# Modelling humanities data with TEI-XML

## SCHOLARLY EDITING AND MANUSCRIPT CATALOGUING IN THE DIGITAL AGE

Dr Katarzyna Anna Kapitan
1 October 2025

# Introductions

- **Teacher**:
  - Dr Katarzyna Anna Kapitan
  - katarzyna.kapitan [at] chartes.psl.eu
- **Students**:
  - Name
  - Background and Study Programme
  - Research Interests
  - Experience with Markup Languages (which ones, what experience)

# Course overview

- **Classes**:
  - 10 x 2 hours between 1 October and 10 December
  - Wednesdays, 15:15-17:15
    - Besides 8 October -> **Time Change: 10 October @ 10:00 on Zoom (the session will be recorded for those who cannot make it)**

# Course overview

▶ **Course prerequisites:**

▶ Laptop with installed

▶ **Oxygen XML Editor** (https://www.oxygenxml.com/xml_editor/download_oxygenxml_editor.html)

▶ **Sublime Text** (https://www.sublimetext.com/download)

# Objectives

▶ Understand the XML structure and its customisation

▶ Use the TEI Guidelines for text editing and manuscript cataloguing

▶ Create XML templates and validation rules

▶ Document encoding choices

▶ Apply XML to individual projects

# Schedule

| | |
|---|---|
| **Session 1** | Markup Languages and Text Encoding |
| **Session 2** | Extensible Markup Language (XML) and Document Type Definition (DTD) |
| **Session 3** | Text Encoding Initiative (TEI) Guidelines |
| **Session 4** | Describing Primary Sources |
| **Session 5** | Transcribing Primary Sources |
| **Session 6** | Editing Primary Sources |
| **Session 7** | Annotating Primary Sources |
| **Session 8** | XPath |
| **Session 9** | Customisation and Documentation |
| **Session 10** | Customisation and Validation |

# Assessment method

▶ **Assessment:**

- Attendance and Participation (20%)
- Encoding Project Portfolio:
  - Assignment 1 (due 20 October): 15%
  - Assignment 2 (due 17 November): 25%
  - Assignment 3 (due 22 December): 40%

# Course Materials

▶ **GitHub Repo:**

https://github.com/KAKDH/TNAH_XML2025

# Text Encoding & Markup Languages

# Markup Languages

▶ **What is a markup language?**

   ▶ Computer language

   ▶ Uses tags to define elements within a document

   ▶ Human-readable

▶ **Examples**:

   ▶ HTML (Hypertext Markup Language)

   ▶ Markdown

   ▶ TeX & LaTeX

   ▶ **XML (Extensible Markup Language)**

   ▶ Scribe, GML (Generalized Markup Language) & SGML (Standard Generalized Markup Language)

# Markup

| | Semantic | Presentational |
|---|---|---|
| **LaTeX** | \selectlanguage{latin}{ad hoc} | \emph{ad hoc} |
| **HTML** | <i lang="la">ad hoc</i> | <i>ad hoc</i> |
| **XML** | <foreign xml:lang="la">ad hoc</foreign> | <hi rend="i">ad hoc</hi> |

# Exercises

▶ **Sync the GitHub repo (**https://github.com/KAKDH/TNAH_XML2025**) and access the Exercises Folder: Week1/Exercises/Ex1**

▶ **GitHub Tutorials:**

▶ https://docs.github.com/fr/repositories/creating-and-managing-repositories/cloning-a-repository

▶ https://dev.to/warish/a-step-by-step-guide-to-cloning-a-github-repository-to-your-local-machine-4ngc

▶ https://dev.to/phrawzty/easily-keep-your-local-git-repo-up-to-date-l52

# Exercise 1 (Week1/Exercises/Ex1)
## Individual Work (10 min)

▶ Which of these files contain markup?

▶ What type of markup is used?

▶ How is the markup expressed?

▶ How is the document structured?

▶ What are these types of markup usually used for, give examples?

▶ List all markup elements used in the documents and explain what are they used for (what do you think they mean)?

**Example 1:**
AM_30_fol.txt
AM_30_fol.md
AM_30_fol.html
AM_30_fol.doc

**Example 2:**
AM_30_fol.txt
AM_30_fol.md
AM_30_fol.tex
AM_30_fol.doc

# Group work (10 min)

▶ *Create two groups. One group consisting of everyone who worked on example 1 and one group consisting of everyone who worked on example 2.*

▶ Compare your answers and prepare a presentation of your example, choose one person who will present your answers.

▶ **Group with example 1:** focus on HTML

▶ **Group with example 2:** focus on TEX

# Plain Text | Formatted Text

- **Plain text** format contains no formatting information,

- **Formatted Text - Rich text** format includes formatting details such as font size, style, colour, etc.

  - **WYSIWYG** (**W**hat **Y**ou **S**ee **I**s **W**hat **Y**ou **G**et) - a formatted document as it will appear on screen or in print without showing the descriptive code.

# Typesetting with TeX & LaTeX

▶ TeX is a typesetting system widely used in academia, especially in mathematics, computer science, engineering, linguistics, and many more, among them: **digital scholarly editing.**

▶ **WYSIWYM** (What You See Is What You **Mean**) is a paradigm for editing a structured document separating presentation from content.

▶ TeX commands commonly start with a backslash and are grouped with curly braces.

  ▶ **\myCommand{ My Content goes here }**

## Saga af Hrómu⟨n⟩de Greipszýne

Cápituli .I.

(S)á kongr rieþe fyrir ⌐Gordom⌐ ⌐ı⌐ ⌐danmorc⌐ er Olafr hiet, hann var sonr Gnóþar Asmundar, hann var frǫgr maþr. Broþr .ij. kári oc ørnulfr, voru landvarnar menn kongz, hermenn mikler. Þar Bió eirn rýkur Boandi, sá hiet Greýpr. hann átte þá kono, er Gunnlǫþ hiet, dótter Hrókz hinz suarta, Þau átto .ix. sono er sva hieto. Hrólfr, Hake, Gautr(,) Þrostr, Angantýr, Logi, Hrómundr. Helge. Hrókr. Þeir voru aller efnileger menn. Þo var Hrómundr fyrir þeim aullom. hann kunni eigi at hroðast, hann var augna fagr, hárbiartr, oc herþamikill, mikill oc stercr, lýktiz miǫc Hróki móþr fauþr sýnom. Med kongi voru .ij. menn, hiet eirn Býldur, annar voli, Þeir voru Iller oc underfǫruler. kongr matti Þá mikils. Eytt sinn hiellt olafr kongr, austur fyrir noreg med her sinn, oc hielldo aþ Vlfaskerium, herioþo, oc lau viþ eitt Eýland. Kongr býþr Kára oc Ørnulfi aþ ganga uppá Eýuna, oc vita, huort þeir sæe einginn herskip. Þeir gengu uppá landiþ, oc litu .vj. herskip under hǫmrum nocrum. Þar var eirn dreke allskrautligr. Kári kallar til þeira, oc spir huorier fyrir skiponom rieþi, Eirn dólgr stóþ uppá drekanom, oc qvadz Hraungviþr heita. eþr huort er nafn þitt. Kari sagþi til sýn oc sýnz broþurz. oc mælti. Eg veit aungvann verri enn þik, oc þar fyrir skal ek hauggva þic i smá sticki, Hraungviþr. mælti: Ek hefi

---

2 rieþe] red T1768. 2 ⌐ı⌐ ⌐danmorc⌐] ÷ T1768, A345, P67, L222, B4859. 3 Asmundar] + d. A587, A193. 3 frǫgr maþr] maþr frægr A587, A193. 3 sá hiet Greýpr] er Greypur hiet A345. 9 eigi] ei A193, P67, L222. 11 kongr] ÷ A345. 13 austur] norþr A345, vestur P67. 15 sæe] sæo, A345, B4859. 15 einginn] eingi A587, engi A193, T1768, A345, L222, B4859. 16 herskip] skip A345, L222, B4859. 17 spir] + efter P67. 18 huorier] hver A587, A193, hvorir T1768, hvor A345, P67, L222, huerier B4859. 21 skal] vil P67. 21 Hraungviþr.

# LOST BUT NOT FORGOTTEN
## The Saga of Hrómundur and Its Manuscript Transmission

Katarzyna Anna Kapitan

Open Access Book

Digital Editions

Physical Copies

# HTML (Hypertext Markup Language)

▶ HTML is the standard markup language for web pages.

▶ Markup:

    ▶ **<tagname> Content goes here... </tagname>**

    ▶ Start Tag --- Content --- End Tag

```
<!DOCTYPE html>
<html>
  <head>
    <title> This is my title</title>
  </head>
  <body>
    <p>This is a paragraph</p>
  </body>
</html>
```

**Basic elements:**
- <html> – root element of an HTML page
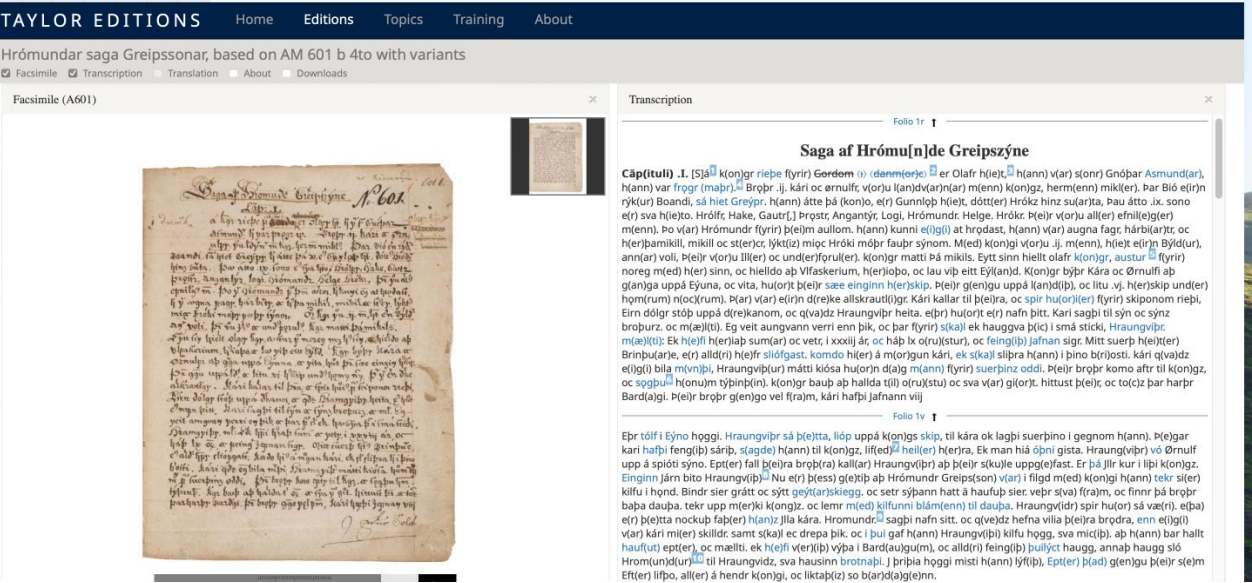- <head> – contains meta information about the document
- <body> – contains the visible page content

**Other elements (examples):**
- <div> - division/section of the page
- <table> - table
- <img> - image
- <ul> – list
- <h1> – heading
- <p> – paragraph

# HTML is everywhere !!!



https://editions.mml.ox.ac.uk/editions/hromundar_A601/

# Clear separation of content and presentation

<div>

<h2>Short Description</h2>

<p>AM 30 fol. is a paper manuscript in folio format comprising 56 leaves gathered into five quires of six conjoint leaves each, with the exception of the last quire, with consists of four conjoint leaves only. The whole codex is made of one type of relatively thick laid paper with a "bull's head" watermark.

</p>

</div>

# Clear separation of content and presentation

HTML + CSS (Cascading Style Sheets)

- ▶ HTML

<h1> My Text </h1>

- ▶ CSS

h1 {

    text-decoration: underline;

    text-align: center;

}

# Exercise 2 (Week1/Exercises/Ex2)

▶ Open AM_30_fol.html with your browser – take a screenshot of the view

▶ Open AM_30_fol.html with Sublime Text

▶ Associate style.css with your html file by adding within <head> the element <link>:

                                                               `<link rel="stylesheet" type="text/css" href="style.css">`

▶ Save the file.

▶ Open AM_30_fol.html in a browser – take a screenshot of the view.

```
1  <!DOCTYPE html>
2  <html>
3  <head>
4    <meta charset="utf-8"/>
5    <meta name="author" content="Katarzyna Anna Kapitan">
6  <title>Copenhagen, Den Arnamagnæanske Samling, AM 30 fol.</title>
7  </head>
8  <body>
9      <h1>Copenhagen, Den Arnamagnæanske Samling, AM 30 fol.</h1>
10 <div>
11 <h2>Short Description</h2>
12 <p>AM 30 fol. is a paper manuscript in folio format comprising 56 leaves gathered into five quires of six conjoint
   leaves each, with the exception of the last quire, with consists of four conjoint leaves only. The whole codex is
   made of one type of relatively thick laid paper with a "bull's head" watermark.  AM 30 fol. contains two main
   texts dealing with the Christianization of the Slavs.</p>
13 </div>
14 <div>
15   <h2>Contents</h2>
16 <p>
17 <ul>
18   <li><i>Chronica Slavorum</i> b
19   <li><i>Chronica Slavorum</i> b
20 </ul>
21 </p>
22 </div>
23 </body>
24 </html>
```

# Copenhagen, Den Arnamagnæanske Samling, AM 30 fol.

## Short Description

AM 30 fol. is a paper manuscript in folio format comprising 56 leaves gathered into five quires of six conjoint leaves each, with the exception of the last quire, with consists of four conjoint leaves only. The whole codex is made of one type of relatively thick laid paper with a "bull's head" watermark. AM 30 fol. contains two main texts dealing with the Christianization of the Slavs.
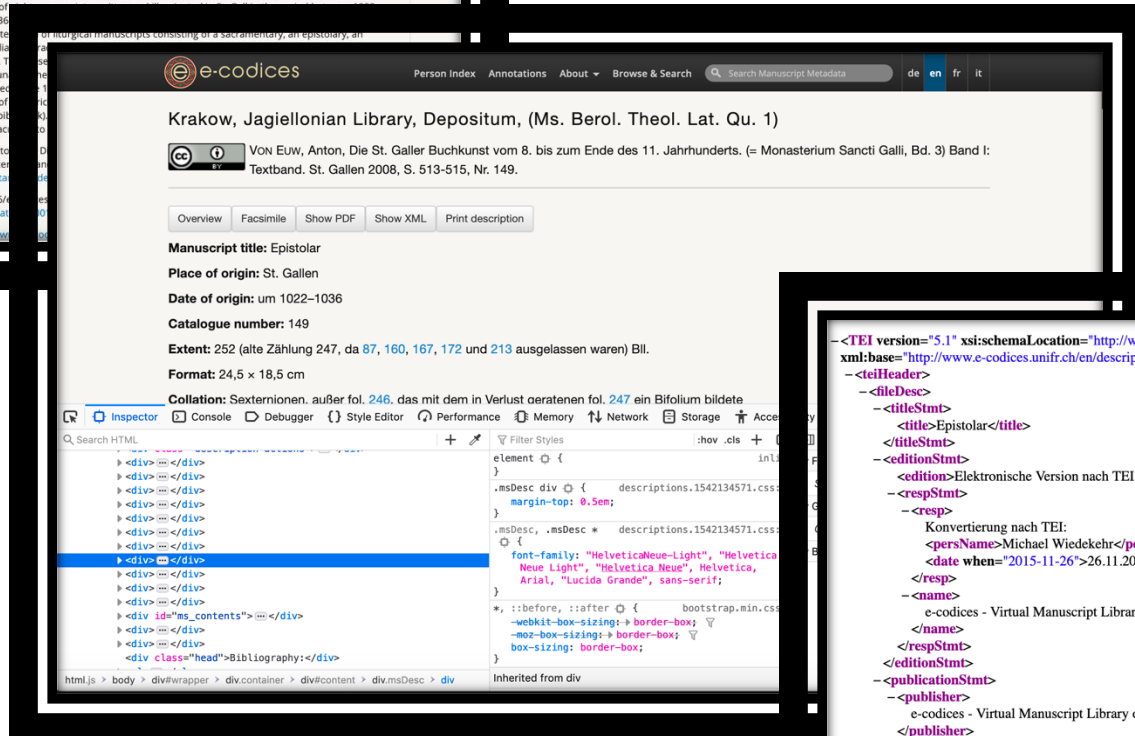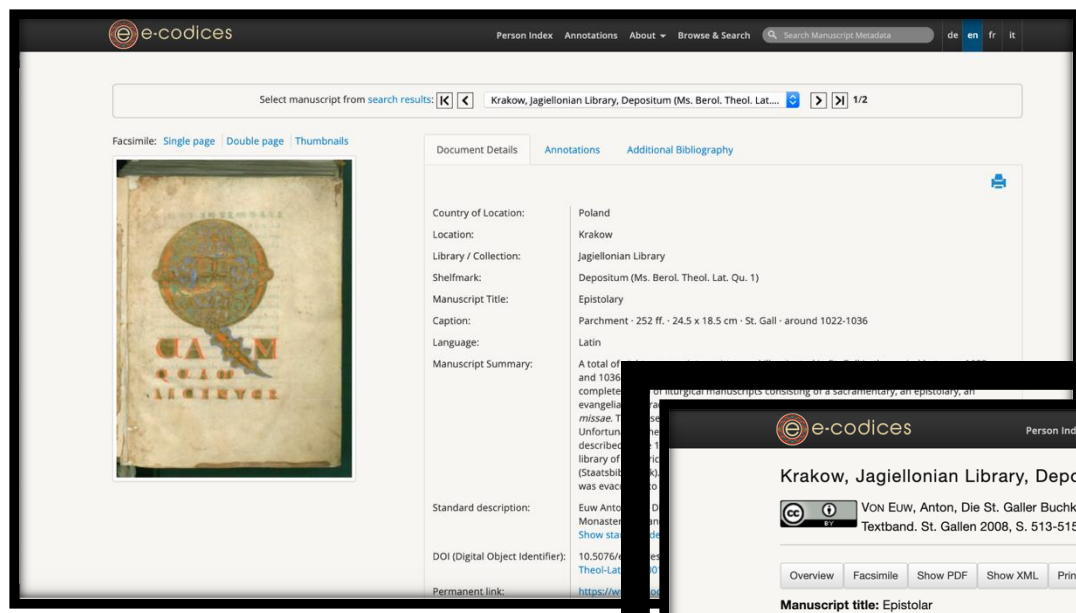
## Contents

- *Chronica Slavorum* by Helmold of Bosau
- *Chronica Slavorum* by Arnold of Lübeck

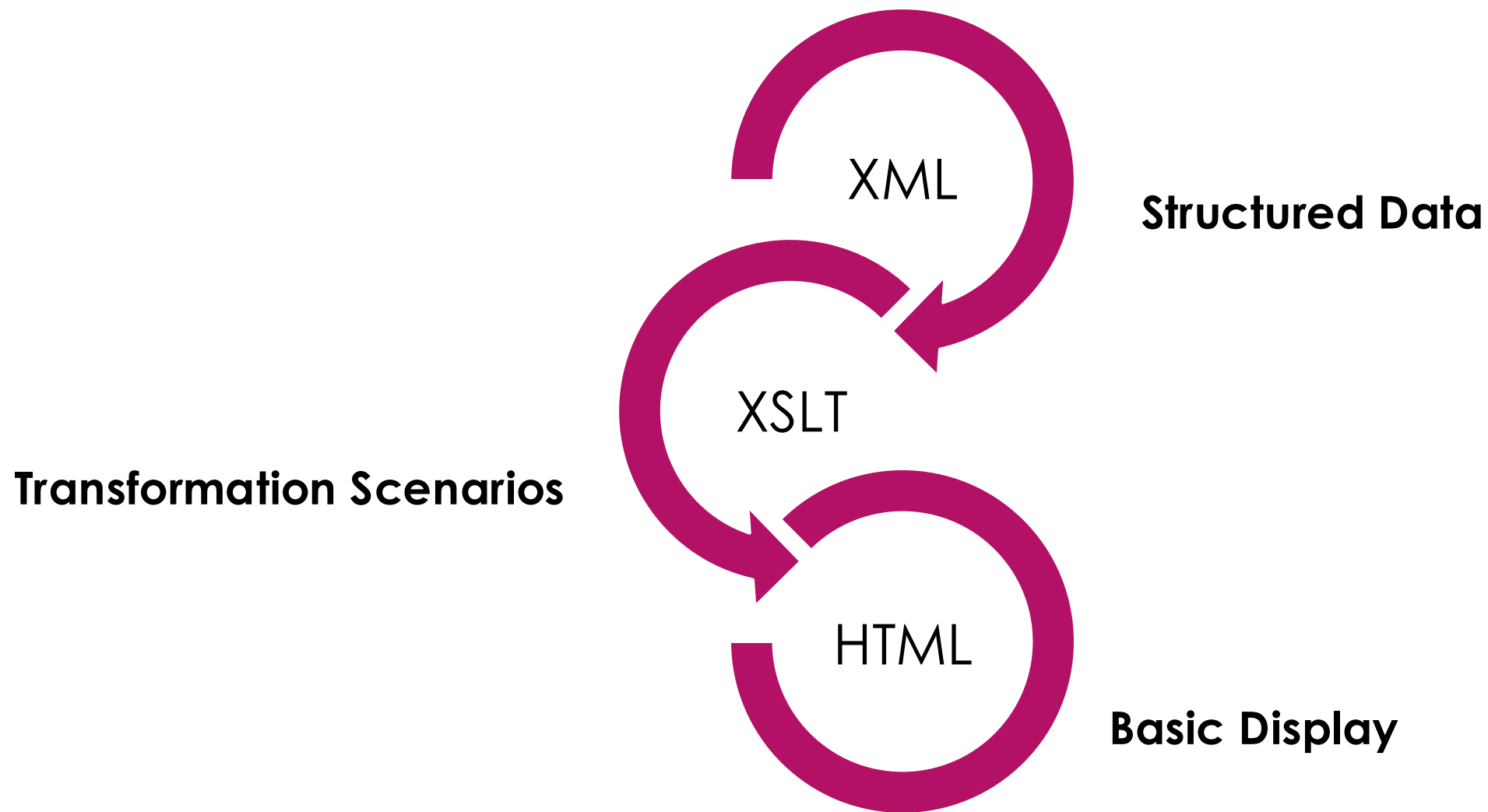# Markup is everywhere
## Behind the e-codices viewer:
## HTML & XML



Kraków, Jagiellonian Library, Ms. Berol. Theol. Lat. Qu. 1
Source: https://www.e-codices.unifr.ch/en/list/one/bj/Berol-Theol-Lat-Qu-0001/

# From Data to Display



XML — **Structured Data**

XSLT

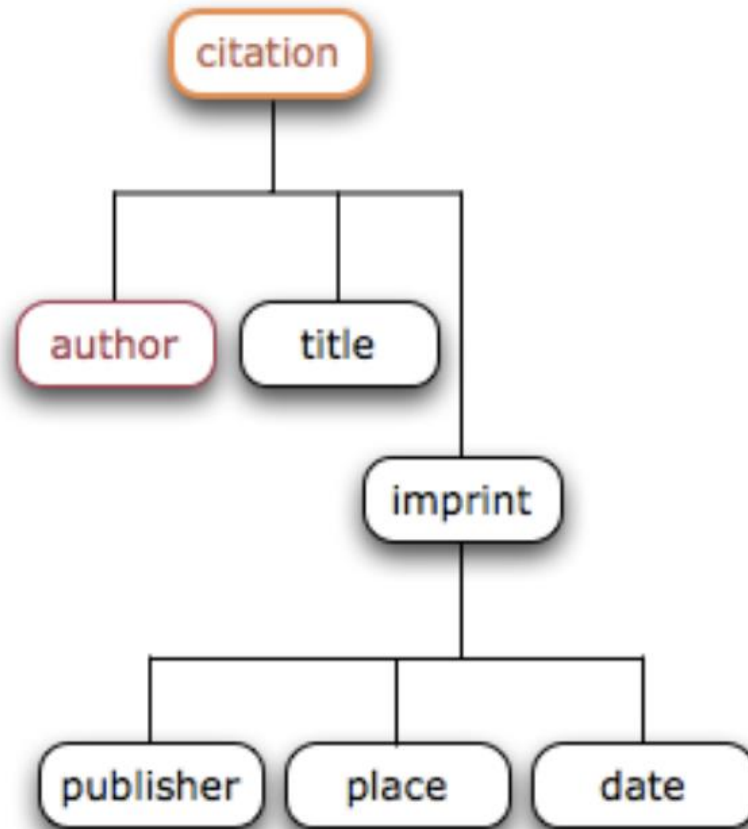**Transformation Scenarios**

HTML — **Basic Display**

# XML
# (Extensible Markup Language)

# XML (Extensible Markup Language)

▶ Storing structured data

▶ International standard, non-proprietary

▶ Standard text format (expressed in plain text)

▶ Easy to parse and read for computer programs.

▶ Widely used to export and share structured data.

▶ Hardware and software independent

Syd Bauman, Julia Flanders, and WWP, Creative Commons Attribution-ShareAlike 3.0 (Unported) license.

Syd Bauman, Julia Flanders, and WWP, Creative Commons Attribution-ShareAlike 3.0 (Unported) license.

# XML Elements & Tags

```
<!DOCTYPE html>
<html>
    <head>
        <title>This is a title</title>
    </head>
    <body>
        <p>This is a paragraph</p>
    </body>
</html>
```
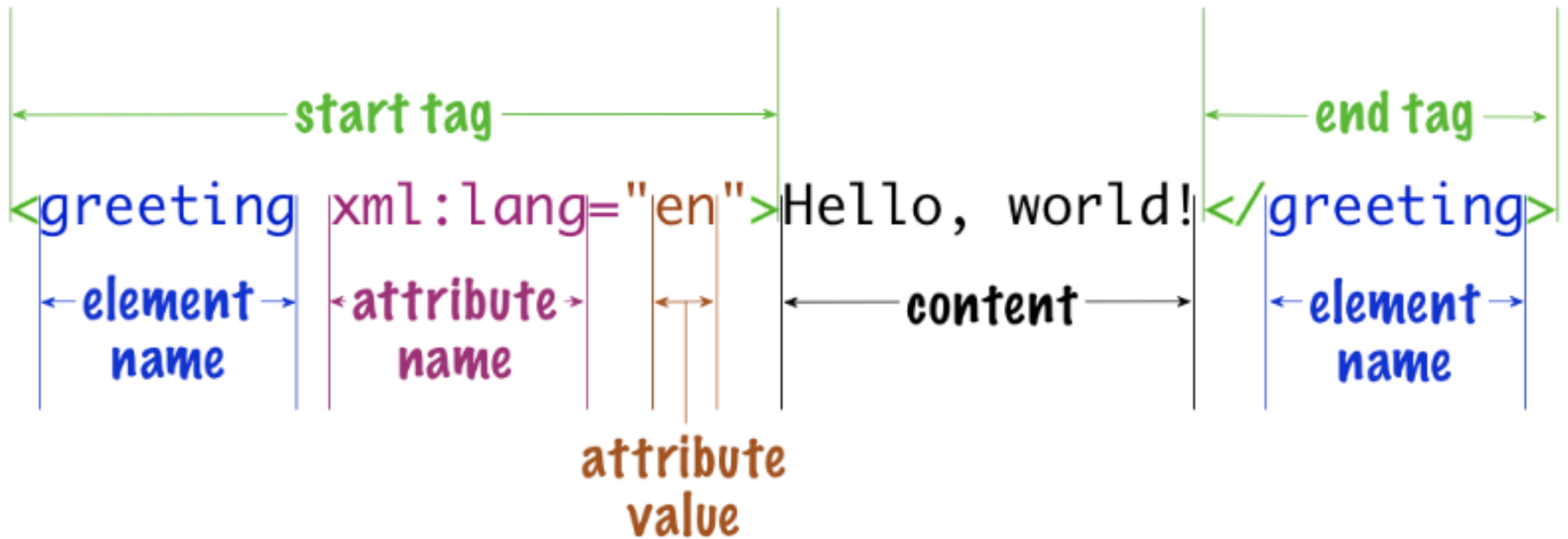
```
<?xml version="1.0" encoding="UTF-8"?>
<myRoot>
    <myContent>
        <content>
            Here is my content
        </content>
    </myContent>
</myRoot>
```

# XML Elements & Tags

- Text is divided into elements (the nouns of the encoding — content objects).

- elements have start-tags and end-tags

  - &lt;heading&gt;My heading&lt;/heading&gt;

- start-tags have < … >

  - &lt;heading&gt;

- end-tags have </ … >

  - &lt;/heading&gt;

# XML Attributes

▶ Attributes are adjectives of XML, they describe the properties of the elements

▶ any number of attributes can be specified on a given start-tag

▶ but only one with a given name

▶ **<person job="musician" age="55">**Paul Simon**</person>**

Syd Bauman, Julia Flanders, and WWP, Creative Commons Attribution-ShareAlike 3.0 (Unported) license.

# Quiz

## I. Which of the following examples are well-formed XML?

1. \<name type="person>Pearl S. Buck\</name>

2. \<name type="person">Toni Morrison\<name>

3. \<name="person">Carl Sagan\</name>

4. \<name type="person">Kurt Vonnegut\</name>

5. \<name type=person>John Cleese\</name>

6. \<name type="person">\<forename>Frances\</forename> \<surname>Perkins\</surname>\</name>

**Course Materials:**
https://github.com/KAKDH/TNAH_XML2025

**Syllabus:**
https://github.com/KAKDH/TNAH_XML2025/blob/main/Kapitan_TEI-XML_Syllabus.pdf

**For Next Class:**
1. Get your Git up and running
2. Get your Oxygen up and running