

1 **Spreadsheet Data Extractor (SDE): A Performance-Optimized,**  
2 **User-Centric Tool for Transforming Semi-Structured Excel**  
3 **Spreadsheets into Relational Data**

4  
5 ANONYMOUS AUTHOR(S)  
6  
7 SUBMISSION ID:  
8

9 Organizations across various sectors frequently struggle to analyze and utilize semi-structured data derived  
10 from spreadsheets due to the lack of defined structure. This paper introduces the Spreadsheet Data Extractor  
11 (SDE), an open-source tool designed to convert semi-structured spreadsheet data into structured formats  
12 without requiring programming knowledge. Building upon previous work, we have enhanced the SDE  
13 with incremental loading of worksheets, accurate rendering of cell dimensions by parsing XML data, and  
14 performance optimizations to handle large datasets efficiently. We compare our tool with existing solutions  
15 and demonstrate its effectiveness through performance evaluations, highlighting its potential to facilitate  
16 efficient and reliable data extraction from diverse spreadsheet formats.

17 CCS Concepts: • **Applied computing → Spreadsheets;** • **Information systems → Data cleaning;** •  
18 **Software and its engineering → Extensible Markup Language (XML);** • **Human-centered computing**  
19 → *Graphical user interfaces;* • **Theory of computation → Data compression.**

20 **ACM Reference Format:**

21 Anonymous Author(s). 2025. Spreadsheet Data Extractor (SDE): A Performance-Optimized, User-Centric Tool  
22 for Transforming Semi-Structured Excel Spreadsheets into Relational Data. In *Proceedings of 2025 International*  
23 *Conference on Management of Data (ACM PODS '25)*. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/XXXXXX.XXXXXX>

24 **1 INTRODUCTION**

25 Spreadsheets are ubiquitous tools used across various domains, including healthcare [3], nonprofit  
26 organizations [10, 12], finance, commerce, academia, and government [6]. Despite their widespread  
27 use, analyzing and utilizing data stored in spreadsheets poses significant challenges due to their  
28 semi-structured nature. Data in spreadsheets are often formatted for human readability, employing  
29 layouts with empty cells, merged cells, hierarchical headers, and multiple tables, which hinder  
30 machine readability and automated data processing.  
31

32 While unstructured data in spreadsheets have advantages—such as an easily comprehensible  
33 hierarchy of metadata for humans—the same features complicate automated data extraction.  
34

35 The objective of this work is to transform semi-structured spreadsheet data into machine-  
36 readable formats. Existing approaches often rely on heuristics, machine learning, or programming-  
37 by-example techniques, which may introduce errors or require significant Aufwand um diese  
38 Fehler zu finden und zu reparieren. To address these challenges, we present the Spreadsheet Data  
39 Extractor (SDE), a tool that enables users to define data hierarchies through cell selection without  
40 any programming knowledge.

---

41 Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee  
42 provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the  
43 full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored.  
44 Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires  
45 prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

46 *ACM PODS '25, June 22–27, 2025, Berlin, Germany*

47 © 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

48 ACM ISBN ...\$15.00

49 <https://doi.org/10.1145/XXXXXX.XXXXXX>

## 50 1.1 Contributions

51 Our contributions are as follows:

- 52 (1) We release the SDE under the open-source GNU General Public License v3.0, promoting  
53 community access and collaboration.
- 54 (2) We implement incremental loading of worksheets to enhance performance, allowing the  
55 tool to handle large Excel files efficiently.
- 56 (3) We accurately render row heights and column widths by parsing XML data, ensuring that  
57 the spreadsheet's visual representation closely matches that of Excel.
- 58 (4) We optimize the rendering engine to draw only the visible cells, significantly improving  
59 performance when dealing with large datasets.
- 60 (5) We integrate the selection hierarchy, worksheet view, and output preview into a unified  
61 interface, streamlining the user experience.

## 63 2 INTRODUCTION

64 Previous work by Alexander Aue, Norbert Röder, and Andrea Ackermann introduced a tool that  
65 facilitates data extraction from Excel files. [1] While effective, this solution faced performance  
66 issues and inaccuracies in rendering cell dimensions, limiting its usability with large and complex  
67 datasets. In this paper, we build upon their foundational work by releasing the software under the  
68 open-source GNU Public License Version 3 (GPLv3) and introducing several key enhancements.  
69 These improvements include implementing incremental loading of worksheets to enhance perfor-  
70 mance, accurately rendering row heights and column widths by parsing XML data, and optimizing  
71 the rendering engine to handle large datasets efficiently. Additionally, we have integrated the  
72 selection hierarchy and worksheet view into a unified interface to improve user experience. These  
73 contributions collectively address the limitations of the existing solution, making data extraction  
74 from complex spreadsheets more efficient and user-friendly.

## 76 3 RELATED WORK

77 The extraction of relational data from semi-structured documents, particularly spreadsheets, has  
78 garnered significant attention due to their ubiquitous use across domains such as business, govern-  
79 ment, and scientific research. Several frameworks and tools have been developed to address the  
80 challenges of converting flexible spreadsheet formats into normalized relational forms suitable for  
81 data analysis and integration. Notable among these are **DeExcelerator**, **XLIIndy**, **FLASHRELATE**,  
82 **Senbazuru**, **TableSense** und den Ansatz von Aue et al. auf dessen arbeit wir aufbauen.

### 84 3.1 Aue et al.'s Converter

85 Aue et al. [1] developed a tool aimed at facilitating data extraction from Excel spreadsheets by  
86 utilizing the Dart 'excel' package to open '.xlsx' files. Users can select cells containing data and  
87 metadata to define the data hierarchy. However, this method encounters performance bottlenecks  
88 as the package requires loading the entire '.xlsx' file into memory before processing, leading to  
89 slow response times, especially with large files. Die Lösung nutzte das Dart-Package excel, um  
90 die .xlsx-Dateien zu öffnen. Dies war jedoch sehr langsam, da das Paket die gesamte .xlsx-Datei  
91 zunächst vollständig einliest. Wir haben daher eine eigene Funktionalität in Dart implementiert,  
92 die die Excel-Arbeitsblätter inkrementell lädt. Additionally, their solution calculates row heights  
93 and column widths based solely on cell content, disregarding the actual dimensions specified in the  
94 Excel file. This results in discrepancies between the tool's rendering and the original spreadsheet.  
95 The tool also renders all cells regardless of their visibility within the viewport, causing significant  
96 performance degradation when handling worksheets with numerous cells.

### 99      3.2 DeExcelerator

100     Eberius et al. [7] introduced **DeExcelerator**, a framework that transforms partially structured  
101    spreadsheets into first normal form relational tables using heuristic-based extraction phases. It  
102    addresses challenges such as table detection, metadata extraction, and layout normalization. While  
103    effective in automating normalization, its reliance on predefined heuristics limits adaptability to  
104    heterogeneous or unconventional spreadsheet formats, highlighting the need for more flexible  
105    approaches.

### 107      3.3 XLIndy

108     Koci et al. [8] developed **XLIndy**, an interactive Excel add-in with a Python-based machine learning  
109    backend. Unlike DeExcelerator's fully automated heuristic approach, XLIndy integrates machine  
110    learning techniques for layout inference and table recognition, enabling a more adaptable and  
111    accurate extraction process. XLIndy's interactive interface allows users to visually inspect extraction  
112    results, adjust configurations, and compare different extraction runs, facilitating iterative fine-tuning.  
113    Additionally, users can manually revise predicted layouts and tables, saving these revisions as  
114    annotations to improve classifier performance through (re-)training. This user-centric approach  
115    enhances the tool's flexibility, allowing it to accommodate diverse spreadsheet formats and user-  
116    specific requirements more effectively than purely heuristic-based systems.

### 118      3.4 FLASHRELATE

119     Barowy et al. [2] presented **FLASHRELATE**, an approach that empowers users to extract structured  
120    relational data from semi-structured spreadsheets without requiring programming expertise.  
121    FLASHRELATE introduces a domain-specific language, **FLARE**, which extends traditional regular  
122    expressions with spatial constraints to capture the geometric relationships inherent in spreadsheet  
123    layouts. Additionally, FLASHRELATE employs an algorithm that synthesizes FLARE programs  
124    from a small number of user-provided positive and negative examples, significantly simplifying the  
125    automated data extraction process.

126     FLASHRELATE distinguishes itself from both DeExcelerator and XLIndy by leveraging programming-  
127    by-example (PBE) techniques. While DeExcelerator relies on predefined heuristic rules and XLIndy  
128    incorporates machine learning models requiring user interaction for fine-tuning, FLASHRELATE  
129    allows non-expert users to define extraction patterns through intuitive examples. This approach  
130    lowers the barrier to entry for extracting relational data from complex spreadsheet encodings,  
131    making the tool accessible to a broader range of users.

### 134      3.5 Senbazuru

135     Chen et al. [4] introduced **Senbazuru**, a prototype Spreadsheet Database Management System  
136    (SSDBMS) designed to extract relational information from a large corpus of spreadsheets. Senbazuru  
137    addresses the critical issue of integrating data across multiple spreadsheets, which often lack explicit  
138    relational metadata, thereby hindering the use of traditional relational tools for data integration  
139    and analysis.

140     Senbazuru comprises three primary functional components:

- 141       (1) **Search**: Utilizing a textual search-and-rank interface, Senbazuru enables users to quickly  
142          locate relevant spreadsheets within a vast corpus. The search component indexes spread-  
143          sheets using Apache Lucene, allowing for efficient retrieval based on relevance to user  
144          queries.
- 145       (2) **Extract**: The extraction pipeline in Senbazuru consists of several stages:

- **Frame Finder:** Identifies data frame structures within spreadsheets using Conditional Random Fields (CRFs) to assign semantic labels to non-empty rows, effectively detecting rectangular value regions and associated attribute regions.
  - **Hierarchy Extractor:** Recovers attribute hierarchies for both left and top attribute regions. This stage also incorporates a user-interactive repair interface, allowing users to manually correct extraction errors, which the system then generalizes to similar instances using probabilistic methods.
  - **Tuple Builder and Relation Constructor:** Generates relational tuples from the extracted data frames and assembles these tuples into coherent relational tables by clustering attributes and recovering column labels using external schema repositories like Freebase and YAGO.
- (3) **Query:** Supports basic relational operations such as selection and join on the extracted relational tables, enabling users to perform complex data analysis tasks without needing to write SQL queries.

Senbazuru's ability to handle hierarchical spreadsheets, where attributes may span multiple rows or columns without explicit labeling, sets it apart from earlier systems like DeExcelerator and XLIIndy. By employing machine learning techniques and providing user-friendly repair interfaces, Senbazuru ensures high-quality extraction and facilitates the integration of spreadsheet data into relational databases.

### 3.6 TableSense

Dong et al. [5] developed **TableSense**, an end-to-end framework for spreadsheet table detection using Convolutional Neural Networks (CNNs). TableSense addresses the diversity of table structures and layouts by introducing a comprehensive cell featurization scheme, a Precise Bounding Box Regression (PBR) module for accurate boundary detection, and an active learning framework to efficiently build a robust training dataset.

While **DeExcelerator**, **XLIIndy**, **FLASHRELATE**, and **Senbazuru** focus primarily on transforming spreadsheet data into relational forms through heuristic, machine learning, and programming-by-example approaches, **TableSense** specifically targets the accurate detection of table boundaries within spreadsheets using deep learning techniques. Unlike region-growth-based methods employed in commodity spreadsheet tools, which often fail on complex table layouts, TableSense achieves superior precision and recall by leveraging CNNs tailored for the unique characteristics of spreadsheet data. However, TableSense focuses on table detection and visualization, allowing users to generate diagrams from the detected tables but does not provide functionality for exporting the extracted data for further analysis.

### 3.7 Comparison and Positioning

While **DeExcelerator**, **XLIIndy**, **FLASHRELATE**, **Senbazuru**, and **TableSense** each offer unique approaches to spreadsheet data extraction, they share certain limitations. Many of these tools are not readily accessible: **FLASHRELATE** and **TableSense** are proprietary, and **Senbazuru**, **XLIIndy**, and **DeExcelerator** are discontinued projects with limited or no source code availability. In contrast, we contribute our spreadsheet data extractor under the GNU General Public License v3.0, allowing the community to access, use, and improve the tool freely.

Moreover, unlike the aforementioned tools that rely on heuristics, machine learning, or AI techniques—which can introduce errors requiring users to identify and correct—we adopt a user-centric approach that gives users full control over data selection and metadata hierarchy definition. While this requires more manual input, it eliminates the uncertainty and potential inaccuracies

associated with automated methods. To streamline the process and enhance efficiency, our tool includes user-friendly features such as the ability to duplicate hierarchies of columns and tables, and to move them over similar structures for reuse, reducing the need for repetitive configurations.

By combining the strengths of manual control with enhanced user interface features and performance optimizations, our tool offers a robust and accessible solution for extracting relational data from complex and visually intricate spreadsheets. These enhancements not only improve performance and accuracy but also elevate the overall user experience, making our tool a valuable asset for efficient and reliable data extraction from diverse spreadsheet formats.

## 4 METHODOLOGY

In this section, we detail the design and implementation of the Spreadsheet Data Extractor (SDE), emphasizing its user-centric approach and performance optimizations. The SDE enables users to transform semi-structured spreadsheet data into structured, machine-readable formats without requiring programming expertise. We achieve this through an intuitive interface that allows for cell selection and hierarchy definition, incremental loading of worksheets, accurate rendering of cell dimensions, and optimized performance for handling large datasets by incrementally loading of worksheets and dadurch dass wir nur solche Zellen rndern, die in der aktuellen ansicht auch sichtbar sind.

### 4.1 User-Centric Data Extraction

The core functionality of the SDE revolves around allowing users to select cells containing data and metadata to define a data hierarchy. This process is facilitated through a graphical interface that displays the spreadsheet and allows for intuitive selection and manipulation of the selection hierarchy.

*4.1.1 Hierarchy Definition.* Users can select individual cells or ranges of cells by clicking and pressing using shift-click for multi-selection. These selections represent either data or metadata.

The selected cells are organized into a hierarchical tree structure, where each node represents a data element, and child nodes represent nested data or metadata. This hierarchy defines how the data will be transformed into a structured format.

*4.1.2 Reusability and Efficiency.* To optimize the extraction process and reduce repetitive tasks, the SDE allows users to duplicate previously defined hierarchies and apply them to similar regions within the spreadsheet. This feature is particularly useful for spreadsheets with repeating structures, such as multiple tables with the same format.

### 4.2 Example Workflow

Consider a spreadsheet containing statistical forecasts of future nursing staff availability in Germany [11]. The SDE interface consists of three main components (Figure 1):

**Hierarchy Panel (Top Left):** Displays the hierarchy of cell selections, initially empty.

**Spreadsheet View (Top Right):** Shows the currently opened Excel worksheet for cell selection.

**Output Preview (Bottom):** Provides immediate feedback on the data extraction based on current selections.

*4.2.1 Step-by-Step Extraction:* The user adds a node to the hierarchy and selects the cell containing the metadata "Nursing Staff" (Figure ??). Diese Zelle steht für eine Metainformation, welche die alle Zellen in diesem Arbeitsblatt gemeinsam haben. Daher sollte sie als erstes ausgewählt werden and should appear at the beginning of each row in the output CSV file.

Define Subtables:



Fig. 1. The SDE Interface Overview.

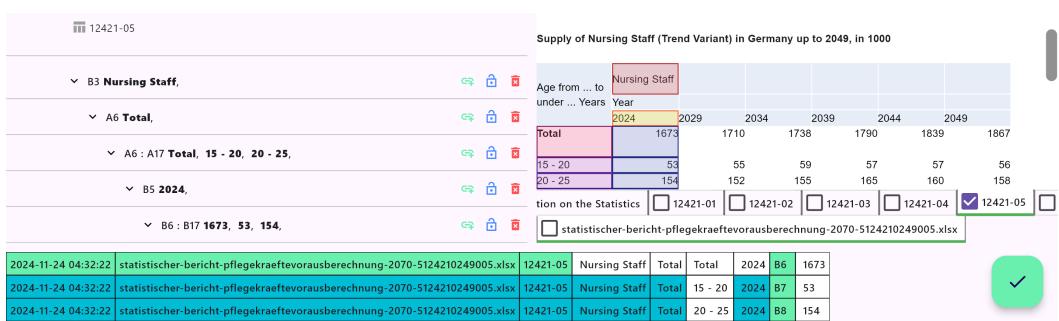


Fig. 2

Within this node, the user adds a child node and selects the cell "Total" which serves both as a table header and a row label. This selection represents the first subtable.

#### Select Row Labels:

The user adds another child node and selects the range of cells containing row labels ("Total," "15-20," "20-25") by clicking the first cell and shift-clicking the last cell.

#### Select Column Data:

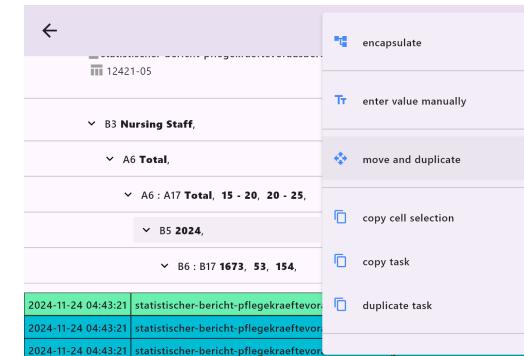
A child node is added under the row labels node, and the user selects the year "2024" and after that another child node is added under the year node and the user selects the corresponding data cells (e.g., "1673," "53," "154").

Die Hierarchie ist nun gefüllt mit 5 Knoten von denen jedes bis auf den letzten einen eingebetteten Kind-Knoten haben. Im Bereich oben rechts werden die selektierten Zellen mit unterschiedlichen Farben für jeden Knoten angezeigt.

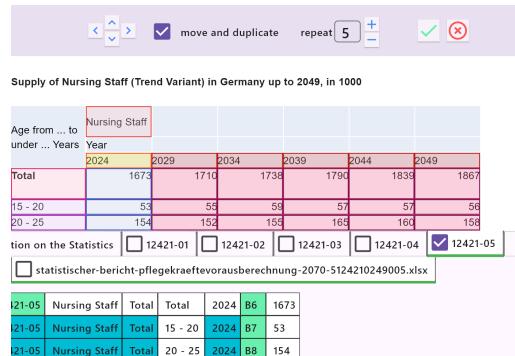
Im unteren Bereich wird die Ausgabe angezeigt. Für jeden Knoten der als Kind eingebettet wird, erscheint in der Ausgabe eine weitere Spalte. Sind in der Selektion nicht nur einzelne Zellen angegeben, so erscheinen die Werte der Liste auch in der Ausgabe als Werte in einer neuen Zeile.

#### Duplicate and Adjust Hierarchies:

To avoid repetitive manual entry for additional years, the user duplicates the hierarchy for "2024" and adjusts the cell selections to include data for subsequent years (e.g., "2025," "2026") using the "Move and Duplicate" feature. Dazu wird der Knoten der ersten Spalte "2024" selektiert und darauf rechtsgeklickt. Ein Popup öffnet sich in dem die Aktion "move and duplicate" auftaucht, die nun geklickt werden sollte, wie in Abbildung 3a zu sehen ist. This feature allows the user to shift cell selections horizontally or vertically and specify the number of repetitions, streamlining the process of capturing similar data structures.



(a)



(b)

Fig. 3

Daraufhin öffnet sich in der App Bar oben rechts eine Reihe von schaltflächen, die erlauben die Zellenselektionen des Knotens sowie aller Kinderknoten zu verschieben wie in Abbildung 3b zu sehen ist. Durch das drücken auf die Schaltfläche zum verschieben der Selektion um eine Einheit nach rechts, ist die nächste Spalte ausgewählt, jedoch wird damit auch die Selektion der ersten Spalte wieder aufgehoben, da die Selektion verschoben wurde. Damit auch die erste Spalte erhalten bleibt, kann die Checkbox "move and duplicate" aktiviert werden. Dadurch wird die verschobene Selektion zusätzlich zur ursprungsselektion erstellt. Die Änderungen werden allerdings erst übernommen wenn auf den Akzeptieren button geklickt wird. Die nächsten SPalten könnten auf die gleich Art und Weise ebenfalls selektiert werden. Doch das geht auch schneller, denn statt nur einmal die selektion zu verschieben und gleichzeitig zu duplizieren, kann auch das eingabefeld "repeat" mit so vielen Wiederholungen gefüllt werden, wie es Spalten gibt. Mit der Eingabe der Nummer 5 wird damit die Selektion der ersten Spalte 5 mal um eine Einheit nach rechts verschoben und dabei bei jeden Schritt dupliziert.

### Finalize Selection:

The user reviews the selections in the spreadsheet view, where each selection is highlighted in a different color corresponding to its node in the hierarchy. Erst nachdem der Nutzer die verschobenen und duplizierten selektionen in der Worksheet ansieht überprüft hat und den Akzeptieren button geklickt hat, werden die Knoten in der Hierachie wie gewünscht angelegt. Das Ergebnis dieser Operation ist in Abbildung 4 zu sehen.

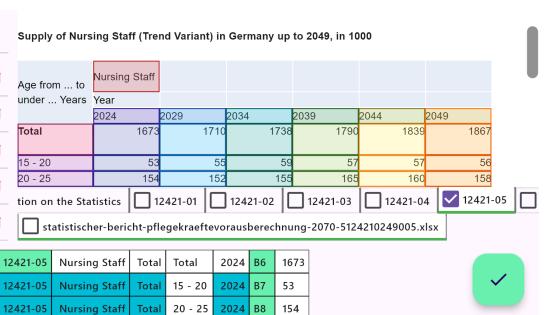
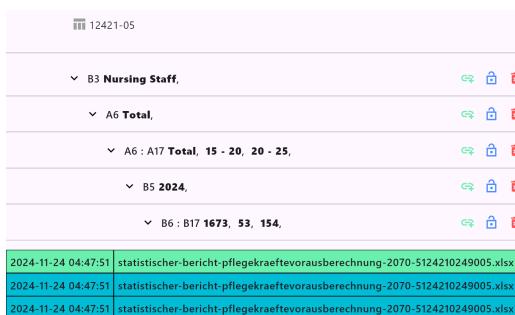


Fig. 4

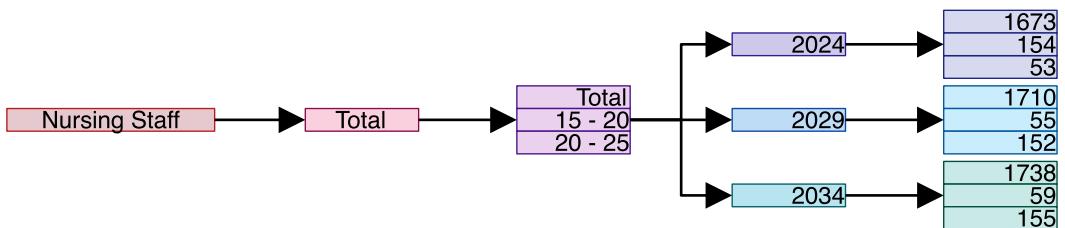
344      **Cross Product Transformation:**

345      Damit ist die erste Tabelle vollständig selektiert. Der daraus resultierende Graph ist in Abbildung  
 346      5 zu sehen. Zur Vereinfachung ist das Beispiel auf die ersten 3 Spalten und die ersten 3 Zeilen der  
 347      Tabelle beschränkt.

348      Once the hierarchy is defined, the SDE applies a cross product operation to generate a relational  
 349      format from the selection graph. This involves:

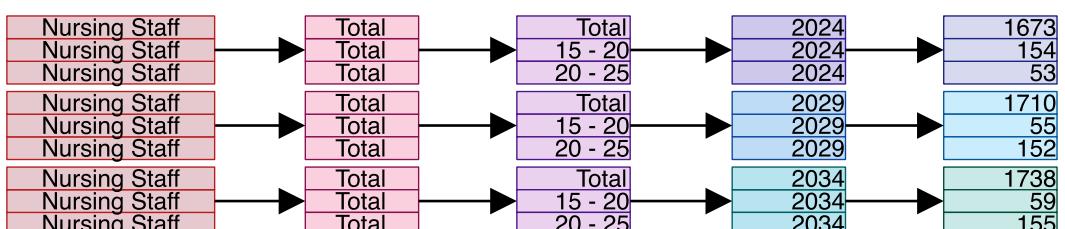
350      Node Duplication: Nodes with multiple edges (e.g., the row labels "Total," "15-20," "20-25") are  
 351      duplicated so that each path in the graph represents a unique data row.

352      Value Replication: Nodes with single values (e.g., the year "2024") are replicated to match the  
 353      number values of the associated node , ensuring alignment in the output.



363      Fig. 5. Illustration of the cross product transformation.

364      Um aus dem Selektionsgraphen ein relationales Format zu generieren wird darauf das cross  
 365      product angewendet. Was dabei genau passiert ist in Abbildung 6 zu sehen. Dort wo der Graph  
 366      mehrere Kanten hat, was in dem Beispiel für den Knoten mit der Liste der Werte "Total", "15-20",  
 367      "20-25" gilt, wird der Knoten so häufig dupliziert, wie es Kanten gibt, sodass jeder Knoten nur noch  
 368      eine Kante zum nächsten Knoten hat. Darüber hinaus werden Knoten, die keine Liste von Werten  
 369      enthalten, sondern nur einen Wert enthalten, in eine Liste umgewandelt und ihr Wert wird so  
 370      häufig dupliziert, bis die Anzahl dem Verknüpften Knoten entspricht. So wird in dem Beispiel der  
 371      Knoten mit dem Wert "2024" in eine Liste mit den Werten "2024", "2024", "2024" umgewandelt, damit  
 372      er die gleiche Anzahl an Werten hat wie der verknüpfte Knoten mit den Werten "1673", "154", "53".  
 373



383      Fig. 6

384      Das gleiche was für die Spalten bereits so gut funktioniert hat, kann nun auch für die Untertabellen  
 385      wiederholt werden, wie in Abbildung 7 zu sehen ist.

386      Durch das Selektieren des Knotens mit dem Wert "Total" und das Drücken auf die Schaltfläche  
 387      "move and duplicate" kann die Selektion der Untertabelle "Total" auf die anderen Untertabellen  
 388      angewendet werden. Dazu muss die Tabelle um so viele Zeilen nach unten verschoben werden, bis  
 389      sie mit der Untertabelle überlappt. Dabei gibt es nur ein kleines Problem. Denn zu den Unterknoten  
 390      des Knotens "Total" gehören auch die Spaltenheader dazu. Würden diese Spaltenheader in der  
 391

393 <

394    v B3 Nursing Staff.

395       v A6 Total.

396          v A6 : A17 Total, 15 - 20, 20 - 25,

397          v B5 2024,

398             v B6 : B17 1673, 53, 154,

399          v C5 2029,

400             v C6 : C17 1710, 55, 152,

401          v D5 2034,

402             v D6 : D17 1738, 59, 155,

403          v E5 2039,

404             v E6 : E17 1790, 57, 165,

405          v F5 2044,

406             v F6 : F17 1839, 57, 160,

407       v G5 2049,

408          v G6 : G17 1867, 56, 158,

409          v G7 1887, 56, 158,

410          v G8 1887, 56, 158,

411          v G9 1887, 56, 158,

412          v G10 1887, 56, 158,

413          v G11 1887, 56, 158,

414          v G12 1887, 56, 158,

415          v G13 1887, 56, 158,

416          v G14 1887, 56, 158,

417          v G15 1887, 56, 158,

418          v G16 1887, 56, 158,

419          v G17 1887, 56, 158,

420          v G18 1887, 56, 158,

421          v G19 1887, 56, 158,

422          v G20 1887, 56, 158,

423          v G21 1887, 56, 158,

424          v G22 1887, 56, 158,

425          v G23 1887, 56, 158,

426          v G24 1887, 56, 158,

427          v G25 1887, 56, 158,

428          v G26 1887, 56, 158,

429          v G27 1887, 56, 158,

430          v G28 1887, 56, 158,

431          v G29 1887, 56, 158,

432          v G30 1887, 56, 158,

433          v G31 1887, 56, 158,

434          v G32 1887, 56, 158,

435          v G33 1887, 56, 158,

436          v G34 1887, 56, 158,

437          v G35 1887, 56, 158,

438          v G36 1887, 56, 158,

439          v G37 1887, 56, 158,

440          v G38 1887, 56, 158,

441          v G39 1887, 56, 158,

442          v G40 1887, 56, 158,

443          v G41 1887, 56, 158,

444          v G42 1887, 56, 158,

445          v G43 1887, 56, 158,

446          v G44 1887, 56, 158,

447          v G45 1887, 56, 158,

448          v G46 1887, 56, 158,

449          v G47 1887, 56, 158,

450          v G48 1887, 56, 158,

451          v G49 1887, 56, 158,

452          v G50 1887, 56, 158,

453          v G51 1887, 56, 158,

454          v G52 1887, 56, 158,

455          v G53 1887, 56, 158,

456          v G54 1887, 56, 158,

457          v G55 1887, 56, 158,

458          v G56 1887, 56, 158,

459          v G57 1887, 56, 158,

460          v G58 1887, 56, 158,

461          v G59 1887, 56, 158,

462          v G60 1887, 56, 158,

463          v G61 1887, 56, 158,

464          v G62 1887, 56, 158,

465          v G63 1887, 56, 158,

466          v G64 1887, 56, 158,

467          v G65 1887, 56, 158,

468          v G66 1887, 56, 158,

469          v G67 1887, 56, 158,

470          v G68 1887, 56, 158,

471          v G69 1887, 56, 158,

472          v G70 1887, 56, 158,

473          v G71 1887, 56, 158,

474          v G72 1887, 56, 158,

475          v G73 1887, 56, 158,

476          v G74 1887, 56, 158,

477          v G75 1887, 56, 158,

478          v G76 1887, 56, 158,

479          v G77 1887, 56, 158,

480          v G78 1887, 56, 158,

481          v G79 1887, 56, 158,

482          v G80 1887, 56, 158,

483          v G81 1887, 56, 158,

484          v G82 1887, 56, 158,

485          v G83 1887, 56, 158,

486          v G84 1887, 56, 158,

487          v G85 1887, 56, 158,

488          v G86 1887, 56, 158,

489          v G87 1887, 56, 158,

490          v G88 1887, 56, 158,

491          v G89 1887, 56, 158,

492          v G90 1887, 56, 158,

493          v G91 1887, 56, 158,

494          v G92 1887, 56, 158,

495          v G93 1887, 56, 158,

496          v G94 1887, 56, 158,

497          v G95 1887, 56, 158,

498          v G96 1887, 56, 158,

499          v G97 1887, 56, 158,

500          v G98 1887, 56, 158,

501          v G99 1887, 56, 158,

502          v G100 1887, 56, 158,

503          v G101 1887, 56, 158,

504          v G102 1887, 56, 158,

505          v G103 1887, 56, 158,

506          v G104 1887, 56, 158,

507          v G105 1887, 56, 158,

508          v G106 1887, 56, 158,

509          v G107 1887, 56, 158,

510          v G108 1887, 56, 158,

511          v G109 1887, 56, 158,

512          v G110 1887, 56, 158,

513          v G111 1887, 56, 158,

514          v G112 1887, 56, 158,

515          v G113 1887, 56, 158,

516          v G114 1887, 56, 158,

517          v G115 1887, 56, 158,

518          v G116 1887, 56, 158,

519          v G117 1887, 56, 158,

520          v G118 1887, 56, 158,

521          v G119 1887, 56, 158,

522          v G120 1887, 56, 158,

523          v G121 1887, 56, 158,

524          v G122 1887, 56, 158,

525          v G123 1887, 56, 158,

526          v G124 1887, 56, 158,

527          v G125 1887, 56, 158,

528          v G126 1887, 56, 158,

529          v G127 1887, 56, 158,

530          v G128 1887, 56, 158,

531          v G129 1887, 56, 158,

532          v G130 1887, 56, 158,

533          v G131 1887, 56, 158,

534          v G132 1887, 56, 158,

535          v G133 1887, 56, 158,

536          v G134 1887, 56, 158,

537          v G135 1887, 56, 158,

538          v G136 1887, 56, 158,

539          v G137 1887, 56, 158,

540          v G138 1887, 56, 158,

541          v G139 1887, 56, 158,

542          v G140 1887, 56, 158,

543          v G141 1887, 56, 158,

544          v G142 1887, 56, 158,

545          v G143 1887, 56, 158,

546          v G144 1887, 56, 158,

547          v G145 1887, 56, 158,

548          v G146 1887, 56, 158,

549          v G147 1887, 56, 158,

550          v G148 1887, 56, 158,

551          v G149 1887, 56, 158,

552          v G150 1887, 56, 158,

553          v G151 1887, 56, 158,

554          v G152 1887, 56, 158,

555          v G153 1887, 56, 158,

556          v G154 1887, 56, 158,

557          v G155 1887, 56, 158,

558          v G156 1887, 56, 158,

559          v G157 1887, 56, 158,

560          v G158 1887, 56, 158,

561          v G159 1887, 56, 158,

562          v G160 1887, 56, 158,

563          v G161 1887, 56, 158,

564          v G162 1887, 56, 158,

565          v G163 1887, 56, 158,

566          v G164 1887, 56, 158,

567          v G165 1887, 56, 158,

568          v G166 1887, 56, 158,

569          v G167 1887, 56, 158,

570          v G168 1887, 56, 158,

571          v G169 1887, 56, 158,

572          v G170 1887, 56, 158,

573          v G171 1887, 56, 158,

574          v G172 1887, 56, 158,

575          v G173 1887, 56, 158,

576          v G174 1887, 56, 158,

577          v G175 1887, 56, 158,

578          v G176 1887, 56, 158,

579          v G177 1887, 56, 158,

580          v G178 1887, 56, 158,

581          v G179 1887, 56, 158,

582          v G180 1887, 56, 158,

583          v G181 1887, 56, 158,

584          v G182 1887, 56, 158,

585          v G183 1887, 56, 158,

586          v G184 1887, 56, 158,

587          v G185 1887, 56, 158,

588          v G186 1887, 56, 158,

589          v G187 1887, 56, 158,

590          v G188 1887, 56, 158,

591          v G189 1887, 56, 158,

592          v G190 1887, 56, 158,

593          v G191 1887, 56, 158,

594          v G192 1887, 56, 158,

595          v G193 1887, 56, 158,

596          v G194 1887, 56, 158,

597          v G195 1887, 56, 158,

598          v G196 1887, 56, 158,

599          v G197 1887, 56, 158,

600          v G198 1887, 56, 158,

601          v G199 1887, 56, 158,

602          v G200 1887, 56, 158,

603          v G201 1887, 56, 158,

604          v G202 1887, 56, 158,

605          v G203 1887, 56, 158,

606          v G204 1887, 56, 158,

607          v G205 1887, 56, 158,

608          v G206 1887, 56, 158,

609          v G207 1887, 56, 158,

610          v G208 1887, 56, 158,

611          v G209 1887, 56, 158,

612          v G210 1887, 56, 158,

613          v G211 1887, 56, 158,

614          v G212 1887, 56, 158,

615          v G213 1887, 56, 158,

616          v G214 1887, 56, 158,

617          v G215 1887, 56, 158,

618          v G216 1887, 56, 158,

619          v G217 1887, 56, 158,

620          v G218 1887, 56, 158,

621          v G219 1887, 56, 158,

622          v G220 1887, 56, 158,

623          v G221 1887, 56, 158,

624          v G222 1887, 56, 158,

625          v G223 1887, 56, 158,

626          v G224 1887, 56, 158,

627          v G225 1887, 56, 158,

628          v G226 1887, 56, 158,

629          v G227 1887, 56, 158,

630          v G228 1887, 56, 158,

631          v G229 1887, 56, 158,

632          v G230 1887, 56, 158,

633          v G231 1887, 56, 158,

634          v G232 1887, 56, 158,

635          v G233 1887, 56, 158,

636          v G234 1887, 56, 158,

637          v G235 1887, 56, 158,

638          v G236 1887, 56, 158,

639          v G237 1887, 56, 158,

640          v G238 1887, 56, 158,

641          v G239 1887, 56, 158,

642          v G240 1887, 56, 158,

643          v G241 1887, 56, 158,

644          v G242 1887, 56, 158,

645          v G243 1887, 56, 158,

646          v G244 1887, 56, 158,

647          v G245 1887, 56, 158,

648          v G246 1887, 56, 158,

649          v G247 1887, 56, 158,

650          v G248 1887, 56, 158,

651          v G249 1887, 56, 158,

652          v G250 1887, 56, 158,

653          v G251 1887, 56, 158,

654          v G252 1887, 56, 158,

655          v G253 1887, 56, 158,

656          v G254 1887, 56, 158,

657          v G255 1887, 56, 158,

658          v G256 1887, 56, 158,

659          v G257 1887, 56, 158,

660          v G258 1887, 56, 158,

661          v G259 1887, 56, 158,

662          v G260 1887, 56, 158,

663          v G261 1887, 56, 158,

664          v G262 1887, 56, 158,

665          v G263 1887, 56, 158,

666          v G264 1887, 56, 158,

667          v G265 1887, 56, 158,

668          v G266 1887, 56, 158,

669          v G267 1887, 56, 158,

670          v G268 1887, 56, 158,

671          v G269 1887, 56, 158,

672          v G270 1887, 56, 158,

673          v G271 1887, 56, 158,

674          v G272 1887, 56, 158,

675          v G273 1887, 56, 158,

676          v G274 1887, 56, 158,

677          v G275 1887, 56, 158,

678          v G276 1887, 56, 158,

679          v G277 1887, 56, 158,

680          v G278 1887, 56, 158,

681          v G279 1887, 56, 158,

682          v G280 1887, 56, 158,

683          v G281 1887, 56, 158,

684          v G282 1887, 56, 158,

685          v G283 1887, 56, 158,

686          v G284 1887, 56, 158,

687          v G285 1887, 56, 158,

688          v G286 1887, 56, 158,

689          v G287 1887, 56, 158,

690          v G288 1887, 56, 158,

691          v G289 1887, 56, 158,

692          v G290 1887, 56, 158,

693          v G291 1887, 56, 158,

694          v G292 1887, 56, 158,

695          v G293 1887, 56, 158,

696          v G294 1887, 56, 158,

697          v G295 1887, 56, 158,

698          v G296 1887, 56, 158,

699          v G297 1887, 56, 158,

700          v G298 1887, 56, 158,

701          v G299 1887, 56, 158,

702          v G300 1887, 56, 158,

703          v G301 1887, 56, 158,

704          v G302 1887, 56, 158,

705          v G303 1887, 56, 158,

706          v G304 1887, 56, 158,

707          v G305 1887, 56, 158,

708          v G306 1887, 56, 158,

709          v G307 1887, 56, 158,

710          v G308 1887, 56, 158,

711          v G309 1887, 56, 158,

712          v G310 1887, 56, 158,

713          v G311 1887, 56, 158,

714          v G312 1887, 56, 158,

715          v G313 1887, 56, 158,

716          v G314 1887, 56, 158,

717          v G315 1887, 56, 158,

718          v G316 1887, 56, 158,

719          v G317 1887, 56, 158,

720          v G318 1887, 56, 158,

721          v G319 1887, 56, 158,

722          v G320 1887, 56, 158,

723          v G321 1887, 56, 158,

724          v G322 1887, 56, 158,

725          v G323 1887, 56, 158,

726          v G324 1887, 56, 158,

727          v G325 1887, 56, 158,

728          v G326 1887, 56, 158,

729          v G327 1887, 56, 158,

730          v G328 1887, 56, 158,

731          v G329 1887, 56, 158,

732          v G330 1887, 56, 158,

733          v G331 1887, 56, 158,

734          v G332 1887, 56, 158,

735          v G333 1887, 56, 158,

736          v G334 1887, 56, 158,

737          v G335 1887, 56, 158,

738          v G336 1887, 56, 158,

739          v G337 1887, 56, 158,

740          v G338 1887, 56, 158,

741          v G339 1887, 56, 158,

742          v G340 1887, 56, 158,

743          v G341 1887, 56, 158,

744          v G342 1887, 56, 158,

745          v G343 1887, 56, 158,

746          v G344 1887, 56, 158,

747          v G345 1887, 56, 158,

748          v G346 1887, 56, 158,

749          v G347 1887, 56, 158,

750          v G348 1887, 56, 158,

751          v G349 1887, 56, 158,

752          v G350 1887, 56, 158,

753          v G351 1887, 56, 158,

754          v G352 1887, 56, 158,

755          v G353 1887, 56, 158,

756          v G354 1887, 56, 158,

757          v G355 1887, 56, 158,

758          v G356 1887, 56, 158,

759          v G357 1887, 56, 158,

760          v G358 1887, 56, 158,

761          v G359 1887, 56, 158,

762          v G360 1887, 56, 158,

763          v G361 1887, 56, 158,

764          v G362 1887, 56, 158,

765          v G363 1887, 56, 158,

766          v G364 1887, 56, 158,

767          v G365 1887, 56, 158,

768          v G366 1887, 56, 158,

769          v G367 1887, 56, 158,

770          v G368 1887

442 vergleich zu den Worksheet Datei aber relativ klein und das entpacken ällt daher nicht stark ins  
443 gewicht. Die sheetX.xml Datei lassen wir aber zunächst unverpackt und speichern im Arbeitsspe-  
444 icher lediglich die Verpackte Datei zwischen. Erst wenn eine beliebige Zelle in dem Worksheet von  
445 unserem Programm abgefragt wird, entweder weil der Benutzer den Worksheet durch einen Klick  
446 darauf zur anzeigen gebracht hat, oder weil ein Unitest die Werte einer Excel Datei überprüft, wird  
447 im Hintergrund auch erst die Archivdatei des Arbeitsblattes extrahier. Auf diese Weise wartet der  
448 Nutzer bei öffnen einer Excel Datei nur einen Bruchteil einer Sekunde, bevor er mit dem Programm  
449 weiterarbeiten kann wohingegen Excel bei außergewöhnlich großen Dateien einige Sekunden  
450 benötigen kann, bis das erste Arbeitsblatt angezeigt werden kann.

451 Öffnet der Benutzer ein anderes Arbeitsblatt, so wird auch dieses erst beim Zugriff auf das  
452 Arbeitsblatt gepasst. Diese Vorgehensweise ermöglicht es uns jede Excel Datei nahezu sofort  
453 dem Benutzer zur Anzeigen zu bringen und ihm keine Wartezeiten im Öffnen der Excel Dateien  
454 zuzumuten.

455

#### 456 4.4 Darstellung der Arbeitsblätter

457 Damit der Benutzer keine Schwierigkeiten hat, sich im Worksheet zurechtzufinden, legen wir  
458 großen wert darauf, die Arbeitsblätter so anzuzeigen, dass sie der Anzeigen in Excel sehr nahe  
459 kommen.

460

461 4.4.1 *Anzeigen der Höhe und Breite.* Unsere Lösung extrahiert die Informationen über die Spal-  
462 tenbreiten aus dem *width* Attribut der *col-knoten* und die Zeilenhöhen aus dem *ht* Attribut der *row*  
463 Knoten aus der sheetX.xml Dateien. Dabei ist darauf zu achten, dass in Excel die Spaltenbreiten und  
464 Zeilenhöhen in Excel Einheiten gemessen werden. Diese müssen mit einem Faktor mit multipliziert  
465 werden, um die tatsächliche Höhe oder Breite in Pixeln zu erhalten. Wichtig dabei zu betrachten  
466 ist, dass dieser Faktor für die Spaltenbreiten und die Zeilenhöhen jeweils unterschiedlich ist. Damit  
467 die Spaltenbreite genauso bereit erscheint, wie sie in Excel erscheint, muss die angegebene Einheit  
468 in der XML Datei mit dem Faktor 7 multipliziert werden. Um die tatsächliche Zeilenhöhe in Pixeln  
469 zu erreichen, muss dagegen der Faktor 4/3 mit dem Wert der in der Excel Datei angegeben wird,  
470 multipliziert werden. Die Werte haben wir durch eine Reihe von Tests ermittelt, indem wir Excel  
471 Dateien mit unterschiedlich großen Zellen erstellt haben und dann in der Worksheet Datei die  
472 Werte der Spaltenbreiten und Zeilenhöhen abgelesen haben. Nachdem wir nach einer Reihe solcher  
473 tests die Werte eine schätzung der Werte ermittelt haben, haben wir die Breiten und Höhen in den  
474 Arbeitsblätter modifiziert und die Excel Datei geöffnet um die Breiten und höhen zu messen und zu  
475 verifizieren, dass unsere Schätzung korrekt war.

476 Wie es zu den Unterschiedlichen Faktoren für die Zeilenhöhe und die Spaltenbreite kommt,  
477 konnten wir dagegen nicht durch Recherche herausfinden.

478 Durch das anwenden dieser Faktoren können wir die Spaltenbreiten und Zeilenhöhen so genau  
479 wie möglich darstellen, wie sie in Excel dargestellt werden.

480

481 4.4.2 *Formatierungen.* Die Formatierungen der Zellen werden in der styles.xml Datei gespe-  
482 ichert. In dieser Datei werden die Formatierungen in Form von Stilen gespeichert, die dann in der  
483 sheetX.xml Datei referenziert werden. Dieser Index verweist auf den entsprechenden *xf*-Knoten im  
484 *cellStyleXfs* Knoten. Diese *xf*-Knoten enthalten wiederum Indeces der Schriftart, Hintergrundfarbe  
485 und Border über die Attribute *fontId*, *fillId* und *borderId*, welche auf die entsprechenden *font*-Knoten  
486 im *fonts*-Knoten, die *fill*-Knoten im *fills*-Knoten und die *border*-Knoten im *borders*-Knoten ver-  
487 weisen. Diese Informationen nutzen wir, um die Zellen so genau wie möglich darzustellen, wie sie  
488 in Excel dargestellt werden.

489

491     4.4.3 *Overflow*. Genau wie in Excel auch achten wir darauf, dass Texte, die nicht in die Zelle  
 492 hineinpassen durch einen Overflow in die anliegenden Zellen hineinragen. Doch das geschieht nur,  
 493 wenn die anliegenden Zellen nicht selbst Werte enthalten. In dem Fall, dass die anliegenden Zellen  
 494 Werte enthalten, wird der Text abgeschnitten. Dieses verhalten armen wir nach um die visuelle  
 495 representation der Excel Datei so genau wie möglich nachzubilden.

496     Abbildung 8 zeigt, wie es aussieht, wenn man das nicht umsetzt. Der Text "Nursing Staff" ragt  
 497 in die Zelle hinein, die rechts daneben liegt. In Excel würde der Text in die Zelle hineinragen, die  
 498 rechts daneben liegt, wenn diese Zelle leer ist. In unserem Programm wird der Text abgeschnitten,  
 499 wenn die Zelle rechts daneben einen Wert enthält.

500  
 501  
 502 [To the table of c](#)  
 503  
 504

## 505 Supply of Nurs

| 508 Age from ... to<br>509 under ... Years | Nursing Staff |      |      |      |      |      |      |          |
|--------------------------------------------|---------------|------|------|------|------|------|------|----------|
|                                            | 510 Year      | 2024 | 2029 | 2034 | 2039 | 2044 | 2049 | 511 1867 |
| 512 Total                                  |               | 1673 | 1710 | 1738 | 1790 | 1839 | 1867 |          |

513 Fig. 8  
 514

## 515 4.5 Performance

516 Um hohe bildwiederholraten auch bei großen Arbeitsblätter zu gewährleisten, haben wir den  
 517 Spreadsheet Data extractor so optimiert, dass er nur die Zellen Zeichnet , die auch tatsächlich  
 518 sichtbar sind. Zu diesem zweck nutzen wird die seit Aug 17, 2023 erstmals erschienene package  
 519 two\_dimensional\_scrollables. [9].

520     Da wir alle Spaltenbreiten und Zeilen Höhen aus den xml Dateien extrahieren, können wir diese  
 521 Nutzen um zu berechnen, wie hoch und wie breit jede zelle ist und indem die breiten bzw höhen  
 522 aufaddiert werden können wir auch berechnen an welcher koordinate sich die Zelle befindet. Die  
 523 Koordinaten und breiten und höhen lkönnen dann genutzt werden um zu berechnen, welche Zellen  
 524 aktuell sichtbar sind und alle anderen Zellen beim Zeichnen ignorieren. Dazu wird das aktuelle  
 525 Scroll Offset in der x und y Achse berücksichtigt. Durch die Höhe und breite des Panels, welches das  
 526 Excel Arbeitsblatt anzeigt, wird ein viewport aufgespannt, über den berechnet werden kann, welche  
 527 Zellen darin sichtbar sind und welche außerhalb dieses viewports liegen. Welche Zellen aktuell  
 528 sichtbar sind Indem die Spaltenbreiten so lange addiert werden, bis die Summe die linke Kante des  
 529 Viewports erreicht. Alle Zellen links davon werden beim Zeichnen ignoriert. Weiterhin werden die  
 530 Spaltenbreite Breiten weiter addiert, bis sie die rechte Kante des Viewports erreicht.Alles Zellen  
 531 rechts davon werden beim Zeichnen ignoriert. Weiterhin werden die Zeilenhöhen entweder aus der  
 532 Standard Höhe.Des Arbeitsblattes oder über die Benutzer definierten Höhen der Zeilen extrahiert  
 533 und so lange aufaddiert, bis die Summe die obere Kante des Viewports erreicht.Alle Zellen, die  
 534 oberhalb dieser Kante liegen, werden beim Zeichnen ignoriert.Die Zeilenhöhen werden weiter  
 535 aufaddiert, bis sie die.Unsere Kante des Viewports erreicht.Nur diese Zellen werden gezeichnet.  
 536 Alle Zellen unterhalb dieser Kante werden ignoriert.  
 537

540 das two\_dimensional\_scrollables package bietet eine Funktion, die es erlaubt, die sichtbaren  
 541 Zellen zu berechnen und nur diese zu zeichnen. In dieser Funktion werden das horizontalOffset  
 542 und das verticalOffset und die viewportWidth und die viewportHeight als Parameter übergeben,  
 543 die die Position des Viewports in der x und y Achse beschreiben und die genutzt werden kann um  
 544 daran die sichtbaren Zellen zu berechnen. Der Algorithmus zum Berechnen der sichtbaren Zellen  
 545 ist in Algorithmus 1 dargestellt.

546

---

**Algorithm 1** Layout Spreadsheet Cells in Grid
 

---

548

```

1: Initialize Indices
2:   leadingColumnIndex  $\leftarrow$  column index based on horizontalOffset
3:   leadingRowIndex  $\leftarrow$  row index based on the verticalOffset
4:   trailingColumnIndex  $\leftarrow$  column index based on horizontalOffset + the viewportWidth
5:   trailingRowIndex  $\leftarrow$  row index based on the verticalOffset + the viewportHeight
6: Calculate Offsets
7:   leadingColumnOffset  $\leftarrow$  sum of widths from the first column to leadingColumnIndex
8:   leadingRowOffset  $\leftarrow$  sum of heights from the first row to leadingRowIndex
9:   horizontalLayoutOffset  $\leftarrow$  leadingColumnOffset - horizontalOffset
10: for each columnIndex from leadingColumnIndex to trailingColumnIndex do
11:   verticalLayoutOffset  $\leftarrow$  leadingRowOffset - verticalOffset
12:   for each rowIndex from leadingRowIndex to trailingRowIndex do
13:     child  $\leftarrow$  build the child for the columnIndex and rowIndex or obtain the cached one
14:     layout the child using the current horizontalLayoutOffset and verticalLayoutOffset
15:     if the row for rowIndex exists in the worksheet row definitions then
16:       verticalLayoutOffset  $\leftarrow$  verticalLayoutOffset + height of the row for rowIndex
17:     else
18:       verticalLayoutOffset  $\leftarrow$  verticalLayoutOffset + defaultRowHeight
19:     end if
20:   end for
21:   columnWidth  $\leftarrow$  width of the column for columnIndex
22:   horizontalLayoutOffset  $\leftarrow$  horizontalLayoutOffset + columnWidth
23: end for

```

---

572

573

## 5 EVALUATION

574

Der Ansatz des Spreadsheet Data Extractors unterscheidet sich auf dem Converter von Alexander Aue et al. auf dem wir aufbauen nicht verändert. Die Effektivität dieses Ansatzes wurde bereits untersucht. Sie haben die Extraktion von Daten aus über 500 Excel Dateien evaluiert. Die Zeit, die für jede Datei benötigt wurde, wurde aus einer Stichprobe von 331 verarbeiteten Excel Dateien bestimmt, die 3.093 Arbeitsblätter umfassen. Im Durchschnitt benötigten die studentischen Hilfskräfte 15 Minuten pro Datei und 95 Sekunden pro Arbeitsblatt.

581

582

583

584

585

586

587

Wir konzentrieren uns auf die Verbesserung der Benutzererfahrung und die Optimierung der Leistung des Spreadsheet Data Extractors. Die Benutzererfahrung durch die Darstellung der Excel-Arbeitsblätter ähnlich wie sie auch in Excel dargestellt werden vereinfacht und die Anzahl der nötigen Klicks durch die Vereinigung der Benutzeroberflächen für die Selektionshierarchie, dem Worksheet-View und der Vorschau der Ausgabe in einem View verringert. Die Leistung wurde durch das inkrementelle Laden der Excel Dateien und das nur Zeichnen der sichtbaren Zellen verbessert.

588

## 589    5.1 beschleunigung beim öffnen der Datei

590 Zum evaluieren der beschleunigung beim öffnen der Datei haben wir eine Reihe von Tests durchge-  
 591 führ. Wir haben den gesamten Datensatz an Exceldateien von Destatis heruntergeladen und die  
 592 größte Datei identifiziert. Mithilfe eines VBA Scriptes haben wir diese Exceldatei geöffnet und  
 593 Werte in einem Arbeitsblatt ausgelesen. Dieser Test wurde 10 mal wiederholt und die Zeit gemessen,  
 594 die benötigt wurde, um die Datei zu öffnen und die Werte auszulesen. Ein Equivalent dieses Scriptes  
 595 haben wir als Unit Test geschrieben welcher dieselbe Datei mit der für den Spreadsheet Data  
 596 Extractor implementierten Funktionen öffnet und dieselben Zellen aussliest Die Ergebnisse dieser  
 597 Laufzeiten sind in Abbildung putput zu sehen.

598 Der Spreadsheet Data Extractor öffnete das Arbeitsblatt im median in 120 Millisekunden, im  
 599 Durchschnitt 178 Millisekunden und im ersten durchlauf 668 Millisekunden, was durch sein  
 600 kann, weil die Datei noch nicht im Cache war und erst von der Festplatte geladen werden musste.  
 601

602 Excel öffnete das Arbeitsblatt im median in 40 und 281 Millisekunden, im Durchschnitt 41  
 Sekunden und 138 Millisekunden.

603 Das dart package brauchte im ersten durchlauf 13 Minuten und 15 Sekunden zum öffnen des  
 604 Arbeitsblattes. Die 9 anderen Durchläufe konnten nicht abgeschlossen werden da während des  
 605 zweiten durchlaufs eine Out of Memory Exception geworfen wurde.

606 Damit öffnete der Spreadsheet Data Extractor das Arbeitsblatt over two orders of magnitude  
 607 faster than Excel and nearly four orders of magnitude faster than the excel package used in prior  
 608 work.

609 Die Ergebnisse sind in 9 zu sehen.



610  
 611 Fig. 9  
 612

## 623    6 OUTLOOK

624 Wir haben vor den Spreadsheet Data Extractor weiter zu verbessern. Dazu gehört die Implemen-  
 625 tierung von Features, die die Benutzererfahrung weiter verbessern wie zum Beispiel die korrekte  
 626 Darstellung von Texten für die in Excel das horizontale Text Alignment "Center Across Selection"  
 627 verwendet wurde.

628 Paralel soll das Tool an weiteren Daten erprobt werden, um die Effektivität des Tools zu evaluieren.

### 631    6.1 Center Across Selection

632 Auch wenn wir uns mühe geben die Anzeige der Excel Dateien so genau wie möglich nachzubilden,  
 633 gibt es immer noch einige Unterschiede zwischen der Anzeige in Excel und der Anzeige in unserem  
 634 Spreadsheet Data Extractor. Die Unterschiede sind unbeabsichtigt und könnten behoben werden.  
 635 Einer solcher Unterschiede ist in Abbildung 10 zu sehen. In Excel werden die Text "Nursing Staff"  
 636 und "Year" über den gesamten Spaltenkopf hinweg zentriert angezeigt, während er in unserem  
 637

638 Spreadsheet Data Extractor linksbündig angezeigt wird. Dieser Unterschied ist darauf zurück-  
 639 zuführen, dass wir zwar merged cells parsen, die für gewöhnlich für das zentrieren solcher texte  
 640 über mehrere zellen hinweg verwendet werden, doch in diesem Arbeitsblatt wurden keine merged  
 641 cells verwendet um diese visuelle representation herzustellen. Stattdessen hat die erste Zelle in  
 642 der styles.xml datei zugeordneten xml knoten alignment für das attribut horizontal den Wert  
 643 centerContinuous. Bisherige Memühungen herraufzufinden über welche zellenkoordinaten sich  
 644 diese horizontale zentrierung erstreckt sind jedoch bisher gescheitert. Weitere tests oder recherche  
 645 ist notwendig um herraufzufinden, wie solche zentrierten Texte in der xml struktur der excel datei  
 646 gespeichert sind.

647  
 648 [To the table of contents](#)  
 649

## 650 651 Supply of Nursing Staff (Trend Variant) in Germany up to 2049, in 1000

| 653<br>654 Age from ...<br>655 to under ...<br>656 Years | Nursing Staff |      |      |      |      |      |
|----------------------------------------------------------|---------------|------|------|------|------|------|
|                                                          | Year          |      |      |      |      |      |
|                                                          | 2024          | 2029 | 2034 | 2039 | 2044 | 2049 |

657  
 658 Fig. 10  
 659  
 660

## 661 6.2 Evaluation an Real-World Daten

662 Bisher wurde die Basisversion des Tools an einem Datensatz des Agricultural Structure Survey on  
 663 land use and livestock in Germany für 2020 erprobt. Interessant wäre nun die neue Version des  
 664 Tools an den Daten vor 2020 zu testen, um zu sehen, ob die Verbesserungen, die wir vorgenommen  
 665 haben, auch die Effektivität des Tools verbessern. Dazu können wir auf die in den ausgabe csv  
 666 Dateien dokumentierten Timestamps zurückgreifen, um sie dann bei einem erneuten test mit den  
 667 neuen timestamps zu vergleichen um so herraufzufinden ob die studentischen Hilfskräfte schneller  
 668 mit dem neuen Tool arbeiten können.  
 669

## 670 7 CONCLUSION

671 In this paper, we presented the Spreadsheet Data Extractor, a tool that allows users to extract  
 672 relational data from complex Excel files. By adopting a user-centric approach that gives users full  
 673 control over data selection and metadata hierarchy definition, we provide a robust and accessible  
 674 solution for data extraction. Our tool offers user-friendly features such as the ability to duplicate  
 675 hierarchies of columns and tables, and to move them over similar structures for reuse, reducing  
 676 the need for repetitive configurations. We also contribute our tool under the GNU General Public  
 677 License v3.0, allowing the community to access, use, and improve the tool freely. By combining the  
 678 strengths of manual control with enhanced user interface features and performance optimizations,  
 679 our tool offers a valuable asset for efficient and reliable data extraction from diverse spreadsheet  
 680 formats.  
 681

## 682 REFERENCES

- 683 [1] Andrea Ackermann Alexander Aue. 2024. Converting data organised for visual perception into machine-readable  
 684 formats. In 44. GIL-Jahrestagung, *Biodiversität fördern durch digitale Landwirtschaft*. Gesellschaft für Informatik eV,  
 685 179–184.

- 687 [2] Daniel W Barowy, Sumit Gulwani, Ted Hart, and Benjamin Zorn. 2015. FlashRelate: extracting relational data from  
688 semi-structured spreadsheets using examples. *ACM SIGPLAN Notices* 50, 6 (2015), 218–228.
- 689 [3] D.J. Berndt, J.W. Fisher, A.R. Hevner, and J. Studnicki. 2001. Healthcare data warehousing and quality assurance.  
690 *Computer* 34, 12 (2001), 56–65. <https://doi.org/10.1109/2.970578>
- 691 [4] Zhe Chen, Michael Cafarella, Jun Chen, Daniel Prevo, and Junfeng Zhuang. 2013. Senbazuru: A prototype spreadsheet  
692 database management system. *Proceedings of the VLDB Endowment* 6, 12, 1202–1205.
- 693 [5] Haoyu Dong, Shijie Liu, Shi Han, Zhouyu Fu, and Dongmei Zhang. 2019. Tablesense: Spreadsheet table detection with  
694 convolutional neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 69–76.
- 695 [6] Angus Dunn. 2010. Spreadsheets-the Good, the Bad and the Downright Ugly. *arXiv preprint arXiv:1009.5705* (2010).
- 696 [7] Julian Eberius, Christoper Werner, Maik Thiele, Katrin Braunschweig, Lars Dannecker, and Wolfgang Lehner. 2013.  
697 DeExelerator: a framework for extracting relational data from partially structured documents. In *Proceedings of the  
698 22nd ACM international conference on Information & Knowledge Management*. 2477–2480.
- 699 [8] Elvis Koci, Dana Kuban, Nico Luettig, Dominik Olwig, Maik Thiele, Julius Gonsior, Wolfgang Lehner, and Oscar  
700 Romero. 2019. Xlindy: Interactive recognition and information extraction in spreadsheets. In *Proceedings of the ACM  
701 Symposium on Document Engineering 2019*. 1–4.
- 702 [9] Kate Lovett. 2023. two\_dimensional\_scrollables package - Commit 4c16f3e. <https://github.com/flutter/packages/commit/4c16f3ef40333aa0aebe8a1e46ef7b9fef9a1c1f> Accessed: 2023-08-17.
- 703 [10] Gursharan Singh, Leah Findlater, Kentaro Toyama, Scott Helmer, Rikin Gandhi, and Ravin Balakrishnan. 2009. Nu-  
704 metric paper forms for NGOs. In *2009 International Conference on Information and Communication Technologies and  
705 Development (ICTD)*. IEEE, 406–416.
- 706 [11] Statistisches Bundesamt (Destatis). 2024. *Statistischer Bericht - Pflegekraeftevorausberechnung - 2024 bis 2070*. Technical  
707 Report. Statistisches Bundesamt, Wiesbaden, Germany. <https://www.destatis.de/DE/Themen/Gesellschaft-Umwelt/Bevoelkerung/Bevoelkerungsvorausberechnung/Publikationen/Downloads-Vorausberechnung/statistischer-bericht-pflegekraeftevorausberechnung-2070-5124210249005.html> Statistical Report - Projection of Nursing Staff -  
708 2024 to 2070.
- 709 [12] Hannah West and Gina Green. 2008. Because excel will mind me! the state of constituent data management in small  
710 nonprofit organizations. In *Proceedings of the Fourteenth Americas Conference on Information Systems*. Association for  
Information Systems, AIS Electronic Library (AISeL). <https://aisel.aisnet.org/amcis2008/336>

711 Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735