

«Tables are tricky»

Testing Text Encoding Initiative (TEI) Guidelines for FAIR upcycling of digitised historical statistics

Gabi Wüthrich, Universitätsbibliothek Zürich



1. Project Background

Digitising Demographic Statistics of Zurich

- Willi Bretscher Fellowship Project at Zentralbibliothek (Central Library, ZB) Zurich by Dr. Joël Floris (2022): «Digitising demographic data from Zurich, 1910-1925, to quantify and contextualize the Spanish flu»
- Facsimile incl. OCR published on Zurich Open Platform (ZOP) by ZB
- This project: Pilot paper «From book to machine. XML as a platform-independent, machine-readable table format» as part of CAS UZH in Data Management and Information Technologies

Inspiration: Annual Fiscal Accounts of Basle

- Digital edition of the Basel annual accounts 1535-1611
- HTML, facsimile, and table view parallelly available
- Based on XML processing in accordance with TEI and RDF (Resource Description Framework) standard
- Early example for FAIR data edition



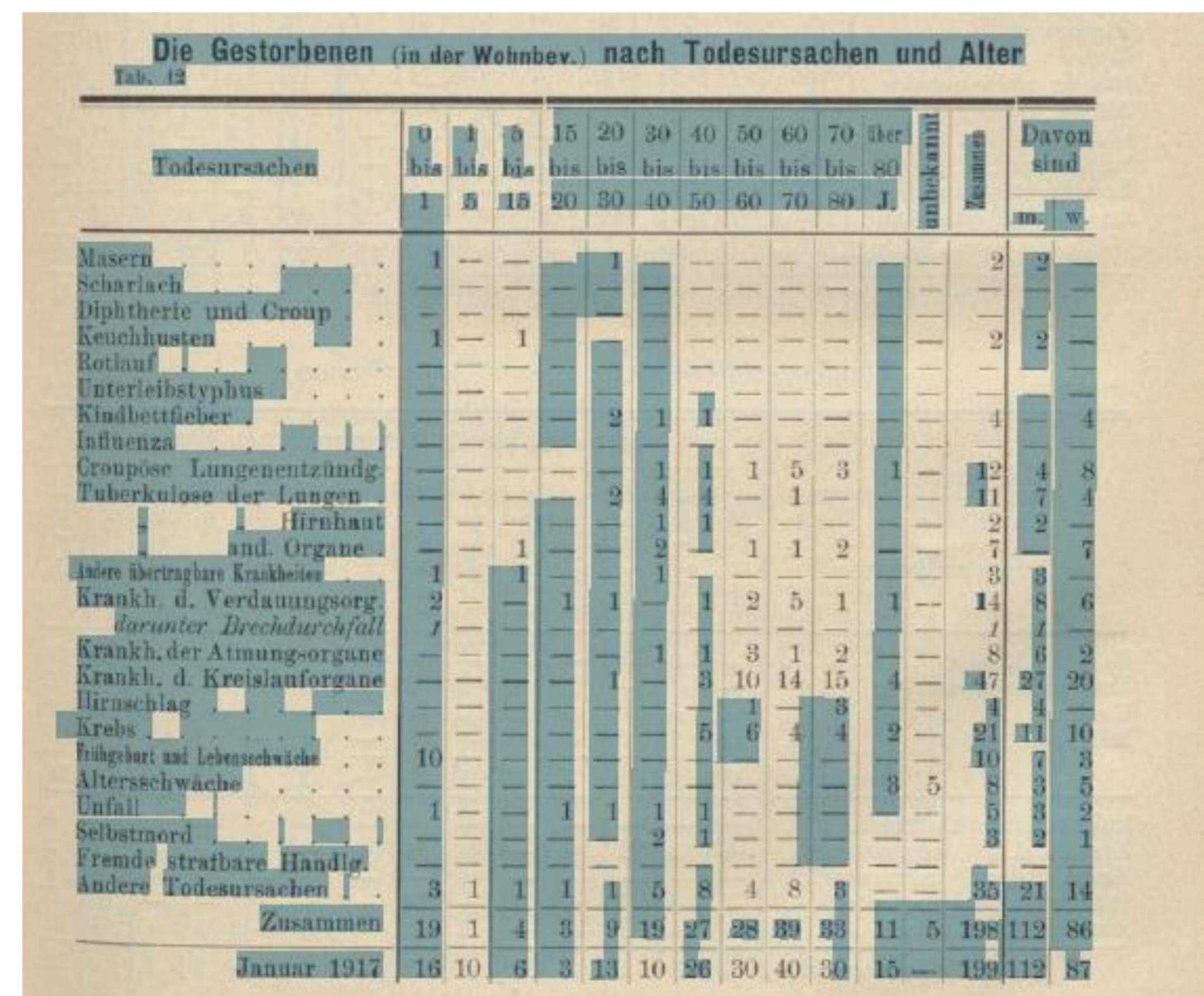
2. Text Encoding Initiative (TEI)

- Consortium developing and maintaining standards for the digital representation of texts since 1987
- Current version 4.9.0. of the guidelines published Jan 25th
- 24 chapters that explain TEI and define text elements (XML-based), including tables in chapter 15
- Structural elements: TEI root, TEI header, text body
- Advantages
 - Focus on the meaning of words and texts, not on layout
 - Software-independent
 - Supported by an open scientific community
- Disadvantages of tables
 - Layout and arrangement are more important than in continuous text, and convey meaning
 - Table processing in XML notations other than TEI



3. Challenges

- OCR recognition in digital copies must be adjusted for tables
- Image quality of digitisation needs to be sufficient
- OCR and TEI are flow text oriented
- Manual transfer of table to TEI structured XML content is prone to errors
 - Table recognition and XML coding are correspondingly deficient for structured text
 - Chicken-and-egg problem: More (XML) table templates for better text recognition supported by LLMs?



Tab.

Er iwu.

0,15 115[20]30|40|50 60 |70 jüber, 5| 5 Davon

Todesursachen bis bis bis bis bis nn bis bis 80& = sind

1 15/20|30|40|50 |5S|IE ee

Masern . See u 1 — 1

Scharlach : — | 1-1 1-2-[1212]2| — 1| —

Diphtherie und Croup . d: Hz _ -|1-1|-/-| . 2| 2 —

Keuchhusten . — | 1 — 1 — 1 —

Unterleibstypus ; — |-111- eis wur. —| 11 —

Kindbettfieber . | -112=+/-| -11 — 1

Influenza : Seesen 1 —| 21 -| 2

Croupöse Lungenentzündg. 1| 11773177 = 2 1| —| 9 3) 6

Tuberkulose der Lungen .|—| 1| 8 3| 2| 2) 4/—| ut 21| 13) 8

5 wsBrnaut Eee De 1 41.21 2

ee ieDel en sl 3 2

Anderer übertragbare Krankheiten 3|—| 11 —| 1| 1| —|=| 5 2 3

Krankh. d. Verdauungsorg. 3| —| 1 2/1| 3) 11 —| 1165

darunter Brechhäufigkeit 83| —| —| —| Ba een 8.267

Krankh. der Atmungsorgane 6| 2 —| 11-11 2| 4 —| 1 611

Krankh. d. BIER 1| —| 17.177174 [1578 ef ej7832| 151717

'Hirnschlag. ö | —| 1-1 1-12| 1.1) — Al- 2:52

Krebs .| —| 1| 12| 1) 54! 6) 7) —| 2724| °6) 18

Frühgeburt ni Tebenschwäche 10| —| eek er 201- 7,48

Altersschwäche .| —| 1 a 71| —| 22| 8) 14

Unfall | —| —| —| —| 6| 2 4

Selbstmord .| —| 1| 21| —| 11| —| 11| —| 8.21 211

Fremde strafbare Handig. er el est | ee en —

Anderer Todesursachen . DEab 1) 4| 4,2| 4| 4) —| 26| 12) 14

Zusammen 32| 2| 322 17 19 17.40 ai 10 — 210) 921118

Januar 1913 31| 2| 5| 14| 13| 29| 27| 31| 32| 9 — 200| 110| 90

Sum row

Label cells

Row D cells

Sum cells

Page number

7

Table head

Tab. 11

Die Gestorbenen nach dem Alter

Januar 1914 date

Altersjahre

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8

1 2 3 4 5 6 7 8</