values), *labels* (keys),  
 102 and *categories* (concepts). They are accessible through TOM-Access API, an  
 103 application programming interface.

104 Spreadsheet-to-TOM converter puts cells of an arbitrary table into the  
 105 physical layer of a TOM-instance.

106 Table Analysis and Interpretation (TAI) Core recovers the logical  
 107 layer of the TOM-instance from its physical layer by performing one of two op-  
 108 tions. Rule Engine option executes a ruleset in an appropriate format with a  
 109 JSR-94 compatible rule engine. The rule engine matches asserted facts (avail-  
 110 able data of the TOM-instance) against rules to create new facts (recovered  
 111 data of the TOM-instance) and assert them into the working memory. CRL2J  
 112 option provides CRL, our domain-specific language designed for expressing  
 113 table analysis and interpretation rules. It uses CRL-to-Java translator to  
 114 automatically generate Java source code from CRL rules and compile it to  
 115 Java bytecode. The generated Java program recovers missing data of the

<sup>13</sup><https://github.com/tabbydoc/tabbyxl/wiki/crl-rules>

116 TOM-instance.

117 **Table Canonicalizer** generates a canonicalized version of the processed  
118 table from the recovered data of TOM logical layer.

### 119 3.2. Software Functionality

120 The software platform enables developing programs for spreadsheet data  
121 extraction and transformation, supporting the following functions (actions)  
122 for cleansing, analysis, and interpretation of an arbitrary table represented  
123 as an instance of TOM.

124 *Cell cleansing* actions are intended to correct an inaccurate layout and  
125 content of a hand-coded table: **merge** combines two adjacent cells when they  
126 share one border; **split** divides a merged cell that spans  $n$ -tiles (row-column  
127 intersections) into  $n$ -cells; **set text** modifies a textual content of a cell; **set**  
128 **indent** modifies a text indentation of a cell. *Role analysis* actions aim at  
129 recovering entries and labels as functional data items presented in a table:  
130 **set mark** annotates a cell with a user-defined tag that can be used in subse-  
131 quent table analysis; **new entry** (**new label**) creates an entry (label) from a  
132 cell content with the use of an optional string processing. *Structural analysis*  
133 actions enable recovering pairs of two kinds: entry-label and label-label: **add**  
134 **label** associates an entry with a label; **set parent** binds two labels as a  
135 parent and its child. *Interpretation* actions serve to recover label-category  
136 pairs: **set category** associates a label with a category; **group** places two  
137 labels to one group that can be considered as an undefined category.

138 The actions aim at mapping table data of syntactic layer to semantic one  
139 of TOM. They are driven by some rulesets or programs performed with one of  
140 the implemented options of TAI-Core. Our previous paper [2] explains the  
141 listed actions in detail.

## 142 4. Implementation

143 TabbyXL implements the presented architecture as follows.

144 TOM, the table object model, is a set of Java classes corresponding to  
145 naming conventions of JavaBeans specification. Their public interfaces define  
146 TOM-Access API. Data items of a TOM-instance (cells, entries, labels, and  
147 categories) are objects of these classes. This implementation enables us to  
148 assert these objects as facts into the working memory of any rule engine that  
149 is compatible with JSR-94 specification.

150 The CRL2J option is implemented as a CRL-to-Java translator, includ-  
151 ing the following components: (i) CRL-parser generated by ANTLR develop-  
152 ment tools, (ii) Java classes to represent CRL Rule Object Model, and (iii)

153 utility programs to compile generated Java-programs for spreadsheet data  
154 extraction and transformation.

155 The **Rule Engine** option supports Drools as a rule engine by default.  
156 This expects rules expressed in DRL (the general-purpose rule language that  
157 is native for Drools) or in a dialect of CRL that is implemented as a domain-  
158 specific language (DSL) corresponding to Drools requirements. In the last  
159 case, CRL rules presented in DSLR format are automatically translated into  
160 DRL format through DSL-specification that defines CRL-to-DRL mappings.  
161 Unlike the pure CRL, this dialect supports DRL attributes in rule declara-  
162 tions. The option enables involving any JSR-94 rule engine specified in a  
163 configuration file. Our tests confirm that the Jess rule engine can also be  
164 used to execute rules represented in CLP format.

## 165 5. Empirical Results

166 This experiment exemplifies the use of our software platform for develop-  
167 ment and execution of a ruleset to extract data from arbitrary tables that are  
168 produced by different authors but pertain to the same document genre. The  
169 performance evaluation is based on Troy200 [29] dataset. It contains 200  
170 arbitrary tables as CSV files collected from 10 different sources of the same  
171 genre (government statistical websites). We added accompanying ground-  
172 truth data to automate the performance evaluation [30].

173 We designed a tested ruleset that transforms Troy200 arbitrary tables  
174 into the relational form. It was implemented in three formats: CRL, DSLR  
175 (CRL-dialect as DSL for Drools), and CLP (Jess). The ruleset services to re-  
176 cover functional data items (entries and labels) and their relationships (entry-  
177 label and label-label pairs). All of its implementations were ran by TabbyXL  
178 to automatically transform the tested tables into the relational form. The  
179 recovered relational tables were the same for all three cases.

180 The performance evaluation was automatically carried out by comparing  
181 the ground-truth data with the results of the ruleset runs. The corrected  
182 functional data items and their relationships were compared with recovered  
183 ones. We adapted the standard metrics, *recall* and *precision*, as follows:

$$\text{recall} = |R \cap S| / |S| \quad \text{precision} = |R \cap S| / |R| \quad (1)$$

184 Where  $R$  is a set of instances (entries, labels, entry-label pairs, or label-  
185 label pairs) in a target table and  $S$  is a set of instances in the corresponding  
186 source table. Table 1 presents the values of these metrics for each type of  
187 automatically recovered instances. Among 200 tested tables, only 25 are  
188 processed with errors (total 1256 false negatives in 25 tables, and 498 false

Table 1: The performance evaluation results of the tested ruleset for recovering data items and their relationships from tables of Troy200 dataset.

metrics	entries	labels	entry-label pairs	label-label pairs
recall	0.9813 $\frac{16602}{16918}$	0.9965 $\frac{4842}{4859}$	0.9773 $\frac{34270}{35066}$	0.9389 $\frac{1951}{2078}$
precision	0.9996 $\frac{16602}{16609}$	0.9364 $\frac{4842}{5171}$	0.9965 $\frac{34270}{34389}$	0.9784 $\frac{1951}{1994}$

positives in 14 ones). Only one table was not processed. This resulted in about 72% of all errors.

All data and steps to reproduce this performance evaluation are available as the published dataset [30] and as a part of TabbyXL documentation<sup>14</sup>. We also prepared a Dockerfile<sup>15</sup> to build Docker image that contains the required software and data to reproduce the presented empirical results.

We compared our tested ruleset with MIPS [23], a state-of-the-art method for segmenting a table into typical functional cell regions. The *accuracy* of table segmentation we obtained was 0.9950 against 0.9899 MIPS that authors reported on the same dataset (Troy200). We also developed and tested an additional ruleset to compare it with the published results of Senbazuru [31] on extracting header hierarchies (label-label pairs) from spreadsheet tables. The testing was performed on a subset of 200 SAUS tables randomly selected. The *F*-score obtained by TabbyXL is 0.8657 against 0.8860 reported by Senbazuru [31] on SAUS tables. Our previous paper [2] presents this comparison in more detail.

Notice that the demonstrated results are close to these state-of-the-art solutions. However, both competitors use a table structure restricted by three predefined functional regions: single-column stub, “pyramid-like” head, and body. Unlike them, our platform allows processing tables with other user-defined layout features (e.g. cut-ins in body, footers, or an inverted “pyramid-like” head).

## 6. Illustrative Example

Fig. 3 illustrates a simple example of a transformational task that consists in converting tables similar to ones (*a* and *c*) to the relational form (*b* and *d*). These tables satisfy the following assumptions:  $1, \dots, n$  are entries;  $a_1, \dots, a_m$  are column labels of the category *A*;  $b_1, \dots, b_k$  are row labels of

<sup>14</sup><https://github.com/tabbydoc/tabbyxl/wiki/performance-evaluation>

<sup>15</sup><https://hub.docker.com/r/tabbydoc/tabbyxl>



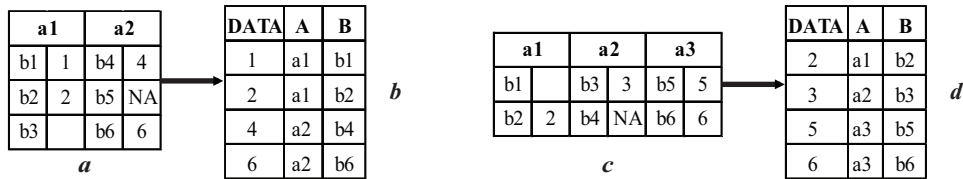


Figure 3: An illustrative example of the transformation of arbitrary tables with the same layout features (*a* and *c*) to their canonicalized versions (*b* and *d*).

```

a  when cell c: c.text.matches("NA")
    then set text "" to c

c  when cell c: (c1 % 2) == 1
    then new label c

    when
      entry e
e   label l: cell.rt == e.cell.rt, cell.cl == e.cell.cl - 1
    then add label l to e

f  when label l: cell.rt == 1
    then set category "A" to l

b  when cell c: (c1 % 2) == 0, !blank
    then new entry c

    when
      entry e
d   label l: cell.cr == e.cell.cr
    then add label l to e

g  when label l: cell.rt > 1
    then set category "B" to l

```

Figure 4: A reference ruleset for transforming tables of the illustrative example shown in Fig. 3: data cleansing — (*a*), entry generation — (*b*), label generation — (*c*), associating entries with column labels — (*d*), associating entries with row labels — (*e*), categorizing column labels — (*f*), and categorizing row labels — (*g*).

216 the category *B*. A reference ruleset designed for this task is shown in Fig. 4.  
 217 Note the demonstrated version of the ruleset is simplified by excluding `rule`  
 218 and `end` keywords and some newline characters. Its complete version is avail-  
 219 able at GitHub<sup>16</sup>. All steps to reproduce this example are presented in the  
 220 documentation<sup>17</sup>. A more complicated example is also available in the doc-  
 221 umentation<sup>18</sup>.

## 222 7. Applications

223 We used TabbyXL in two real applications. The first of them aimed at  
 224 developing a web-based statistical atlas of the Irkutsk region. Our platform  
 225 enabled populating a database with tabular data extracted from statistical  
 226 reports distributed by the Irkutsk Regional Committee of the Russian Federal  
 227 State Statistics Service. The original tables from the statistical reports were

<sup>16</sup><https://github.com/tabbydoc/examples/blob/master/tabbyxl/example1/results/crl2j/rules.crl>

<sup>17</sup><https://github.com/tabbydoc/tabbyxl/wiki/example-1>

<sup>18</sup><https://github.com/tabbydoc/tabbyxl/wiki/example-2>

228 presented as Word objects. They were converted to Excel spreadsheet format.  
229 Then a ruleset we developed in CRL format was executed by the CRL2J  
230 option. It recovered functional roles and relationships inside each table. As  
231 a result, the original tables were transformed into the canonical form by  
232 running the ruleset. Finally, they were converted to CSV format, aggregated,  
233 and loaded into a database of the statistical atlas.

234 The second application was dedicated software development from gener-  
235 ating ontologies from data of arbitrary tables used in Industrial Safety  
236 Inspection (ISI) services. The platform provided extracting data from arbi-  
237 trary tables presented in ISI reports. Such reports were originally presented  
238 as PDF documents. They described technical diagnostics, analysis (including  
239 interpretation) of diagnostics results, calculation of durability and residual  
240 resource, etc. Their content and layout were specified by some corporative  
241 standards. In the considered case, all tables could be separated into two types  
242 by their layout forms. We used two rulesets in CRL format for transforming  
243 tables of both arbitrary forms to the canonical form. The extracted data  
244 (relational tables) enabled us to generate fragments of a conceptual model  
245 as class diagrams in UML notation. Finally, such fragments were aggreg-  
246 ated into one conceptual model to construct a knowledge base on inspected  
247 objects.

248 The detail description of both applications (including the workflows, rule-  
249 sets, samples of real data, and steps to reproduce) is presented as a part of  
250 the platform documentation<sup>19</sup>.

## 251 8. Impact

252 Data analysis needs structured data. However, in practice, data often  
253 are available only in a weakly structured form, such as arbitrary spreadsheet  
254 tables. For example, there is a large volume of tabular data presented in  
255 statistical reports, financial statements, safety data sheets, or business credit  
256 assessments. Many applications of data science and business intelligence  
257 potentially can use such data.

258 The data extraction from semi-structured tabular documents such as  
259 spreadsheet workbooks can be a time-consuming process. When there is  
260 a necessity to process a large volume of arbitrary tables with various layouts  
261 then manual processing should be reduced as much as possible. In recent  
262 years, this challenge attracted the attention of the scientific community in  
263 the area of document analysis and data management. However, a few efforts

---

<sup>19</sup><https://github.com/tabbydoc/tabbyxl/wiki/statistical-atlas>  
<https://github.com/tabbydoc/tabbyxl/wiki/industrial-safety-inspection>

264 of the community were devoted to develop and publish software for tasks of  
265 the table understanding.

266 TabbyXL shows new possibilities of software development for spreadsheet  
267 data extraction and transformation. The existing software packages support  
268 only a couple of predefined types of table layouts. They usually embed  
269 some rules in their internal algorithms. In contrast to them, our platform  
270 implements a general table model and allows extending its functionality with  
271 user-defined rules.

272 Such software can significantly facilitate the data extraction from semi-  
273 structured tabular documents. Particularly, the preliminary implementation  
274 of our platform was used for filling up a data warehouse with socio-economic  
275 data on Mongolian provinces. The current version of the software platform  
276 was used in two real ETL (Extract, Transform, Load) workflows for extract-  
277 ing data from arbitrary spreadsheet tables (Section 7). In one case, it enabled  
278 populating the database of the web-based statistical atlas from tabular data  
279 of government statistical reports. In another case, it facilitated generating  
280 ontologies from data of arbitrary tables used in industrial safety inspection.

281 We believe that the design principles of our platform can be used as a  
282 basis for the development of software for the conversion of tabular data from  
283 weakly or semi-structured sources to databases.

## 284 **9. Conclusions**

285 The contribution of the presented software consists of the following re-  
286 sults. The implemented platform provides developing and executing rule-  
287 based spreadsheet data extraction and transformation programs. In contrast  
288 to the existing ETL tools, our platform supports arbitrary tables where an  
289 implicit semantics is concealed by a complex of layout, style, and content  
290 features.

291 The novelty of the software platform architecture consists in providing  
292 two rule-based ways to implement goal-oriented workflows. In the first case, a  
293 ruleset for table analysis and interpretation is expressed in a general-purpose  
294 rule language and executed by a JSR-94-compatible rule engine (e.g. Drools  
295 or Jess). In the second case, a ruleset expressed in CRL is translated to an  
296 executable Java program.

297 As an important part of the architecture, our two-layered table object  
298 model proposes a novel approach to associate functional roles with data items  
299 but not cells. Unlike other models, we assume that functional data items can  
300 be placed anywhere in a table. Thus, this provides processing such specific  
301 table layouts.

302 CRL, our domain-specific rule-based language, determines queries (con-  
303 ditions) and operations (actions) that are necessary to develop programs for  
304 spreadsheet data transformation from an arbitrary to relational form. CRL  
305 rules map a physical structure of cells (layout, style and text features) to  
306 a logical structure (linked functional data items such as entries, labels, and  
307 categories). In comparison with general-purpose rule languages (e.g. DRL  
308 or Jess), CRL enables expressing rulesets without any instructions for man-  
309 agement of the working memory such as updates of modified facts or blocks  
310 on the rule re-activation. This allows end-users to focus more on the logic of  
311 table analysis and interpretation than on the logic of the rule management  
312 and execution.

313 While the competitive solutions are limited by a few predefined table  
314 layout types, our software platform supports the table analysis and inter-  
315 pretation for various layouts. The empirical results showed that one ruleset  
316 (program) developed and performed with the software platform can process  
317 arbitrary tables originating from different sources of the same genre.

318 The limitation of this work is that we do not involve the use of spreadsheet  
319 formulas. In practice, many arbitrary tables contain formulas. In particular,  
320 they can be used to detect and validate derived data in applications of the  
321 table transformation. It would be interesting for further work to incorporate  
322 formulas in user-defined rules to recover derived data items and their internal  
323 relationships.

324 The current version of the software platform implements the table inter-  
325 pretation limited by grouping and associating labels with user-defined cat-  
326 egories. Further work can overcome this limitation by adding a new func-  
327 tionality based on the named entity recognition and linking. We believe that  
328 the linking of extracted tabular data with the global structure of linked open  
329 data (LOD cloud) would enable them to be interpreted in terms of external  
330 ontologies in third-party software applications.

### 331 **Acknowledgements**

332 This work was supported by the Russian Science Foundation [grant num-  
333 ber 18-71-10001].

### 334 **Conflict of Interest**

335 The authors declare that there is no conflict of interest.

336 **References**

- 337 [1] A. Shigarov, Table understanding using a rule engine, *Expert Systems with Applications* 42 (2) (2015) 929–937.  
338 doi:10.1016/j.eswa.2014.08.045.  
339
- 340 [2] A. O. Shigarov, A. A. Mikhailov, Rule-based spreadsheet data trans-  
341 formation from arbitrary to relational tables, *Information Systems* 71  
342 (2017) 123–136. doi:10.1016/j.is.2017.08.004.
- 343 [3] E. Koci, M. Thiele, O. Romero, W. Lehner, Table identification and  
344 reconstruction in spreadsheets, in: *Proc. 29th Int. Conf. Advanced In-*  
345 *formation Systems Engineering*, 2017, pp. 527–541. doi:10.1007/978-3-  
346 319-59536-8\_33.
- 347 [4] Z. Chen, S. Dadiomov, R. Wesley, G. Xiao, D. Cory, M. Cafarella,  
348 J. Mackinlay, Spreadsheet property detection with rule-assisted active  
349 learning, in: *Proc. ACM on Conf. on Information and Knowledge Man-*  
350 *agement*, 2017, pp. 999–1008. doi:10.1145/3132847.3132882.
- 351 [5] W. Dou, S. Han, L. Xu, D. Zhang, J. Wei, Expandable group  
352 identification in spreadsheets, in: *Proc. 33rd ACM/IEEE Int.*  
353 *Conf. on Automated Software Engineering*, 2018, pp. 498–508.  
354 doi:10.1145/3238147.3238222.
- 355 [6] F. Hermans, M. Pinzger, A. Deursen, Detecting and refactoring code  
356 smells in spreadsheet formulas, *Empirical Softw. Engg.* 20 (2) (2015)  
357 549–575. doi:10.1007/s10664-013-9296-2.
- 358 [7] W. Dou, C. Xu, S. C. Cheung, J. Wei, CACheck: Detecting and re-  
359 pairing cell arrays in spreadsheets, *IEEE Transactions on Software En-*  
360 *gineering* 43 (3) (2017) 226–251. doi:10.1109/TSE.2016.2584059.
- 361 [8] D. W. Barowy, E. D. Berger, B. Zorn, ExceLint: Automatically finding  
362 spreadsheet formula errors, *Proc. ACM Program. Lang.* 2 (OOPSLA)  
363 (2018) 148:1–148:26. doi:10.1145/3276518.
- 364 [9] P. Koch, B. Hofer, F. Wotawa, On the refinement of spreadsheet smells  
365 by means of structure information, *Journal of Systems and Software* 147  
366 (2019) 64–85. doi:10.1016/j.jss.2018.09.092.
- 367 [10] D. W. Barowy, S. Gulwani, T. Hart, B. Zorn, FlashRelate: Extract-  
368 ing relational data from semi-structured spreadsheets using examples,  
369 *SIGPLAN Not.* 50 (6) (2015) 218–228. doi:10.1145/2813885.2737952.

- 370 [11] R. Singh, S. Gulwani, Transforming spreadsheet data types  
371 using examples, *SIGPLAN Not.* 51 (1) (2016) 343–356.  
372 doi:10.1145/2914770.2837668.
- 373 [12] Z. Jin, M. R. Anderson, M. Cafarella, H. V. Jagadish, Foofah: Trans-  
374 forming data by example, in: *Proc. ACM Int. Conf. Management of*  
375 *Data*, 2017, pp. 683–698. doi:10.1145/3035918.3064034.
- 376 [13] D. Amalfitano, A. R. Fasolino, P. Tramontana, V. De Simone,  
377 G. Di Mare, S. Scala, A reverse engineering process for inferring data  
378 models from spreadsheet-based information systems: An automotive in-  
379 dustry experience, in: *Data Management Technologies and Applica-*  
380 *tions*, 2015, pp. 136–153. doi:10.1007/978-3-319-25936-9\_9.
- 381 [14] J. Cunha, J. P. Fernandes, J. Mendes, J. Saraiva, Embedding, evolution,  
382 and validation of model-driven spreadsheets, *IEEE Transactions on Soft-*  
383 *ware Engineering* 41 (3) (2015) 241–263. doi:10.1109/TSE.2014.2361141.
- 384 [15] J. Cunha, M. Erwig, J. Mendes, J. Saraiva, Model inference for spread-  
385 sheets, *Autom Softw Eng* 23 (3) (2016) 361–392. doi:10.1007/s10515-  
386 014-0167-x.
- 387 [16] D. Ritze, C. Bizer, Matching web tables to dbpedia - A feature utility  
388 study, in: *Proc. 20th Int. Conf. on Extending Database Technology*,  
389 2017, pp. 210–221. doi:10.5441/002/edbt.2017.20.
- 390 [17] Z. Zhang, Effective and efficient semantic table interpretation using  
391 TableMiner<sup>+</sup>, *Semantic Web* 8 (6) (2017) 921–957. doi:10.3233/SW-  
392 160242.
- 393 [18] M. de Vos, J. Wielemaker, H. Rijgersberg, G. Schreiber, B. Wielinga,  
394 J. Top, Combining information on structure and content to automat-  
395 ically annotate natural science spreadsheets, *Int. J. Human-Computer*  
396 *Studies* 103 (2017) 63–76. doi:10.1016/j.ijhcs.2017.02.006.
- 397 [19] T. D. Cao, I. Manolescu, X. Tannier, Extracting linked data from statis-  
398 tic spreadsheets, in: *Proc. Int. Workshop Semantic Big Data*, 2017, pp.  
399 5:1–5:5. doi:10.1145/3066911.3066914.
- 400 [20] A. Swidan, F. Hermans, Semi-automatic extraction of cross-table data  
401 from a set of spreadsheets, in: S. Barbosa, P. Markopoulos, F. Paternò,  
402 S. Stumpf, S. Valtolina (Eds.), *End-User Development*, 2017, pp. 84–99.  
403 doi:10.1007/978-3-319-58735-6\_6.

- 404 [21] S. Yang, J. Guo, R. Wei, Semantic interoperability with heteroge-  
405 neous information systems on the internet through automatic tab-  
406 ular document exchange, *Information Systems* 69 (2017) 195–217.  
407 doi:10.1016/j.is.2016.10.010.
- 408 [22] S. Yang, R. Wei, A. Shigarov, Semantic interoperability for electronic  
409 business through a novel cross-context semantic document exchange ap-  
410 proach, in: *Proc. ACM Symposium on Doc. Eng.*, 2018, pp. 28:1–28:10.  
411 doi:10.1145/3209280.3209523.
- 412 [23] D. W. Embley, M. S. Krishnamoorthy, G. Nagy, S. Seth, Converting  
413 heterogeneous statistical tables on the web to searchable databases,  
414 *Int. J. Document Analysis and Recognition* 19 (2) (2016) 119–138.  
415 doi:10.1007/s10032-016-0259-1.
- 416 [24] Z. Chen, Information extraction on para-relational data, Ph.D. thesis,  
417 University of Michigan, US (2016).
- 418 [25] J. Cunha, J. a. Saraiva, J. Visser, From spreadsheets to rela-  
419 tional databases and back, in: *Proc. ACM SIGPLAN Workshop*  
420 *Partial Evaluation and Program Manipulation*, 2009, pp. 179–188.  
421 doi:10.1145/1480945.1480972.
- 422 [26] J. Eberius, C. Werner, M. Thiele, K. Braunschweig, L. Dannecker,  
423 W. Lehner, DeExcelerator: a framework for extracting relational  
424 data from partially structured documents, in: *Proc. 22nd ACM Int.*  
425 *Conf. on Information & Knowledge Management*, 2013, pp. 2477–2480.  
426 doi:10.1145/2505515.2508210.
- 427 [27] Z. Chen, M. Cafarella, J. Chen, D. Prevo, J. Zhuang, Senbazuru: A pro-  
428 totype spreadsheet database management system, *Proc. VLDB Endow.*  
429 6 (12) (2013) 1202–1205. doi:10.14778/2536274.2536276.
- 430 [28] Z. Chen, M. Cafarella, Automatic web spreadsheet data extraction, in:  
431 *Proc. 3rd Int. Workshop Semantic Search Over the Web*, 2013, pp. 1:1–  
432 1:8. doi:10.1145/2509908.2509909.
- 433 [29] G. Nagy, TANGO-DocLab web tables from international statistical sites  
434 (Troy\_200), 1, ID: Troy\_200\_1 (2016).  
435 URL [http://tc11.cvc.uab.es/datasets/Troy\\_200\\_1](http://tc11.cvc.uab.es/datasets/Troy_200_1)
- 436 [30] A. Shigarov, V. Khristyuk, TabbyXL: Dataset for the performance eval-  
437 uation of a software platform for rule-based spreadsheet data extraction

438 and transformation, Mendeley Data, v5 (2018).  
439 URL <http://dx.doi.org/10.17632/ydcr7mcrtf.5>

440 [31] Z. Chen, M. Cafarella, Integrating spreadsheet data via accu-  
441 rate and low-effort extraction, in: Proc. 20th ACM SIGKDD Int.  
442 Conf. Knowledge Discovery and Data Mining, 2014, pp. 1126–1135.  
443 doi:10.1145/2623330.2623617.



444 **Required Metadata**

445 **Current code version**

Nr.	Code metadata description	Description
C1	Current code version	1.0.4
C2	Permanent link to code/repository used of this code version	<i>https : //github.com/tabbydoc/tabbyxl/releases/tag/v1.0.4</i>
C3	Legal Code License	Apache License 2.0
C4	Code versioning system used	Git
C5	Software code languages, tools, and services used	Java
C6	Compilation requirements, operating environments & dependencies	Java Development Kit 8 or more, Apache Maven
C7	If available Link to developer documentation/manual	<i>https : //github.com/tabbydoc/tabbyxl/wiki</i>
C8	Support email for questions	shigarov@icc.ru

Table 2: Code metadata (mandatory)

446 **Current executable software version**

Nr.	(executable) Software metadata description	Description
S1	Current software version	1.0.4
S2	Permanent link to executables of this version	<i>https : //github.com/tabbydoc/tabbyxl/releases/tag/v1.0.4</i>
S3	Legal Software License	Apache License 2.0
S4	Computing platform/Operating System	Linux, OS X, Microsoft Windows
S5	Installation requirements & dependencies	Java SE Runtime Environment 8 or more
S6	If available, link to user manual - if formally published include a reference to the publication in the reference list	<i>https : //github.com/tabbydoc/tabbyxl/wiki</i>
S7	Support email for questions	shigarov@icc.ru

Table 3: Software metadata (optional)