

# Text Encoding and Semantic Representation

---

XML / TEI (1)

marilena.daquino2@unibo.it | [https://github.com/marilenadaquino/tesr\\_dhdk](https://github.com/marilenadaquino/tesr_dhdk)



# Work at home

---

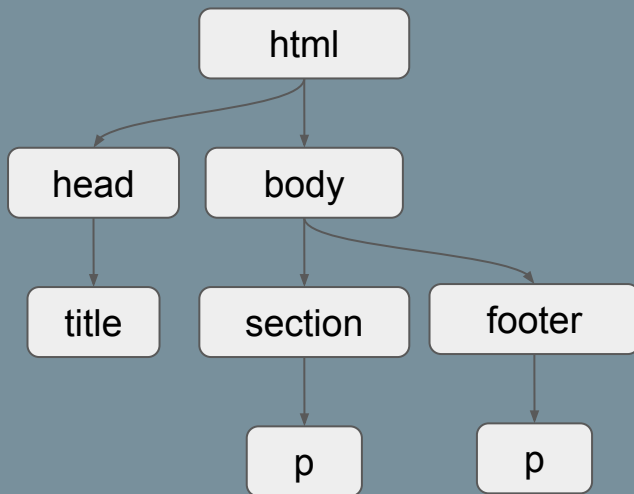
## Practice markup

Given the following tree, create a XML document.  
This time, do not add the prologue!

Fill elements with some text when appropriate.

Tip: open the xml file of both the exercise and the homework in a browser (e.g. Chrome)

NOW... Change the file extension to .html  
and open the new file in a browser



# TEI

---



# In pills

Stands for Text Encoding Initiative.

It is both a schema to annotate literary texts and a community of scholars (mainly philologists). The website includes a set of **guidelines** on how to use the XSD by thematic areas and examples.

The first version of the guidelines was released in 1990 (currently v4.6). The reference guidelines are called P5 Guidelines.

# TEI

---

## Modules

The schema is divided in (thematic) modules, i.e. groupings of XML elements and attributes (~500 elements in total).

Many modules can be used in the same document.

The **core module** includes those elements that are likely to be used to describe any XML document.

Module name	Formal public identifier
analysis	Analysis and Interpretation
certainty	Certainty and Uncertainty
core	Common Core
corpus	Metadata for Language Corpora
dictionaries	Print Dictionaries
drama	Performance Texts
figures	Tables, Formulae, Figures
gaiji	Character and Glyph Documentation
header	Common Metadata
iso-fs	Feature Structures
linking	Linking, Segmentation, and Alignment
msdescription	Manuscript Description
namesdates	Names, Dates, People, and Places
nets	Graphs, Networks, and Trees
spoken	Transcribed Speech
tagdocs	Documentation Elements
tei	TEI Infrastructure
textcrit	Text Criticism
textstructure	Default Text Structure
transcr	Transcription of Primary Sources
verse	Verse

# TEI

---

## Global attributes

Likewise, there are a number of global attributes, meaning they can be used with any element (while non-global attributes can be used only with specific sets of attributes)

<b>@xml:id</b>	(identifier) provides a unique identifier for the element bearing the attribute.
<b>@n</b>	(number) gives a number (or other label) for an element, which is not necessarily unique within the document.
<b>@xml:lang</b>	(language) indicates the language of the element content using a 'tag' generated according to <a href="#">BCP 47</a> .
<b>@rend [att.global.rendition]</b>	(rendition) indicates how the element in question was rendered or presented in the source text.
<b>@style [att.global.rendition]</b>	contains an expression in some formal style definition language which defines the rendering or presentation used for this element in the source text
<b>@rendition [att.global.rendition]</b>	points to a description of the rendering or presentation used for this element in the source text.
<b>@xml:base</b>	provides a base URI reference with which applications can resolve relative URI references into absolute URI references.
<b>@xml:space</b>	signals an intention about how white space should be managed by applications.
<b>@source [att.global.source]</b>	specifies the source from which some aspect of this element is drawn.
<b>@cert [att.global.responsibility]</b>	(certainty) signifies the degree of certainty associated with the intervention or interpretation.
<b>@resp [att.global.responsibility]</b>	(responsible party) indicates the agency responsible for the intervention or interpretation, for example an editor or transcriber.

# TEI

---

## Mandatory elements

The element **teiHeader** includes all the metadata about the digital edition of the text.



```
1 <TEI xmlns="http://www.tei-c.org/ns/1.0">
2   <teiHeader xml:lang="en">
3     <!-- ... -->
4   </teiHeader>
5   <text xml:lang="en">
6     <!-- ... -->
7   </text>
8 </TEI>
```

The root element **TEI** includes the namespace declaration

The element **text** includes the actual text.

# TEI

---

## Mandatory elements

When annotating multiple texts, the root element is **teiCorpus**, with a bespoke **teiHeader**, and there are as many **TEI** elements as the number of documents to be annotated.

```
1 <teiCorpus xmlns="http://www.tei-c.org/ns/1.0">
2   <teiHeader>
3     <!-- corpus-level metadata here -->
4   </teiHeader>
5   <TEI>
6     <teiHeader>
7       <!-- metadata specific to this text here -->
8     </teiHeader>
9     <text><!-- ... --></text>
10  </TEI>
11  <TEI>
12    <teiHeader>
13      <!-- metadata specific to this text here -->
14    </teiHeader>
15    <text><!-- ... --></text>
16  </TEI>
17 </teiCorpus>
```



# Disclaimer

---

Many solutions to the same problem create a new problem

Other than mandatory elements, in most cases there is not just one way to annotate the text, but several variants are possible (e.g. you may wrap sentences in an element or you can annotate substrings in multiple elements). Unfortunately, such a variety is not that good for exchange purposes :(

# Disclaimer

---

500 elements are a lot...

We are not going to see the full elementset, just some of the most important ones. You are free to study more elements according to your (encoding) needs.

# TEI

---

## teiHeader

**teiHeader** includes the metadata of the digital edition, such as people's responsibilities, the bibliographic citation, and the strategies used to encode the text (e.g. naming conventions, critical choices in the transcription).

Only **fileDesc** is mandatory.

```
1 <teiHeader>
2   <fileDesc></fileDesc>
3   <encodingDesc></encodingDesc>
4   <profileDesc></profileDesc>
5 </teiHeader>
```

# TEI

## fileDesc

**fileDesc** (mandatory) describes the current electronic file, with the exception of **sourceDesc**, which describes the original source.

**titleStmt**, **publicationStmt** and **sourceDesc** are **mandatory**.

[Guidelines](#)

```
1 <teiHeader>
2   <fileDesc>
3     <titleStmt>
4       <title><!-- title of the resource --></title>
5     </titleStmt>
6     <editionStmt>
7       <p><!-- the edition of the resource --></p>
8     </editionStmt>
9     <extent><!-- the size of the resource --></extent>
10    <publicationStmt>
11      <p><!-- the distribution of the resource --></p>
12    </publicationStmt>
13    <seriesStmt>
14      <p><!-- any series to which the resource belongs --></p>
15    </seriesStmt>
16    <notesStmt>
17      <note><!-- other aspects of the resource --></note>
18    </notesStmt>
19    <sourceDesc>
20      <p><!-- the source from which the resource was derived --></p>
21    </sourceDesc>
22  </fileDesc>
23 </teiHeader>
```

# TEI

---

## titleStmt

**title** includes the title of the digital edition (derived from the original title)

It can include also the **author** of the original work and all the people that contributed to the digital edition - each recorded in **respStmt**, including the role (**resp**) and the **name** of the person.

Guidelines

```
1 <titleStmt>
2   <title>Resistance: digital edition</title>
3   <author>Muse</author>
4   <respStmt>
5     <resp>compiled by</resp>
6     <name>Marilena Daquino</name>
7   </respStmt>
8 </titleStmt>
```

# TEI

---

## editionStmt

**editionStmt** includes information about the current edition of the digital text (not the edition of the original text).

**edition** includes a string and/or children elements. **respStmt** can be used to record people and roles.

[Guidelines](#)

```
1 <editionStmt>
2   <edition n="1">Private edition, <date>November 2023</date></edition>
3   <respStmt>
4     <resp>Annotations by</resp>
5     <name>Marilena Daquino</name>
6   </respStmt>
7 </editionStmt>
```

# TEI

---

## extent

Describes the extent of the digital file,  
e.g. in terms of MB, GB or other  
units of measures.

Notice the usage of the element  
measure and its attributes.

Guidelines

```
1 <!-- a string -->
2 <extent>between 1 and 2 Mb</extent>
3
4 <!-- or one or more measurements -->
5 <extent>
6   <measure unit="MiB" quantity="2">About 2 megabytes</measure>
7   <measure unit="pages" quantity="1">1 page of source material</measure>
8 </extent>
```

# TEI

---

## publicationStmt

**publicationStmt** includes information about the publisher of the digital edition and the rights/licenses under which the resource is available.

NB. The digital edition of a text that is under copyright cannot be open access!

Guidelines

```
1 <publicationStmt>
2   <publisher>University of Bologna</publisher>
3   <pubPlace>Bologna</pubPlace>
4   <date>2023</date>
5   <idno type="DOI">10.xxxx</idno>
6   <availability>
7     <p>Open access</p>
8   </availability>
9   <licence target="https://creativecommons.org/licenses/by/4.0/">CC-BY</licence>
10 </publicationStmt>
```



# TEI

---

## sourceDesc

**sourceDesc** includes the bibliographic reference of the original source.

Details can include authors, title, dates, publisher, etc.

[Guidelines](#)

```
1 <sourceDesc>
2   <bibl>
3     <title level="a">Resistance</title>. In
4     <author>Muse</author>,
5     <title level="b">Resistance</title>.
6     <date>2009</date>.
7   </bibl>
8 </sourceDesc>
```

# TEI

---

## encodingDesc

It specifies the methods and **editorial principles** which guided the transcription or encoding of the text, such as corrections (if any error was corrected), normalisation strategies (e.g. en-us spelling), or interpretations (e.g. if you add elements of linguistic analysis).

```
1 <encodingDesc>
2   <projectDesc>
3     <p><!-- The purpose of the digital resource --></p>
4   </projectDesc>
5   <editorialDecl>
6     <correction>
7       <p><!-- If the text includes corrections--></p>
8     </correction>
9     <normalization>
10      <p><!-- If any changes happened to uniform editorial choices --></p>
11    </normalization>
12    <interpretation><!-- If any analytical feature was added to the text--></interpretation>
13  </editorialDecl>
14 </encodingDesc>
```

# Exercise

---

## Create a XML/TEI file

Create a XML/TEI file to describe the following song:

*Bohemian Rhapsody*, Queen, 1975. <https://genius.com/Queen-bohemian-rhapsody-lyrics>

Include the element **teiHeader** and all the mandatory children we have seen.

# TEI

---

## The transcription of the text

The transcribed text is included in the element `<text>`. It may include a front matter `<front>`, a text body `<body>`, and a back matter `<back>`.

[Guidelines](#)

```
1 <TEI xmlns="http://www.tei-c.org/ns/1.0">
2   <teiHeader>
3     <!-- ... -->
4   </teiHeader>
5   <text>
6     <front>
7       <!-- front matter of copy text, if any, goes here -->
8     </front>
9     <body>
10      <!-- body of copy text goes here -->
11    </body>
12    <back>
13      <!-- back matter of copy text, if any, goes here -->
14    </back>
15  </text>
16 </TEI>
```

# TEI

---

## Elements in all TEI documents

Some elements can be used in any TEI document regardless of the type/genre of literary text. These include: paragraphs <p>, quotations <q>, graphically highlighted terms: <hi>, foreign terms: <foreign>, emphasis <emph>, etc...

Guidelines

```
1 <p>This is a paragraph mentioning
2   a <foreign xml:lang="fr">croissant</foreign>,
3   a <hi rend="italic">highlight</hi>,
4   and some <emph>important concept</emph>.
5 </p>
6
7 <q>It is followed by a quotation from one of
   your favourite authors</q>
```

# TEI

---

## Divisions of the text

A prose text can be divided in several structures, e.g. parts, chapters, sections. In this case divisions are represented with the element `<div>` and the attributes `@type` and `@num` further specify their peculiarities and differentiate them. Values are free from naming conventions.

```
1 <body>
2   <div type="part" n="1">
3     <div type="chapter" n="1">
4       <!-- text of part 1, chapter 1 -->
5     </div>
6     <div type="chapter" n="2">
7       <!-- text of part 1, chapter 2 -->
8     </div>
9   </div>
10  <div type="part" n="2">
11    ...
12  </div>
13 </body>
```

# TEI

---

## Headings

Inside a division, titles of all levels can be represented with `<head>` and further specified via the attributes `@rend`, `@type` and more

```
1 <div n="19" type="chap">
2   <head rend="bold" type="main">Chapter 19</head>
3   <p>To say that Deronda was romantic would be to
4     misrepresent him: but under his calm and
5     somewhat self-repressed exterior ...</p>
6 </div>
```

# TEI

---

## Lists

Bullet lists are represented via the element `<list>` and its children `<item>`. Lists can be ordered or unordered. Items can include other lists.

```
1 <list rend="numbered">
2   <item>a butcher</item>
3   <item>a baker</item>
4   <item>a candlestick maker, with
5   <list rend="bulleted">
6     <item>rings on his
7     fingers</item>
8   </list>
9   </item>
10 </list>
```



# Exercise

---

## Modify a XML/TEI file

- Modify your XML, include the element **text** and add the transcription of the lyrics.

# TEI

---

## Deletions, corrections

Omissions, changes (deletions and additions), as well as unclear texts are recorded in elements like `<gap>` (words omitted), `<del>` and `<add>`

Guidelines

```
1 <gap reason="illegible" unit="word" quantity="2"  
  resp="#editor04"/>  
2 <!-- notice the attribute resp -->  
3  
4 <l>  
5   <del rend="overstrike">Inviolable</del>  
6   <add place="below">Inexplicable</add>  
7   splendour of Corinthian white and gold  
8 </l>
```

# TEI

---

## Letters

Letters include peculiar logical elements such as

<opener>, <closer> and <postscript>.

Line breaks in a paragraph can be recorded via <lb>, a milestone element written in the short version.

In the example, also notice the element <unclear> to mark words difficult to read.

Guidelines

```

1 <div type="letter">
2   <opener>
3     <dateline>
4       <placeName>Newport</placeName>
5       <date when="1761-05-27">May ye 27th 1761</date>
6     </dateline>
7     <salute>Gentlemen</salute>
8   </opener>
9   <p>Capt Stoddard's Business
10  <lb/>calling him to Providence, have
11  <lb/>got him to look at Hopkins brigantine
12  <lb/>&amp; if can agree to Purchase her, shall
13  <lb/>be much oblig'd for your further
14  <lb/>assistance herein, &amp; will acquiesce with
15  <lb/>whatever you &amp; he shall Contract
16  <lb/>for – I Thank you for your
17  <lb/>
18    <unclear>Line</unclear> respecting the brigantine &amp;
19  <lb/>leave to Recommend the Bearer
20  <lb/>to you for your advice &amp; Friendship
21  <lb/>in this matter</p>
22  <closer>
23  <salute>I am your most humble servant</salute>
24  <signed>Joseph Wanton Jr</signed>
25 </closer>
26 <postscript>
27   <label>P.S.</label>
28   <p>I have Mollases, Sugar,
29   <lb/>Coffee &amp; Rum, which
30   <lb/>will Exchange with you
31   <lb/>for Candles or Oyl</p>
32 </postscript>
33 </div>

```

# TEI

---

## Verses

Poems can be characterised by ungrouped lines (element `<l>`). When grouped, lines are included in the element `<lg>`. The attribute `@type` can specify the structure or unit, e.g. a sonnet or a stanza.

[Guidelines](#)

```
1 <text>
2 <body>
3   <lg type="quatrain">
4     <l>My Mistres eyes are nothing like the Sunne,</l>
5     <l>Curall is farre more red, then her lips red</l>
6     <l>If snow be white, why then her brests are dun:</l>
7     <l>If haire be wiers, black wiers grown on her head:</l>
8   </lg>
9   ...
10  <lg type="couplet">
11    <l>And yet by heaven I think my love as rare,</l>
12    <l>As any she beli'd with false compare.</l>
13  </lg>
14 </body>
15 </text>
```

# TEI

---

## Performances

Performance texts may include descriptions of the stage and the sequence of speakers' dialogues.

Here like in other situations, the hierarchy of text divisions can be further specified by using the elements `<div1>` and `<div2>` (not mandatory).

Guidelines

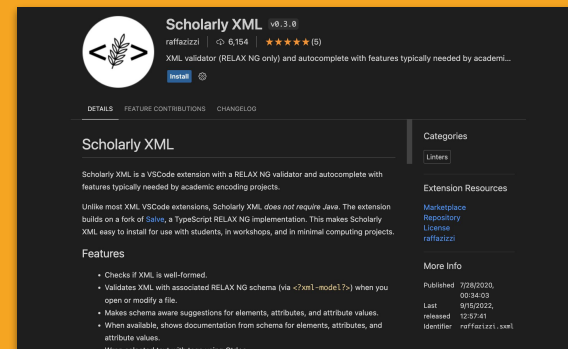
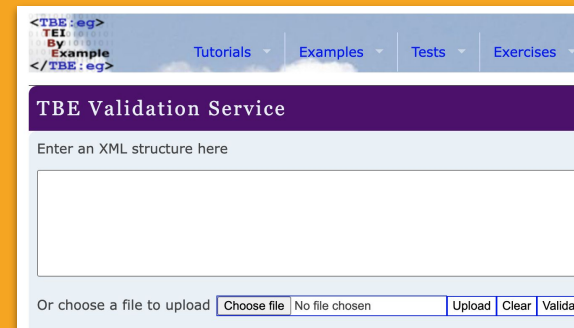
```
1 <body>
2 <div1 type="act" n="1">
3   <head>Act One</head>
4   <div2 type="scene" n="1">
5     <stage>Pa Ubu, Ma Ubu</stage>
6     <sp>
7       <speaker>Pa Ubu</speaker>
8       <p>Pschitt!</p>
9     </sp>
10  </div2>
11 <div2 type="scene" n="2">
12   <stage>A room in Pa Ubu's house, where a magnificent
13     collation is set out</stage>
14 </div2>
15 </div1>
16 <div1 type="act" n="2">
17   <head>Act Two</head>
18   <div2 type="scene" n="1">
19     <head>Scene One</head>
20   </div2>
21   <div2 type="scene" n="2">
22     <head>Scene Two</head>
23   </div2>
24 </div1>
25 </body>
```

# TEI

## Validation

There are several ways to validate a XML/TEI document

- Online validator: tells you whether the schema is respected (valid) and the XML is well-formed
- XML/TEI plugin for VS Code: tells you whether the XML is valid, well-formed, and suggests TEI elements
- Python library: a shell script to validate XML/TEI docs



# TEI

---

## Validation w/ VSCode plugin

- Download and install the plugin for VS Code
- **Read the instructions**
- Create a XML file in VS Code
- Include the minimal structure of TEI doc
- Include this XML/TEI sample
- Modify the elements and check the validity

NB. Relax ng (.rng) is yet another way to encode a XML schema, like XSD and DTD.



```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <?xml-model
3   href="https://vault.tei-
4   c.org/P5/current/xml/tei/custom/schema/relaxng/tei_all.rng"
5   schematypens="http://relaxng.org/ns/structure/1.0"
6   type="application/xml"?>
7 <TEI xmlns="http://www.tei-c.org/ns/1.0">
8   ...
9 </TEI>
```

# Homework

---

## Quiz time!

Answer the questions and get your results immediately.

<https://forms.gle/DNt9cPCxuUZzwxue8>