



UNIVERSITÄT ZU LÜBECK

# Information Systems

CS4130-KP06

**Prof. Dr. Sylvia Melzer**

SoSe2026





# Data Models & Data Representation

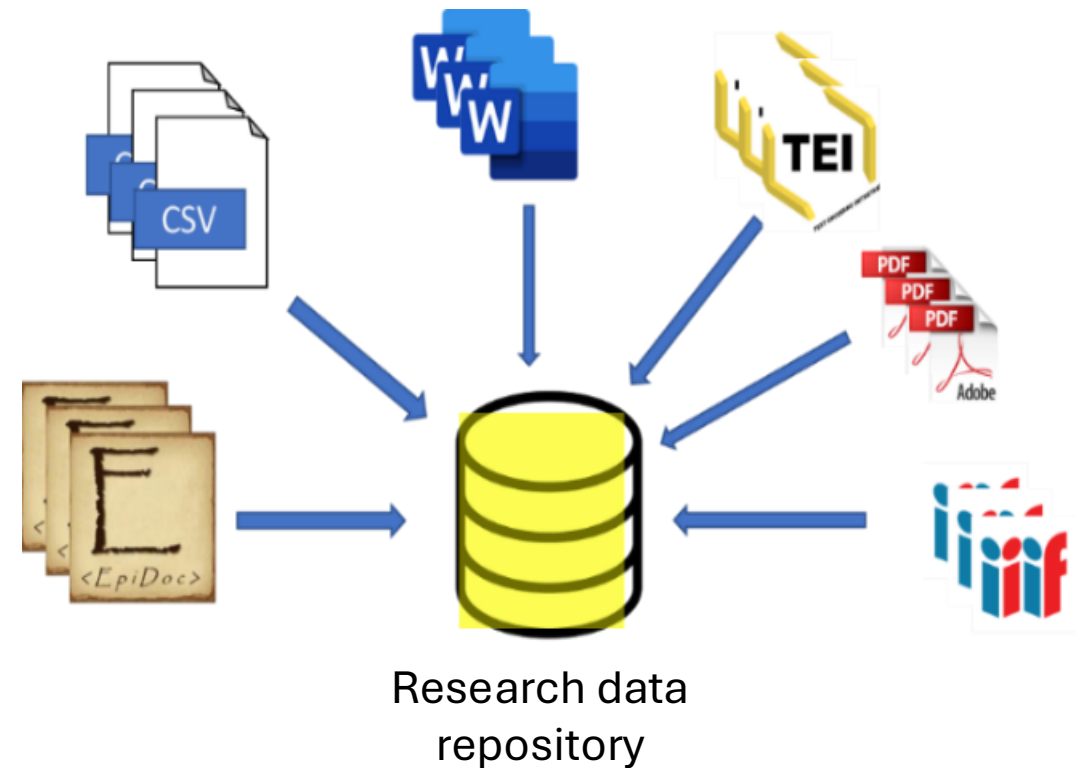
**Information Systems**

# What You Will Learn

- What is a data model and why it matters
- To explain differences between unstructured and structured data
- To recognize TEI and EpiDoc as modeling approaches
- To interpret how representation formats encode meaning
- To understand the role of schemas such as Relax NG
- To transform structured text into tabular formats
- To reflect on underlying data models

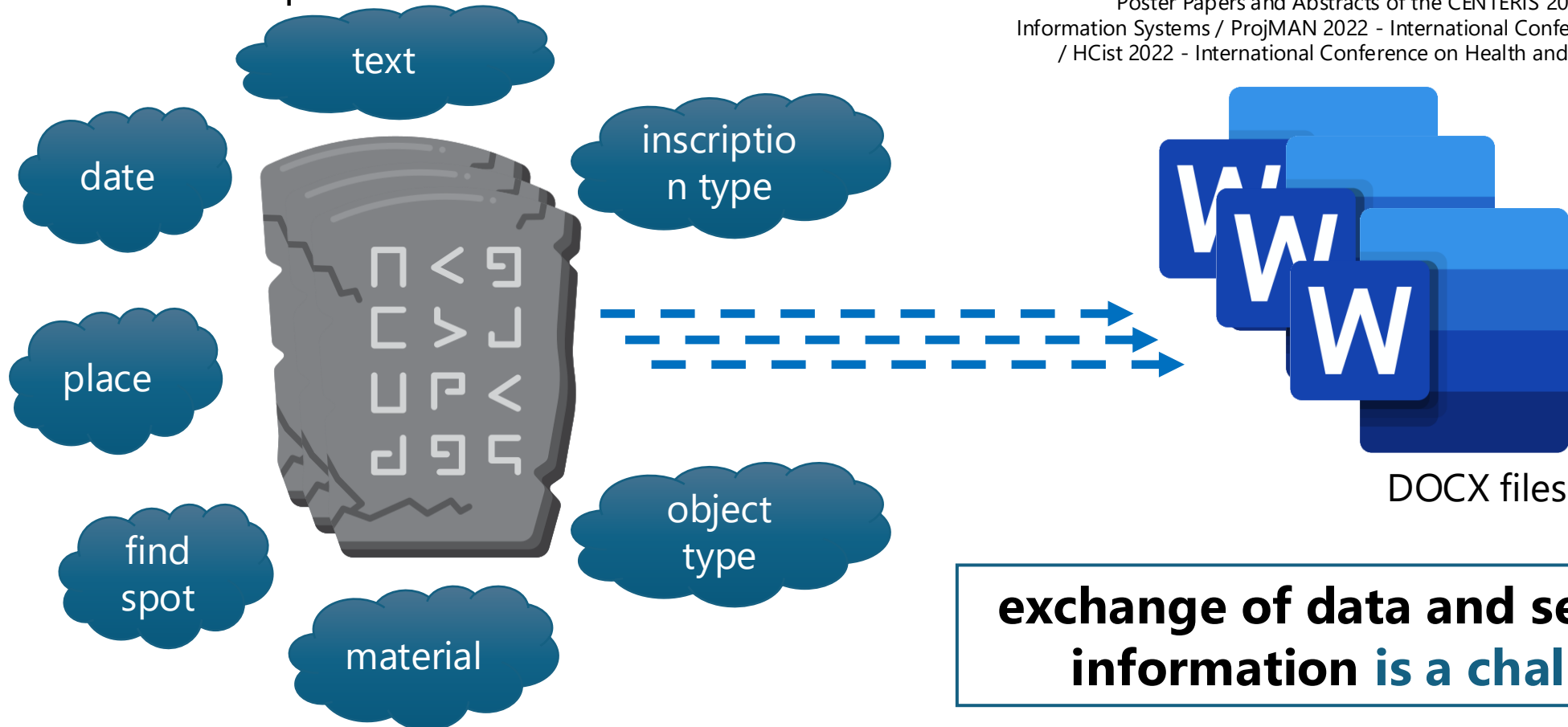
# The Nature of Research Data

- Research data often originates from heterogeneous sources
- It includes texts, images, measurements, and annotations
- Data is frequently created in discipline-specific formats
- Many datasets are not originally designed for reuse
- The structure of data is often implicit rather than explicit
- This makes integration and comparison difficult
- Therefore, additional processing is required to make data usable



# DOCX in Humanities Research

- Focus on presentation



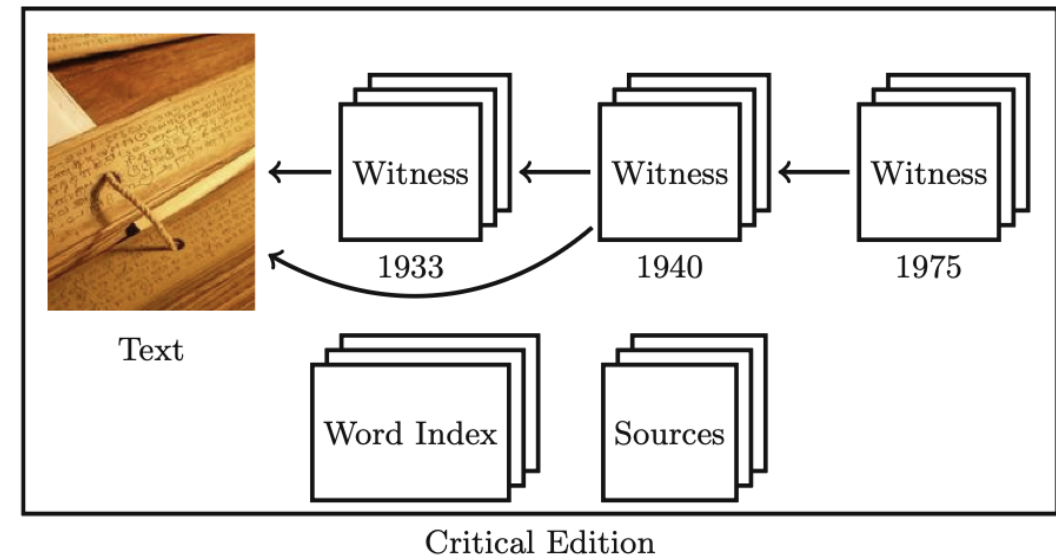
Sylvia Melzer, S. Schiff, F. Weise, K. Harter, and R. Möller, Databasing on demand for research data repositories explained with a large epidoc dataset, Book of Industry Papers, Poster Papers and Abstracts of the CENTERIS 2022 - Conference on ENTERprise Information Systems / ProjMAN 2022 - International Conference on Project MANagement / HCist 2022 - International Conference on Health and Social Care Info , pp. 150--153, 2022. SciKA, Portugal.

# Problem with DOCX

- Focus on visual presentation (WYSIWYG) rather than structure
- Lack of semantic encoding (e.g., persons, places, variants not explicitly defined)
- Data stored in different formats and tools across researchers
- Merging and harmonizing data is time-consuming and error-prone
- Creation of information based on research data (e.g., critical editions) takes years and involves high manual effort
- Difficult to produce digital, FAIR-compliant editions in addition to print
- Limited support for:
  - Automated processing
  - Large-scale analysis
  - Advanced search (e.g., faceted search)
- Result: Data exchange and information retrieval become a major challenge

# Critical Editions as Example

- Critical editions are created over many years of research
- They combine multiple textual witnesses and sources
- Scholars work with different tools and formats
- Data integration becomes complex and time-consuming
- The final output is often a printed edition
- Digital reuse is not inherently supported
- This highlights the need for structured representations



# Need for Digital and FAIR Data

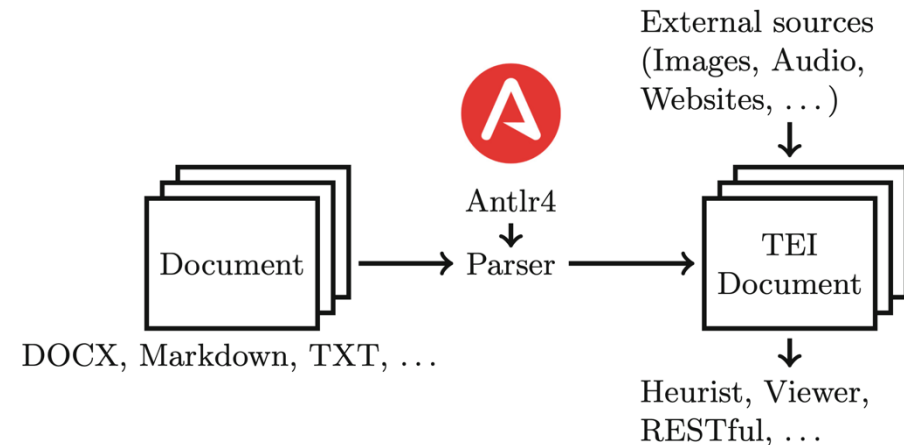
- There is growing demand for digital editions
- Digital data must be searchable and accessible
- Linking between data sources becomes important
- Data should follow FAIR principles
- Reusability is essential for future research
- Machine-readable formats enable automation
- Structured data supports advanced analysis methods

# The EASE-Oriented Approach

- **E**xplorable, **A**ccessible in terms of data quality, **S**eeable in operation in new contexts, and **E**asily checkable for reuse (EASE)
- The EASE-oriented approach focuses on supporting existing research practices
- It does not require scholars to abandon their preferred tools
- Instead, it introduces lightweight integration and transformation mechanisms, the goal is to reduce technical barriers e.g., for humanities researchers
- Data is transformed in the background without disrupting workflows
- This enables gradual adoption of structured data practices
- The approach emphasizes usability, flexibility, and sustainability

# EASE in Practice (Project: Netamil)

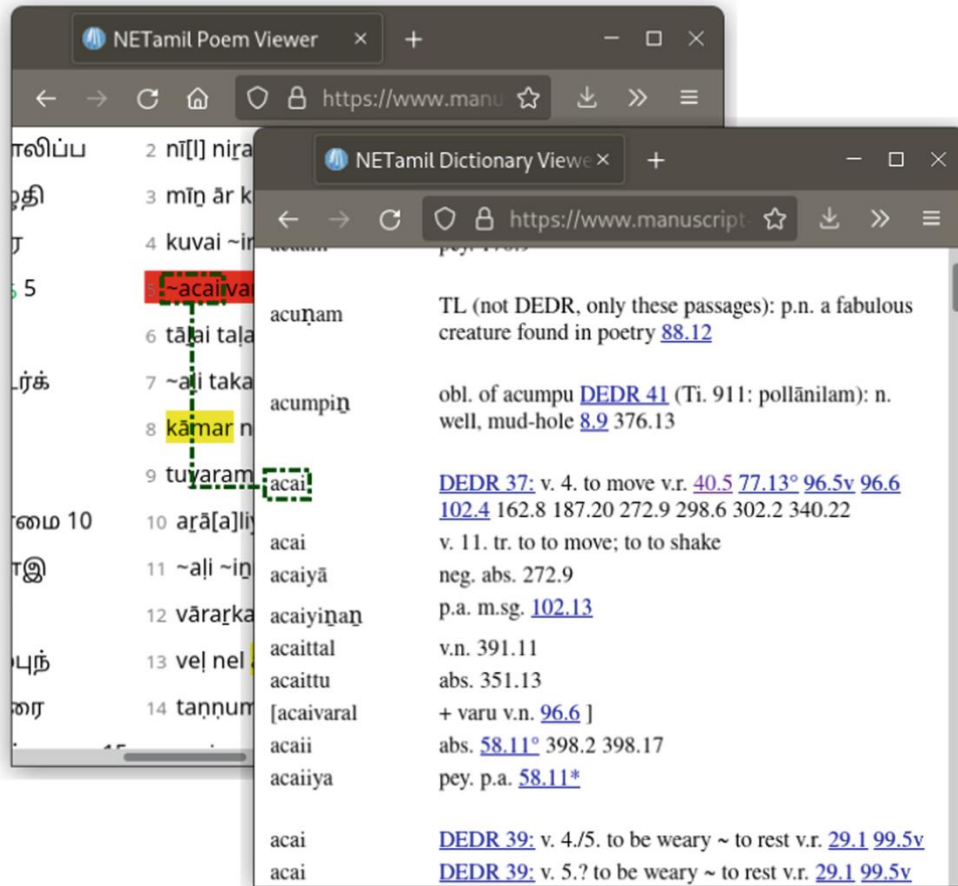
- Documents are created in familiar environments such as Word
- A transformation layer converts them into TEI automatically
- The transformation can be supported by parsers and rules
- Structured data is stored in repositories
- Data can be enriched with external resources
- Outputs can be generated for print and digital formats
- The process remains transparent to the researcher



# EASE and Data Modeling

- EASE enables the transition from implicit to explicit data models (e.g., DOCX → TEI)
- TEI acts as the intermediate structured model
- Data models emerge from transformation processes
- Researchers can analyze and refine these models
- The approach supports iterative improvement of data quality
- It connects representation with modeling practice.

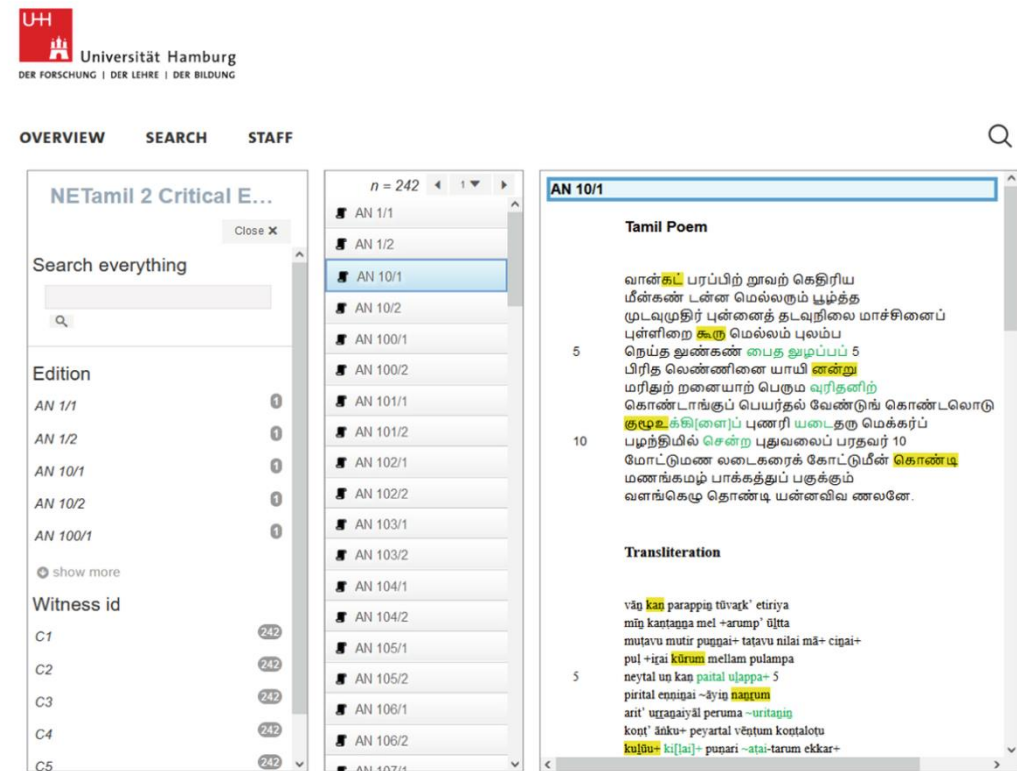
# EASE and Data Modeling



NETamil Poem Viewer

NETamil Dictionary Viewer

acai	TL (not DEDR, only these passages): p.n. a fabulous creature found in poetry <a href="#">88.12</a>
acai	obl. of acumpu <a href="#">DEDR 41</a> (Ti. 911: pollānilam): n. well, mud-hole <a href="#">8.9</a> 376.13
acai	<a href="#">DEDR 37</a> : v. 4. to move v.r. <a href="#">40.5</a> <a href="#">77.13°</a> <a href="#">96.5v</a> <a href="#">96.6</a> <a href="#">102.4</a> 162.8 187.20 272.9 298.6 302.2 340.22
acai	v. 11. tr. to to move; to to shake
acai	neg. abs. 272.9
acai	p.a. m.sg. <a href="#">102.13</a>
acai	v.n. 391.11
acai	abs. 351.13
acai	+ varu v.n. <a href="#">96.6</a> ]
acai	abs. <a href="#">58.11°</a> 398.2 398.17
acai	pey. p.a. <a href="#">58.11*</a>
acai	<a href="#">DEDR 39</a> : v. 4./5. to be weary ~ to rest v.r. <a href="#">29.1</a> <a href="#">99.5v</a>
acai	<a href="#">DEDR 39</a> : v. 5.? to be weary ~ to rest v.r. <a href="#">29.1</a> <a href="#">99.5v</a>



Universität Hamburg  
DER FORSCHUNG | DER LEHRE | DER BILDUNG

OVERVIEW SEARCH STAFF

NETamil 2 Critical E...

Search everything

Edition

- AN 1/1
- AN 1/2
- AN 10/1
- AN 10/2
- AN 100/1
- AN 100/2
- AN 101/1
- AN 101/2
- AN 102/1
- AN 102/2
- AN 103/1
- AN 103/2
- AN 104/1
- AN 104/2
- AN 105/1
- AN 105/2
- AN 106/1
- AN 106/2
- AN 107/1

Witness id

- C1 (242)
- C2 (242)
- C3 (242)
- C4 (242)
- C5 (242)

AN 10/1

Tamil Poem

வாண்குட்பரப்பிற் றாவற் கெதிரிய  
மீன்கண் டன்ன மெல்லும் பூத்த  
முடவுமுதிர் புண்ணைத் தடவுநிலை மாச்சினைப்  
புள்ளிறை **கூறு** மெல்லம் புலம்ப  
நெய்த துண்கண் பைத **ஆழப்ப** 5  
பிரித லெண்ணினை யாயி **என்று**  
மரிதற் றனையாற் பெரும் **வர்தனிற்**  
கொண்டாங்குப் பெயர்தல் வேண்டுங் கொண்டலொடு  
**குடிஉக்கிளை**ப் புணரி யடைதரு மெக்கர்ப்  
பழநிமில் **சென்று** புதுவலைப் பரதவர் 10  
மோட்டுமண லடைகரைக் கோட்டுமீன் **கொண்டி**  
மணங்கமம் பாக்கத்துப் பகுக்கும்  
வளங்கெழு தொண்டி யன்னவிவ ணலனே.

Transliteration

vāṅ **kaṅ** parappiṅ tūvaṅk' eṭriya  
mīṅ kaṅṅga mel + arump' ūṭṭa  
mutavu mutir puṅgai+ taṭavu nilai mā+ ciṅai+  
pu + iṅai **kūrum** mellam pulampa  
neytal un kaṅ **paṅtal uṅṅpa** 5  
pirital enniṅai ~ṅiṅ **naṅṅum**  
ariṅ' uṅṅaiyāḷ peruma ~uritaṅṅi  
kont' āṅku+ peyartal vēṅṅum koṅtalotu  
**kuṅṅu**+ **ku[la]**+ puṅari -atai-tarum ekkar+

# Key Insight of EASE

The goal is not to replace tools, but to extend them.

Structure is added without changing existing workflows.

Data becomes machine-readable without extra burden.

Transformation enables interoperability and reuse.

Modeling becomes part of everyday research practice.

This lowers the barrier to digital methods.

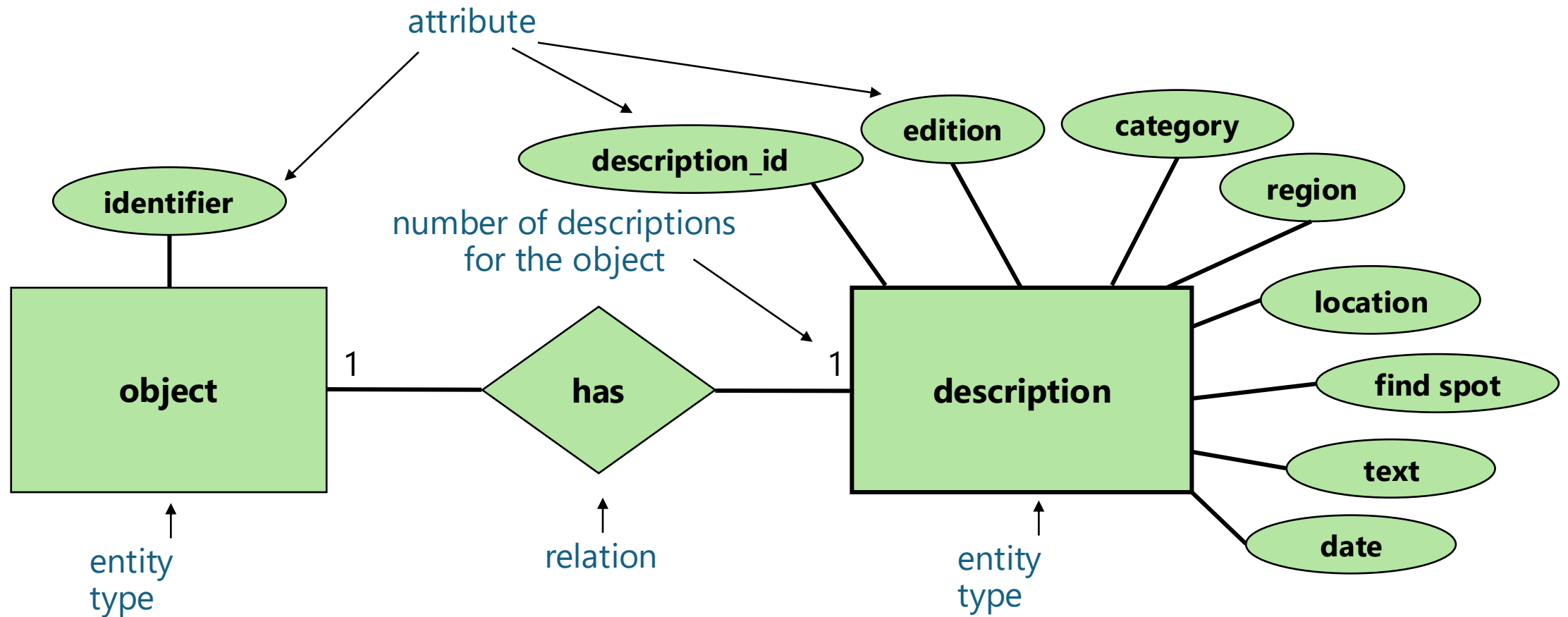
EASE enables sustainable research data ecosystems.

# Representation = Modeling

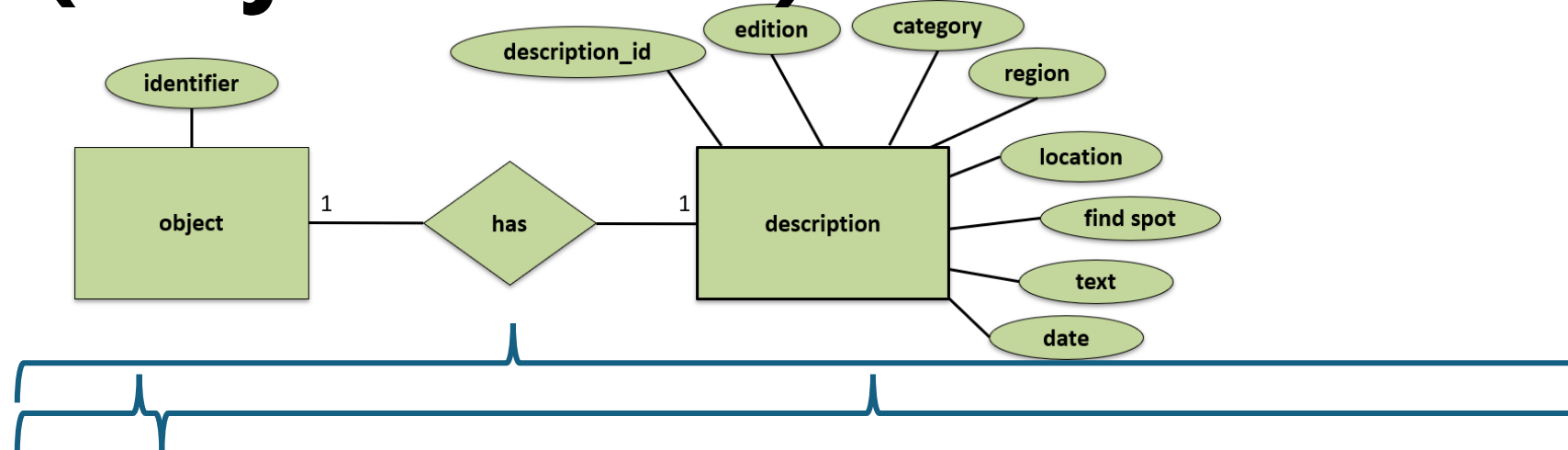
- Every data representation encodes a data model
- A data model is an abstract representation of data
- A format defines what can be expressed in data
- It determines how information is structured
- It also defines relationships between elements
- Representation is therefore not neutral
- It reflects assumptions about the data
- Understanding this is central to data modeling

```
—<ab>  
  <lb n="1"/>  
  Πλανκίαν Μάγναν  
  <lb n="2"/>  
  Άκυλλίαν θυγατέρα Ίου  
  <lb n="3" break="no"/>  
  λίου Σεουήρου και  
—<expan>  
  <abbr>Κλαυ</abbr>  
  <ex>δίας</ex>  
</expan>  
<lb n="4"/>  
  Άκυλλίας, τὴν ἐγ βασιλέ  
  <lb n="5" break="no"/>  
  ων, ἡρωίδα, ὁ δῆμος ὁ Ἄν  
  <lb n="6" break="no"/>  
  κυρανῶν τὴν θυγατέρα τῆς  
  <lb n="7"/>  
  μητροπόλεως.  
</ab>
```

# Data Model (Project: EDAK)



# Data Model (Project: EDAK)



- The table represents the **has**-relation between the **object** and the **description**
- The attributes are the column names
- A data set is in one row

Identifier	Description_id	Edition	Category	Region	Location	Find spot	Text	Date
object_id	00002670	MAMA I, Nr. 319	Epigramm	Lykaonien	Laodikeia	Gözlü	Ἀντίοκος {Ἀντίοχος} τόδε σῆμα φί- λῳ περικαλλέει πεδὶ τε[ύ]- ξεν Ἀρμένιου πόθῳ [χ]άριν θέτο τ[ὺ]δ' ἐπὶ τύμβῳ]. ἀντὶ φιλοστοργ[ί]ης [τ]ε τίτλο[ν] ἀμφὶ δ' ἄρ' αὐτῷ πολλὰ κινυρ[ό]- μενος κατηφιάας κέ προ[—]	

Data from <https://www.epigraphik.uni-hamburg.de>

# Data Model

- Data models ensure consistency across datasets
- They reduce redundancy and duplication
- They support efficient data retrieval
- They improve data quality and reliability
- They enable interoperability between systems
- They make data easier to understand
- They are essential for scalable analysis

# Structured vs Unstructured Data

- Unstructured data consists of free-form text or media
- Structured data follows predefined schemas
- Semi-structured data lies between both extremes
- Structured data is easier to query and analyze
- Unstructured data requires preprocessing
- Representation formats influence structure
- Modeling transforms unstructured into structured data

```

<text>
  <body>
    <div n="AN_1/1" type="tamil" xml:lang="ta">
      <head>Tamil Poem</head>
      <l n="1">வண்டுபடத் ததைந்தகண்ணியொண்கழ</l>
      <l n="2">ஐருவக் குதினாமழவிரோட்டிய</l>
      <l n="3">முருக னற்போர்நெடுவேளாவி</l>
      <l n="4">பறுகேள் டியானைப்<span style="color:#00B050;">
        பொதுளிய</span>யாங்க</l>
      <l n="5">சிறுகா ரோடன்பயினொடுசேர்த்திய</l>
      <l n="6">கற்போற் பிரியலமென்றசொற்றா</l>
      <l n="7">மறந்தனர் கொல்வேதோழிசிறந்த</l>
      <l n="8">வேய்மருள் பணைத்தோ<span style="color:#00B050;">

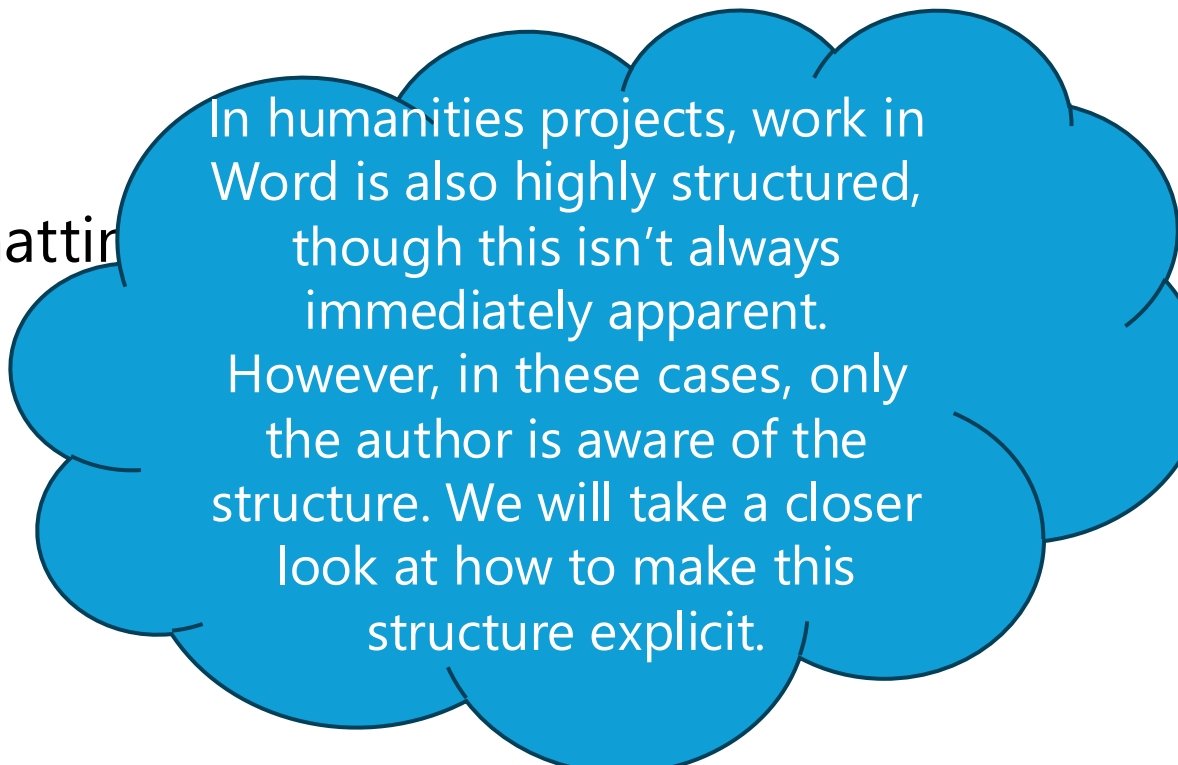
```



Palm leaves with text in Sinhalese from Sri Lanka (ca. mid-20th century). Private collection.

# Case Study – DOCX (Weak Model)

- DOCX focuses on visual layout and formatting
- Structure is implicit rather than explicit
- Formatting conveys meaning indirectly
- Machines cannot reliably interpret formatting
- Data extraction becomes complex
- There is no enforced structure
- This results in weak data modeling



In humanities projects, work in Word is also highly structured, though this isn't always immediately apparent. However, in these cases, only the author is aware of the structure. We will take a closer look at how to make this structure explicit.

# Implicit Data Model

- Entities are not formally defined
- Relationships are not explicitly encoded
- Data meaning depends on human interpretation
- Structure varies across documents
- No standardization is enforced
- Data integration is difficult
- This limits reuse and analysis

As a current research question to be answered here is *why there are some bronze statues in Beijing whose origin can be attributed to the royal family Paolola Šāhi*. For this purpose, it is necessary to evaluate the inscriptions of the above-mentioned Buddhist bronze statues e.g. from the field of epigraphy and philology. Epigraphy is the study of inscriptions to e.g. "clarify the meanings, classify their uses according to dates and cultural contexts, and drawing conclusions about the writing and writers." [6] The Buddhist bronze statues have inscriptions, which are written in Sanskrit language by using two different types of the so-called "Gandhāra-Brāhmī" handwriting: the "round type" is from the 2nd century to the 630 AD, and the "rectangular type" form is from 630 to the 8th century. The study of the written language is called philology. The aim of this study is to determine the meaning of inscriptions. If the history in the 4th-6th centuries and the Buddhist bronze statue inscriptions are studied more closely, it is found that even the Tibetans consider the statues to be holy, but could no longer read and understand the writing. Therefore, in this paper we present the requirements to be solved by using AI methods to answer the research question by linking different research data sources.



# Consequences

- Data processing requires manual effort
  - Automated workflows are limited
  - Errors occur during data integration
  - Redundancies are common
  - Queries cannot be performed efficiently
  - Large-scale analysis is impractical
- **This motivates structured approaches**

## Case Study – TEI (Explicit Model)

- The Text Encoding Initiative (TEI) is a collaborative consortium focused on developing and maintaining a standard for digitally representing texts
- Its primary output is a set of Guidelines outlining encoding methods primarily used in the humanities, social sciences, and linguistics
- Since 1994, these Guidelines have been widely adopted by libraries, museums, publishers, and scholars for presenting texts online, supporting research, education, and preservation efforts
- In addition to the Guidelines, the TEI Consortium offers various resources, training events, project showcases, and a bibliography to support learning and usage of TEI
- Guidelines: <https://tei-c.org/release/doc/tei-p5-doc/en/html/index.html>



# TEI: Overview

## Front Matter

### Title

- i. [Releases of the TEI Guidelines](#)
- ii. [Dedication](#)
- iii. [Preface and Acknowledgments](#)
- ⊕ iv. [About These Guidelines](#)
- ⊕ v. [A Gentle Introduction to XML](#)
- ⊕ vi. [Languages and Character Sets](#)

## Back Matter

- ⊕ Appendix A [Model Classes](#)
- ⊕ Appendix B [Attribute Classes](#)
- ⊕ Appendix C [Elements](#)
- ⊕ Appendix D [Attributes](#)
- ⊕ Appendix E [Datatypes and Other Macros](#)
- ⊕ Appendix F [Bibliography](#)
- ⊕ Appendix G [Deprecations](#)
- ⊕ Appendix H [Prefatory Notes](#)
- Appendix I [Colophon](#)

## Text Body

- ⊕ 1 [The TEI Infrastructure](#)
- ⊕ 2 [The TEI Header](#)
- ⊕ 3 [Elements Available in All TEI Documents](#)
- ⊕ 4 [Default Text Structure](#)
- ⊕ 5 [Characters, Glyphs, and Writing Modes](#)
- ⊕ 6 [Verse](#)
- ⊕ 7 [Performance Texts](#)
- ⊕ 8 [Transcriptions of Speech](#)
- ⊕ 9 [Computer-mediated Communication](#)
- ⊕ 10 [Dictionaries](#)
- ⊕ 11 [Manuscript Description](#)
- ⊕ 12 [Representation of Primary Sources](#)
- ⊕ 13 [Critical Apparatus](#)
- ⊕ 14 [Names, Dates, People, and Places](#)
- ⊕ 15 [Tables, Formulæ, Graphics, and Notated Music](#)
- ⊕ 16 [Language Corpora](#)
- ⊕ 17 [Linking, Segmentation, and Alignment](#)
- ⊕ 18 [Simple Analytic Mechanisms](#)
- ⊕ 19 [Feature Structures](#)
- ⊕ 20 [Graphs, Networks, and Trees](#)
- ⊕ 21 [Non-hierarchical Structures](#)
- ⊕ 22 [Certainty, Precision, and Responsibility](#)
- ⊕ 23 [Documentation Elements](#)
- ⊕ 24 [Using the TEI](#)



# TEI: Header

```
<teiHeader>
  <fileDesc>
    <titleStmt>
      <title>
        <!-- title of the resource -->
      </title>
    </titleStmt>
    <publicationStmt>
      <p>
        <!-- Information about distribution of the resource -->
      </p>
    </publicationStmt>
    <sourceDesc>
      <p>
        <!-- Information about source from which the resource derives -->
      </p>
    </sourceDesc>
  </fileDesc>
</teiHeader>
```

# TEI: Dictionary

```

<body>
<div>
<head>English-French</head>
<entry>
<form>
<orth>cat</orth>
</form>
<!-- ... -->
</entry>
<entry>
<form>
<orth>dog</orth>
</form>
<!-- ... -->
</entry>
<entry>
<form>
<orth>horse</orth>
</form>
<!-- ... -->
</entry>
</div>
<div>
<head>French-English</head>
<entry>
<form>
<orth>chat</orth>
</form>
<!-- ... -->
</entry>
<entry>
<form>
<orth>chien</orth>
</form>
<!-- ... -->
</entry>
<entry>
<form>
<orth>cheval</orth>
</form>
<!-- ... -->
</entry>
</div>
</body>
  
```

## Example: Netamil project

```

1 <l n="1">
2   <entry>
3     <form>
4       <orth>acai</orth>
5     </form>
6     <def>DEDR 37: v. 4. to move v.r. 40.5 77.13° 96.5v 96.6 102.4 162.8 187.20
7       272.9 298.6 302.2 340.22</def>
8   </entry>
9 </l>
  
```

- DEDR 37 refers to entry 37 in the “A Dravidian etymological dictionary” (DEDR: <https://dsal.uchicago.edu/dictionaries/burrow/>)
- v. 4. indicates the 4th verb in DEDR
- to move is the English translation
- v.r. means verbal root
- series of numbers: 40.5, 77.13°, 96.5v, 96.6, 102.4, 162.8, 187.20, 272.9, 298.6, 302.2, 340.22 shows in which poems (number before .) the word occurs and in which line (number after .)
- “o” means primary variant second strand and “v” means secondary variant

# TEI

## TEI as Data Model

- TEI encodes entities explicitly
- It defines relationships within texts
- It structures documents hierarchically
- It integrates “metadata” and content
- It enables machine-readable representation
- It supports consistent data encoding
- It acts as a formal data model

## TEI Structure

- TEI documents have a standardized structure
- The `teiHeader` contains metadata information
- The `text` element contains the content
- Elements define structural units
- Attributes provide additional detail
- Hierarchies reflect document organization
- This structure enables processing

# TEI

## Advantages

- TEI produces machine-readable data
- It supports interoperability between systems
- It enables long-term data preservation
- It facilitates data reuse
- It supports automated processing
- It enables advanced querying
- It improves research efficiency

## Disadvantages

- TEI documents can be complex and difficult to understand for non-experts
- The flexibility of TEI may lead to inconsistent encoding across projects
- Interoperability can be limited when projects use different TEI customizations
- TEI files are not always easily discoverable without proper metadata infrastructure
- Machine readability is high, but practical reuse often requires additional processing
- Long-term accessibility depends on tooling, documentation, and community support

# EpiDoc (Epigraphic Documents in TEI XML)

- EpiDoc is a collaborative international initiative that offers guidelines and tools for encoding scholarly and educational editions of ancient texts
- It employs a specific subset of TEI's standards for digitally representing texts and was originally designed for the publication of digital editions of ancient inscriptions
- Its application has broadened to encompass the publication of papyri and manuscripts, as demonstrated by platforms like Papyri.info
- EpiDoc covers not only the transcription and editorial aspects of the texts but also the historical context and materials of the artifacts on which these texts are found, including manuscripts, monuments, tablets, papyri, and other objects containing text
- Guidelines: <https://epidoc.stoa.org/gl/latest/>

# Case Study – EpiDoc (Explicit Model) – Project EDAK

```

1 <TEI xml:lang="en">
2   <teiHeader>
3     <fileDesc>
4       <titleStmt>
5         <title>SEG 6, 29</title>
6       </titleStmt>
7       <publicationStmt>
8         <authority>Epigraphische Datenbank
          zum antiken Kleinasien –
          Universität Hamburg</authority>
9         <idno type="localID">EDAK00001210
10        </idno>
11        <idno type="filename">EDAK00001210
12        </idno>
13        <availability>[...]</availability>
14      </publicationStmt>
15      <sourceDesc>
16        <msDesc>
17          <msIdentifier/>
18          <msContents>[...]</msContents>
19          <physDesc>[...]</physDesc>
20          <history>
21            <origin>
22              <origPlace>
23                <settlement ref="
  
```

```

1 <body>
2   <div type="edition" xml:lang="gr">
3     <ab>
4       <lb n="1"/><supplied reason="lost"> </
          supplied>
5       <lb n="2" break="no"/> <expan>
          <abbr>Θ<ex> </ex> </abbr></
          expan>
6       <lb n="3"/>Θ
7       <lb n="4" break="no"/> ,<expan><
          abbr> </abbr><ex> </ex></
          expan> ',
8       <lb n="5"/><expan><abbr> </abbr><
          ex> </ex></expan>M '.
9     </ab>
10    </div>
11    <div type="commentary">
12      <p>Grabinschrift für Theodoros; gef. in
          Ankara, einmal im Augustus–Tempel (
          bis 1927).</p>
13    </div>
14    <div type="bibliography">
15      <listBibl>
16        <bibl>
17          <ptr target="EDAK_bibliography.xml
          #SEG"/>
18          <citedRange>6, 29 (zur Datierung)</
          citedRange>
19        </bibl>
20        <bibl>
21          <ptr target="EDAK_bibliography.xml
  
```



7	<publicationStmt>		<abbr>Θ<ex> </ex> </abbr></
8	<authority>Epigraphische Datenbank	6	expan>
	zum antiken Kleinasien –	7	<lb n="3"/>Θ
	Universität Hamburg</authority	8	<lb n="4" break="no"/> ,<expan><
9	>	9	abbr> </abbr><ex> </ex></
	<idno type="localID">EDAK00001210	10	expan> ',
	</idno>	11	<lb n="5"/><expan><abbr> </abbr><
10	<idno type="filename">EDAK00001210	12	ex> </ex></expan>M '.
	</idno>	13	</ab>
11	<availability>[...]</availability>	14	</div>
12	</publicationStmt>	15	<div type="commentary">
13	<sourceDesc>	16	<p>Grabinschrift für Theodoros; gef. in
14	<msDesc>	17	Ankara, einmal im Augustus–Tempel (
15	<msIdentifier/>	18	bis 1927).</p>
16	<msContents>[...]</msContents>	19	</div>
17	<physDesc>[...]</physDesc>	20	<div type="bibliography">
18	<history>	21	<listBibl>
19	<origin>	22	<bibl>
20	<origPlace>	23	<ptr target="EDAK_bibliography.xml
21	<settlement ref="	24	#SEG"/>
	origPlace.xml#	25	<citedRange>6, 29 (zur Datierung)</
	fundort_83"/>	26	citedRange>
22	</origPlace>	27	</bibl>
23	<origDate notBefore="0301"	28	<bibl>
	notAfter="0400">4 Jh.	29	<ptr target="EDAK_bibliography.xml
	n.Chr. (Hondius)</	30	#Anderson1899"/>
	origDate>	31	<citedRange>97 Nr. 80</citedRange>
24	</origin>	32	</bibl>
25	<provenance type="found">	33	<bibl>
26	<placeName ref="findspot.	34	<ptr target="EDAK_bibliography.xml
	xml#fundort_83"/>	35	#Jerphanion1928"/>
27	</provenance>		<citedRange>-287288 Nr. 61 (zu Z. 1,
28	</history>		-45).</citedRange>
29	</msDesc>	27	</bibl>
30	</sourceDesc>	28	</listBibl>
31	</fileDesc>	29	</div>
32	<encodingDesc>[...]</encodingDesc>	30	<div type="translation"> <p/>
33	<profileDesc>[...]</profileDesc>	31	</div>
34	<revisionDesc>[...]</revisionDesc>	32	<div type="apparatus"> <p/>
35	</teiHeader>	33	</div>
		34	</body>



# EpiDoc Stylesheets

## machine-readable format

```
...  
<ab>  
  <lb n="1"/>  
  Πλανκίαν Μάγναν  
  <lb n="2"/>  
  Ἀκυλλίαν θυγατέρα Ἰου  
  <lb n="3" break="no" />  
  λίου Σεουήρου καί  
  <abbr>Κλαυ</abbr>  
  <ex>δίας</ex>  
  <lb n="4"/>  
  Ἀκυλλίας, τὴν ἐν βασιλέ  
  <lb n="5" break="no" />  
  ων, ἡρωίδα, ὁ δῆμος ὁ Ἄν  
  <lb n="6" break="no" />  
  κυρανῶν τὴν θυγατέρα τῆς  
  <lb n="7"/>  
  μητροπόλεως.  
</ab>  
...
```

EpiDoc  
stylesheets

## human-readable format

- 1 Πλανκίαν Μάγναν
- 2 Ἀκυλλίαν θυγατέρα Ἰου-
- 3 λίου Σεουήρου καί Κλαυ(δίας)
- 4 Ἀκυλλίας, τὴν ἐν βασιλέ-
- 5 ων, ἡρωίδα, ὁ δῆμος ὁ Ἄν-
- 6 κυρανῶν τὴν θυγατέρα τῆς
- 7 μητροπόλεως.

# JSON as Data Representation (PATHs Project Example)

- JSON is a lightweight format for representing structured data
- It is widely used in web applications and APIs
- Data is organized as key-value pairs
- It supports nested structures and lists
- It is easy to read for both humans and machines
- JSON is commonly used for data exchange between systems
- It provides a flexible alternative to XML-based representations

```
1 {
2   "head": {
3     "shortsql": "@manuscripts~?bindings|=|1~-200:0",
4     "total_rows":156,
5     "total_pages":6,
6     "table":"paths__manuscripts",
7     "stripped_table":"manuscripts",
8     "table_label":"Manuscripts",
9     "page":1,
10    "no_records_shown":156,
11    "fields": {[...]},
12  },
13  "records": [
14    {
15      "0": {
16        "id": "18",
17        "creator": "1",
18        "cmclid": "CMCL.AS",
19        "tm": "108136",
20        "ldab": "108136",
21        "lcbm": null,
22        "dbmnt": null,
23        "alias": "Codex Bruce",
24        "issinglefrag": "0",
25        "isbookbinding": null,
```

# Bookbinding Data in JSON

- The example shows research data from the PATHS project
- Each record represents one bookbinding entity
- Attributes include stratigraphy, modern history, and book form
- Data is organized into fields and records
- Each record is identified by an internal identifier
- Textual descriptions are stored as string values
- The structure reflects a simplified data model

```
1 {
2   "head": {
3     "shortsql": "@manuscripts~?bindings|=|1~-200:0",
4     "total_rows":156,
5     "total_pages":6,
6     "table":"paths__manuscripts",
7     "stripped_table":"manuscripts",
8     "table_label":"Manuscripts",
9     "page":1,
10    "no_records_shown":156,
11    "fields": {[...]},
12  },
13  "records": [
14    {
15      "0": {
16        "id": "18",
17        "creator": "1",
18        "cmclid": "CMCL.AS",
19        "tm": "108136",
20        "ldab": "108136",
21        "lcbm": null,
22        "dbmnt": null,
23        "alias": "Codex Bruce",
24        "issinglefrag": "0",
25        "isbookbinding": null,
```

# Audition Certificate Data in JSON

```
1 {
2   "text": "[...] لضافلا ملا علا هي قفلا هبجا صبات كلا اذ هي ميجيلعم سبلا هم حرفلؤ ملا طخبعا مسة قبطة روص",
3   "properties": [
4     {
5       "index": 1,
6       "guid": "b8fdb5ba-fcd8-4b8d-b652-29a2bb5066f9",
7       "type": "place",
8       "text": "قشمد",
9       "startIndex": 419,
10      "endIndex": 422,
11      "attributes": {
12        "location": [ "bilad-al-sham", "damascus" ]
13      },
14      "isZeroPoint": false
15    },
16    {
```

```
5 index : 1,  
6 "guid": "b8fdb5ba-fcd8-4b8d-b652-29a2bb5066f9",  
7 "type": "place",  
8 "text": "قشمد",  
9 "startIndex": 419,  
10 "endIndex": 422,  
11 "attributes": {  
12   "location": [ "bilad-al-sham", "damascus"]  
13 },  
14 "isZeroPoint": false  
15 },  
16 {  
17   "index": 2,  
18   "guid": "1a933be2-c462-4a36-a386-217a4dccf6ce",  
19   "type": "person",  
20   "text": "باهش"  
21   "يعفاشلا يقشمد لا يراصناً لا ناو عجنبسابعنبد محمهلا دبعيياً نيدلا لا مكليجاً لا خيشلا نبد محاً سابعلا وبأ نيدلا",  
22   "startIndex": 107,  
23   "endIndex": 218,  
24   "attributes": {  
25     "role": ["reader", "owner"],  
    ...
```

# JSON vs TEI

## JSON

- Typically flatter and more application-oriented
- Easier to process in programming environments
- Commonly used for APIs and data exchange

## TEI

- Richer and more expressive for textual data
- Captures detailed semantics and annotations
- Is used for scholarly encoding and preservation



Both formats represent data models in different ways

# JSON as a Data Model

- JSON structures represent entities as objects
- Keys correspond to attributes of these entities
- Nested objects represent relationships between entities
- Arrays represent collections of entities
- Identifiers are required to link related data
- The structure reflects modeling decisions made during transformation
- JSON therefore encodes a data model, not just a format

# TEI and JSON as Alternative Representations

- The same data can be represented in TEI or JSON
- TEI uses hierarchical XML structures
- JSON uses object-based structures
- Both encode entities, attributes, and relationships
- TEI focuses on semantic richness and textual detail
- JSON focuses on usability in applications and APIs
- Transformation between both formats is possible

# Multiple Representations

- The same research data can be represented in different formats
- DOCX represents data implicitly
- TEI and EpiDoc represent data explicitly and semantically
- JSON represents data in a structured and flexible way
- Each format encodes a different data model
- The choice of format affects analysis possibilities
- Understanding representations is key to data modeling

# Schema & Validation (Relax NG)

## Why Validation?

- Data must follow consistent rules
- Errors must be detected early
- Structure must be enforced
- Validation ensures correctness
- It supports interoperability
- It improves data quality
- It enables reliable processing

# Relax NG

- Elements may be added in an TEI/EpiDoc document in a standardized way with the RELAX NG standard
- RELAX NG is a schema language for XML
- It is also an International Standard ISO/IEC 19757-2 and is part of ISO/IEC 19757 DSDL (Document Schema Definition Languages)
- RELAX NG has an XML syntax
- <https://relaxng.org/>

```
1 <?xml version="1.0" encoding="utf-8"?>
2 <grammar xmlns:xlink="http://www.w3.org/1999/xlink"
3     xmlns:tei="http://www.tei-c.org/ns/1.0"
4     xmlns:teix="http://www.tei-c.org/ns/Examples"
5     xmlns="http://relaxng.org/ns/structure/1.0"
6     datatypeLibrary="http://www.w3.org/2001/XMLSchema-datatypes"
7     ns="http://www.tei-c.org/ns/1.0">
8 <define name="tei_model.pPart.msdesc">
9   <choice>
10    <ref name="tei_catchwords"/>
11    <ref name="tei_dimensions"/>
12    <ref name="tei_heraldry"/>
13    <ref name="tei_locus"/>
14    <ref name="tei_locusGrp"/>
15    <ref name="tei_material"/>
16    <ref name="tei_objectType"/>
17    <ref name="tei_origDate"/>
18    <ref name="tei_origPlace"/>
19    <ref name="tei_secFol"/>
20    <ref name="tei_signatures"/>
21    <ref name="tei_stamp"/>
22    <ref name="tei_watermark"/>
23  </choice>
24 </define>
25 <define name="tei_title">
26   <element name="title">
27     <a:documentation xmlns:a="http://relaxng.org/ns/compatibility/annotations
28       /1.0">(title) contains a title for any kind of work. [3.12.2.2. Titles,
29       Authors, and Editors 2.2.1. The Title Statement 2.2.5. The Series
30       Statement]</a:documentation>
31     <ref name="tei_macro.paraContent"/>
32     <ref name="tei_att.global.attributes"/>
33     <ref name="tei_att.global.attributes"/>
```

# Relax NG

- Elements may be added in an TEI/EpiDoc document in a standardized way with the RELAX NG standard
- RELAX NG is a schema language for XML
- It is also an International Standard ISO/IEC 19757-2 and is part of ISO/IEC 19757 DSDL (Document Schema Definition Languages)
- RELAX NG has an XML syntax
- <https://relaxng.org/>

```
16 <ref name="tei_objectType"/>
17 <ref name="tei_origDate"/>
18 <ref name="tei_origPlace"/>
19 <ref name="tei_secFol"/>
20 <ref name="tei_signatures"/>
21 <ref name="tei_stamp"/>
22 <ref name="tei_watermark"/>
23 </choice>
24 </define>
25 <define name="tei_title">
26 <element name="title">
27 <a:documentation xmlns:a="http://relaxng.org/ns/compatibility/annotations
   /1.0">(title) contains a title for any kind of work. [3.12.2.2. Titles,
   Authors, and Editors 2.2.1. The Title Statement 2.2.5. The Series
   Statement]</a:documentation>
28 <ref name="tei_macro.paraContent"/>
29 <ref name="tei_att.global.attributes"/>
30 <ref name="tei_att.typed.attribute.subtype"/>
31 <ref name="tei_att.canonical.attributes"/>
32 <ref name="tei_att.dateable.attributes"/>
33 <optional>
34 <attribute name="type">
35 <a:documentation xmlns:a="http://relaxng.org/ns/compatibility/
   annotations/1.0">classifies the title according to some convenient
   typology.
36 Sample values include: 1] main; 2] sub (subordinate); 3] alt (alternate); 4] short; 5] desc
   (descriptive)</a:documentation>
37 <data type="token">
38 <param name="pattern">[^\p{C}\p{Z}]+</param>
39 </data>
40 </attribute>
41 </optional>
42 [...]
43 </element>
44 </define>
```

# Function of Relax NG

- It defines structural constraints
- It validates XML documents
- It detects invalid elements
- It enforces modeling rules
- It ensures consistency across datasets
- It supports automated workflows
- It strengthens data reliability

# Same Standard, Different Data Models (EpiDoc Case)

- Different projects can use the same standard such as EpiDoc
- Examples include CGRN and EDAK, both based on TEI/EpiDoc
- Despite using the same framework, their schemas can differ
- Each project defines its own elements, attributes, and constraints
- Modeling decisions depend on research questions and domain needs
- This leads to variations in structure and interpretation of data
- Therefore, interoperability is not guaranteed by the standard alone

# Same Standard, Different Data Models (EpiDoc Case)

```
▼ <origin>
  ▼ <p>
    <origDate notBefore="-0550" notAfter="-0500">ca. 550-500 BC</origDate>
  </p>
  ► <p>
    ...
  </p>
</origin>
```

```
▼ <origin>
  ► <origPlace>
    ...
  </origPlace>
  <origDate notBefore="0101" notAfter="0200">A) 2nd cent. CE (Taeuber, letters)</origDate>
  <origDate notBefore="0301" notAfter="0400">B) 4th cent. CE (Taeuber, letters)</origDate>
</origin>
```

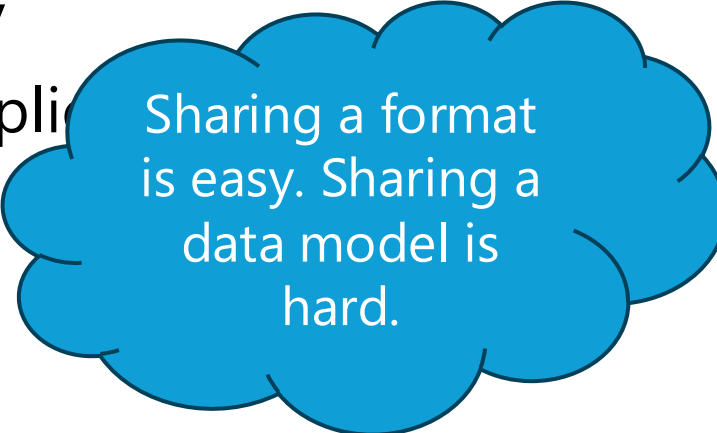
# Same Standard, Different Data Models (EpiDoc Case)

```
▼ <provenance>
  ▼ <p>
    ▼ <placeName type="ancientFindspot" key="Ephesos" n="Asia_Minor_and_Anatolia">
      <ref target="https://pleiades.stoa.org/places/599612" type="external">Ephesos</ref>
    </placeName>

  ▼ <origPlace>
    ▼ <placeName type="ancientFindspot">
      ▼ <settlement>
        <ref target="https://pleiades.stoa.org/places/599612/">Ephesos</ref>
      ▼ <location>
        <geo>37.9407625 27.340307</geo>
      </location>
    </settlement>
    ▼ <region>
      <ref target="https://pleiades.stoa.org/places/550597/ionia/">Ionia</ref>
    </location>
  </region>
</placeName>
</origPlace>
```

# Implications for Data Modeling

- A standard provides a common vocabulary but not a fixed model
- Customization allows flexibility but introduces heterogeneity
- Data integration requires understanding project-specific schemas
- Validation (e.g., Relax NG) enforces local consistency only
- Cross-project comparison becomes more complex
- Shared conventions are needed for true interoperability
- Modeling decisions must be documented and made explicit

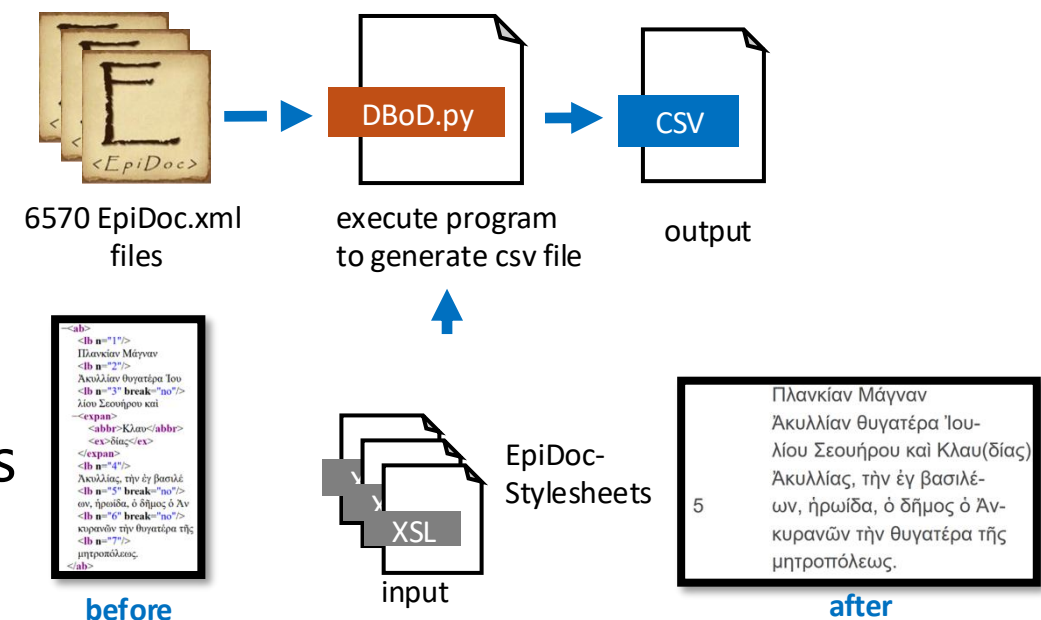


Sharing a format  
is easy. Sharing a  
data model is  
hard.

# From EpiDoc/TEI to Structured Data

## Transformation Concept

- EpiDoc and TEI data can be transformed into other formats
- XML structures can be flattened into tables
- Data extraction focuses on entities
- Relationships must be preserved
- Transformation requires modeling decisions
- Tools support automated conversion
- This enables further processing



# From EpiDoc/TEI to Structured Data

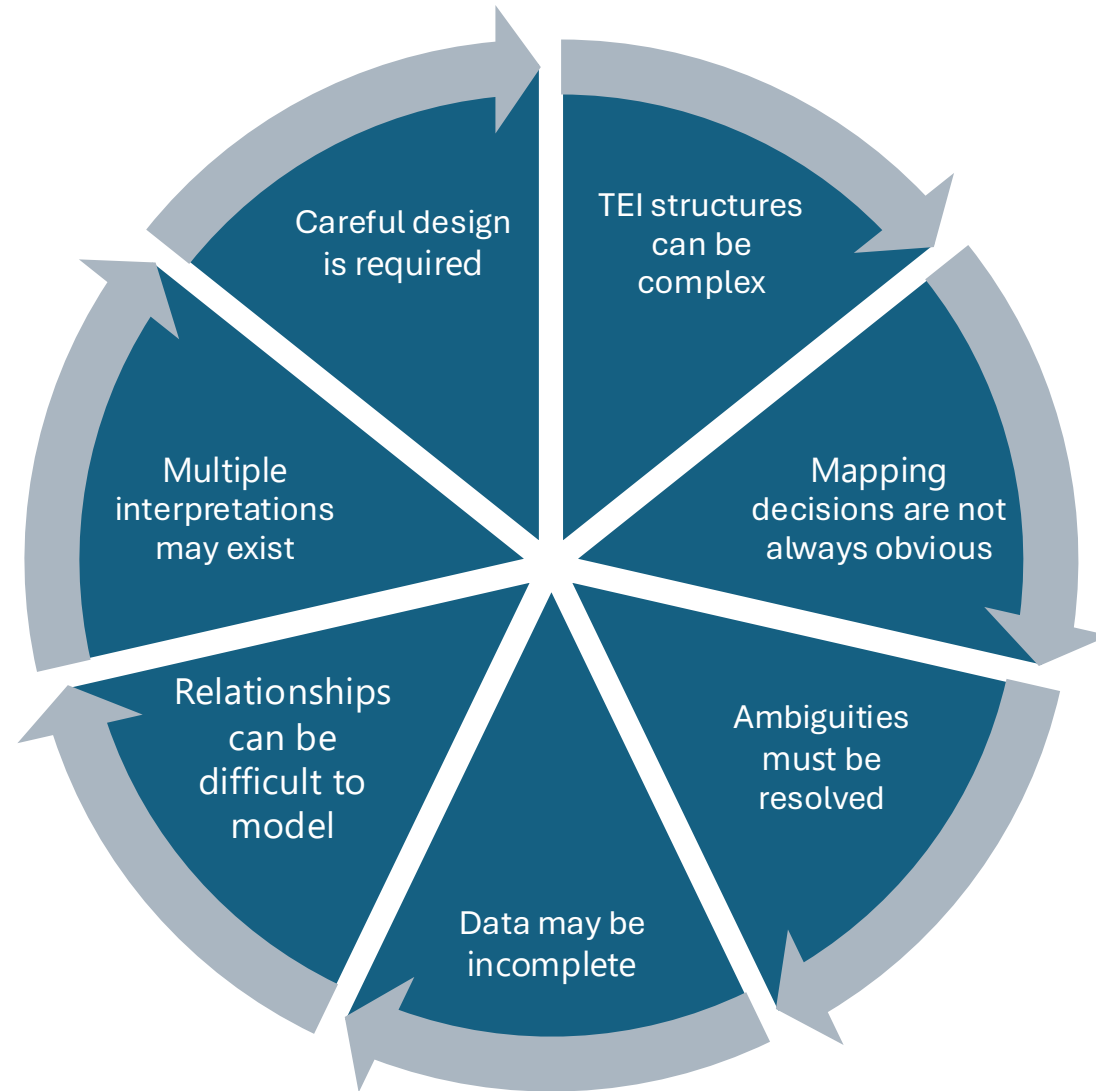
## Extracting Data

- Entities such as person names, object type and date can be extracted
- Places can be identified and structured
- Relationships between entities are mapped
- Attributes become table columns
- Identifiers are required for linking
- Data must be cleaned and standardized
- This produces structured datasets

# TEI to CSV

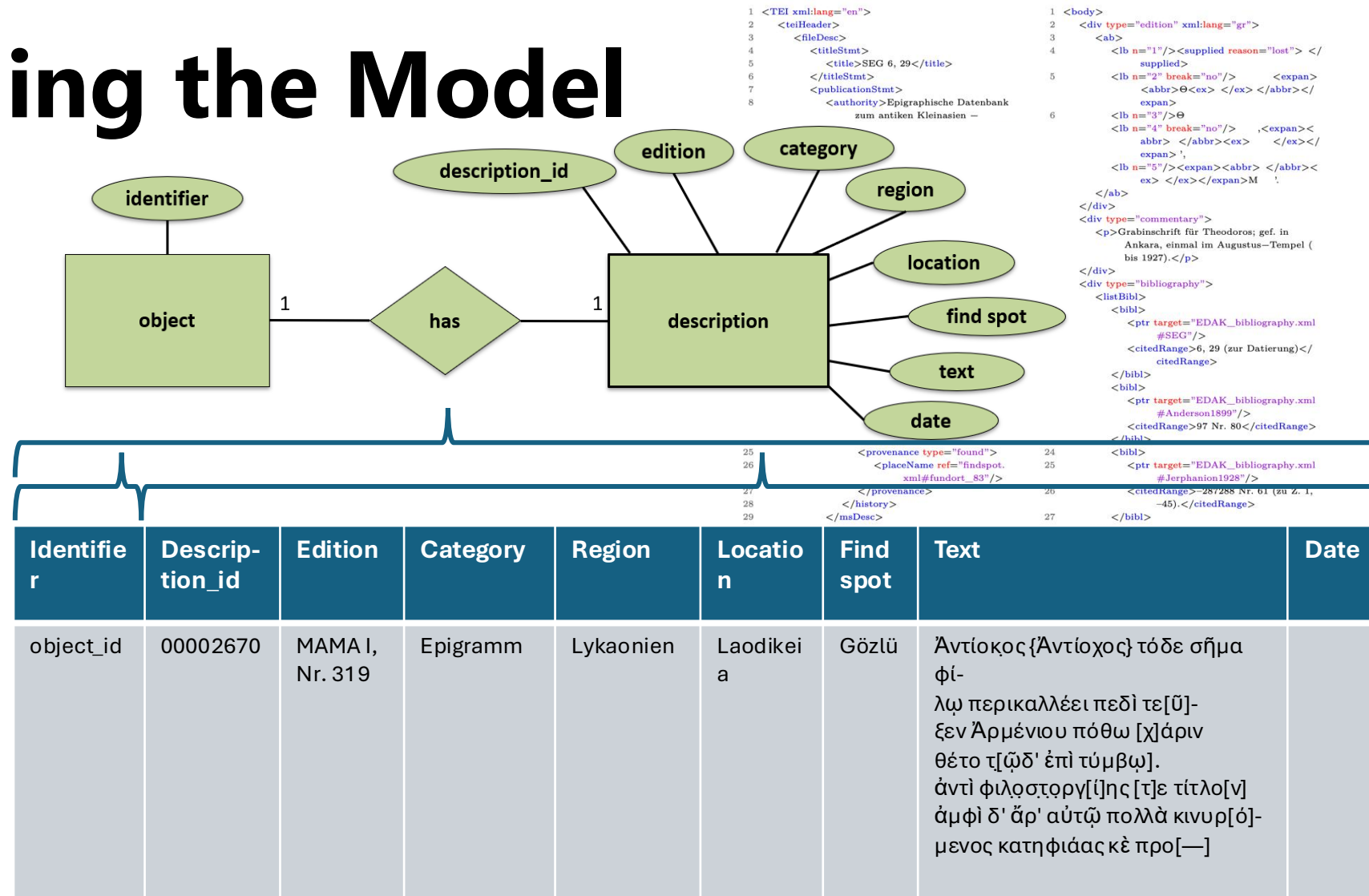
- CSV represents tabular data
- Hierarchical data must be simplified
- Each row represents one entity instance
- Columns represent attributes
- Relationships require additional tables
- Information loss must be managed
- CSV enables database integration

# Challenges



# Understanding the Model

- Identify entities within the dataset
- Determine attributes of each entity
- Analyze relationships between entities
- Define keys and identifiers
- Evaluate consistency of data
- Consider normalization aspects
- Reflect on modeling decisions



# Why Normalization Here?

- Data extracted from TEI is often redundant or inconsistent
- CSV representations may duplicate information across rows
- Databases require well-structured data models
- Normalization helps to improve data quality
- It reduces redundancy and avoids inconsistencies
- It prepares data for efficient querying
- It connects transformation with database design

# First Normal Form (1NF)

- Each field contains atomic values
- No repeating groups are allowed
- Data is organized in rows and columns
- Each record is uniquely identifiable
- Structure becomes consistent
- Queries become easier

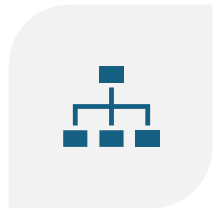
# Second Normal Form (2NF)

- Data must already satisfy 1NF
- No partial dependencies are allowed
- Attributes depend on the full key
- Redundancy is reduced
- Data integrity is improved
- Relationships are clarified

# Third Normal Form (3NF)

- Data must satisfy 2NF
- No transitive dependencies are allowed
- Attributes depend only on the key
- Redundancy is minimized
- Consistency is improved
- Updates become simpler

# TEI vs Relational Models



TEI USES  
HIERARCHICAL  
STRUCTURES



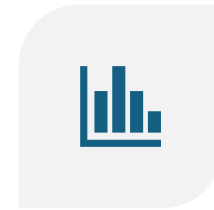
DATABASES USE  
TABULAR  
STRUCTURES



TRANSFORMATION  
IS REQUIRED  
BETWEEN BOTH



TEI CAPTURES RICH  
SEMANTICS



DATABASES  
SUPPORT  
EFFICIENT QUERIES



BOTH  
COMPLEMENT  
EACH OTHER

# Example Projects

- Projects (EDAK, NETamil, ACP, PATH) differ in domain and scope.
- All require structured data models.
- Interoperability is a common goal.
- These examples illustrate real-world applications.

# Conclusion

- Research data must be structured
- Representation formats define models
- TEI and EpiDoc provide rich semantic encoding
- JSON provides structured, application-oriented representation
- Relax NG ensures correctness
- Transformation enables analysis
- Data models improve quality
- Structured data enables FAIR principles

# Discussion

- Is TEI a data format, a data model, or both?
- What is lost and what is gained when transforming TEI into CSV or JSON?
- Which representation is best suited for analysis?