**TEI-XML Basics**

**XML vs TEI**

XML (Extensible Markup Language) is a coding language that allows us to store, transmit, and reconstruct data in a way that is both human-readable and machine-readable. For our purposes, this means we can markup our data (i.e., the text) with notations. TEI (Text Encoding Initiative) is a flavor of XML, which means it runs on the same structure as XML, but it includes elements that are specific for the needs of marking up texts.

Elements (also known as tags) in XML and TEI need to be formed correctly. If they don't, then the code will not be valid and will generate error codes.
- ❖ Generally speaking, elements require two parts: an opening element and a closing element. A closing element will come second and will have a backslash at the beginning of the element.
  - ➢ Ex. <persName>Henry</persName>.
- ❖ If an element does not need to be wrapping a text, sometimes it will only have one element. When this is the case, the backslash will be included at the end of the element.
  - ➢ Ex. <space rend="indent1"/>
- ❖ The majority of elements that we will find in our TEI-XML documents are transformed from the paragraph styles and character styles we use in InDesign. When we are auditing our TEI-XML documents, it helps to have the corresponding InDesign document open.
  - ➢ It is important to note that, when elements are being drawn from **paragraph styles**, a new element will open after every end of paraph mark. This is why there are three separate opening and closing elements for <p style="body"> when we have three paragraphs in a row all styled with the Body paragraph style.
- ❖ Tags often need to be in a certain order in order for a document to be well-formed and validate.
  - ➢ Many tags are part of a parent-child relationship.
    - ■ Parent tags are containers that can hold a variety of smaller tags.
    - ■ Child tags are smaller tags that need to be surrounded by larger parent tags.
    - ■ Most of these parent-child relationships are established in our code by our transformation script.
  - ➢ We also nest tags that don't necessarily have parent-child relationships within each other.
    - ■ Move from general to specific: an element that applies to a full line of text should be on the outside, and an element that applies to only one word in that line should be within.
      - ● Example: Text that should be bold but then there is a Person!

- Render as &lt;hi rend="bold"&gt;Text that should be bold but then there is a &lt;persName&gt;Person&lt;/persName&gt;!&lt;/hi&gt;

It is important to note that neither XML nor TEI are display languages. This means that any styling we need to do for the text to display properly on our website is connected to a CSS (cascading style sheet) where we define how we want styles to look.

Below, we'll go through some of the basics of what to expect from an XML document. For a (*very*) detailed overview of TEI, its structure, and all its elements, take a look through the TEI Guidelines available on tei-c.org.

**Navigating an XML File**

When you open an XML file in oXygen, the first few lines will typically say something like:
&lt;?xml version="1.0" encoding="utf-8"?&gt;
&lt;TEI xmlns="http://www.tei-c.org/ns/1.0" version="3.3.0"&gt;

These first two lines define the entire document's structure. They indicate which version of XML we are using, as well as which version of TEI we are using. Knowing which version we are following allows the program to check whether the rest of the content validates.

The rest of the document is structured in two overarching containers: the &lt;teiHeader&gt; and &lt;text&gt; elements.

&lt;teiHeader&gt;
The TEI header holds all important metadata specific to its particular XML file. This information is important for our metadata transparency. At this stage (intermediate, pre-launch), we do not need to audit this information.

&lt;text&gt;
The text container holds all major information and data for the XML document. The opening &lt;text&gt; element will always be followed by &lt;body xmlns="http://www.tei-c.org/ns/1.0"&gt;, a file-specific &lt;anchor&gt; element, and a &lt;div&gt; element.

```
270 ▽   <text>
271 ▽     <body xmlns="http://www.tei-c.org/ns/1.0">
272         <anchor xml:id="Hudson-Caxton-03-Gen-Intro"/>
273 ▽       <div>
```

The &lt;div&gt; elements are text division elements. They indicate large subdivisions within the &lt;body&gt; of the &lt;text&gt; elements. We use them primarily to separate the contents of the text.

Usually, there will be two main <div> elements: one for the main content and one for footnotes. Exceptions to this are prose texts with chapters, where each chapter receives its own <div> along with its own <div> for footnotes, and facing page texts. Each <div> element should receive a type attribute with a one word description of the type of content to be found within.

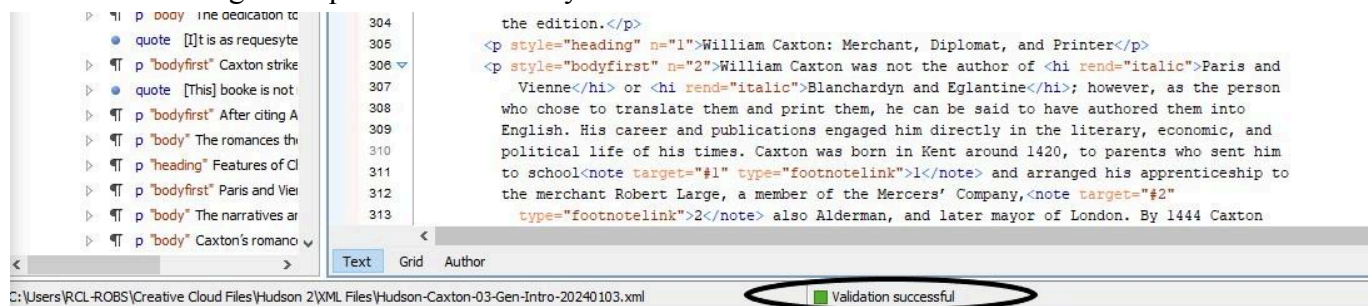❖ Ex. <div type="footnotes">

**<u>Using oXygen XML Editor</u>**

We use the program oXygen XML Editor to view, audit, and edit our TEI-XML files. This program is the standard in the field for editing XML files. It has a lot of functionality to it, which comes with many, many buttons. We only need to use a few of these buttons for our purposes.
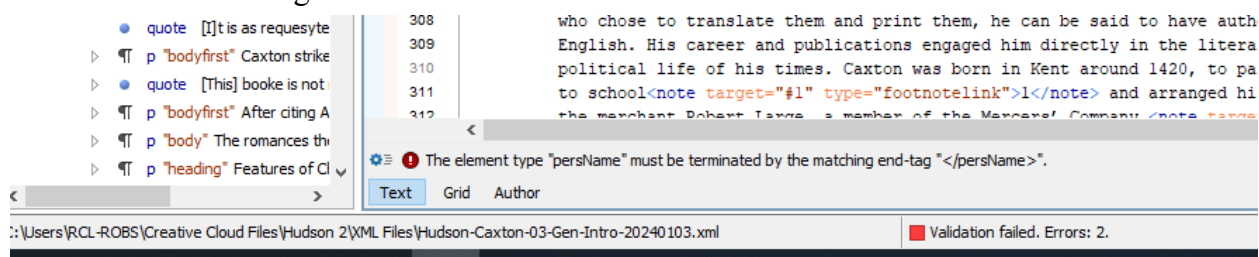
<u>Validation</u>
If an XML file is correct and correctly formed (aka "well-formed"), then it validates. You can check this a couple ways in oXygen.

❖ oXygen displays a status update on whether a document validates as you make changes within it. Look at the bar along the bottom of the screen.
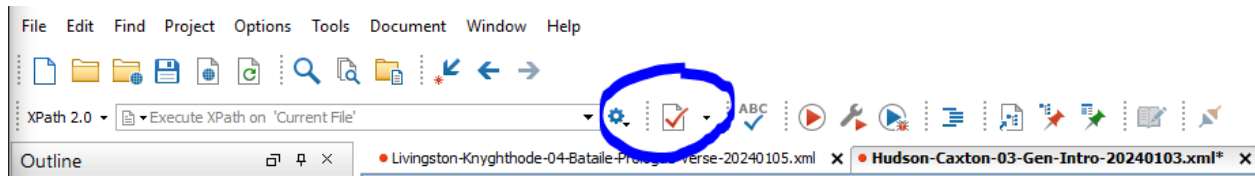❖ The status update on validation lives in the middle of the bar. There will be one of two results.
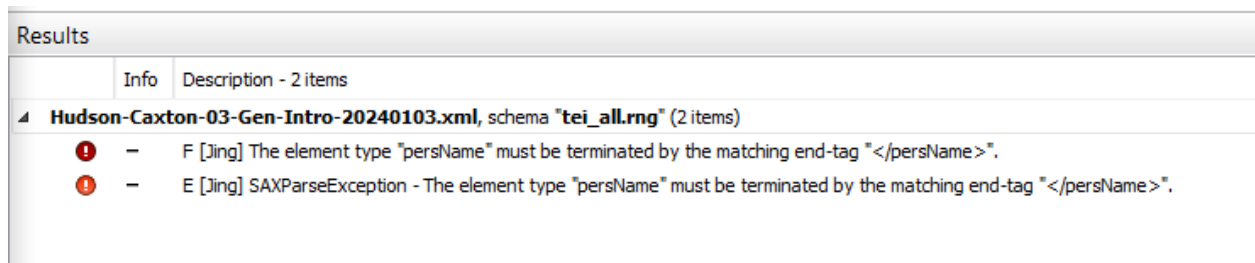➢ A green square and it will say Validation Successful.



➢ A red square that says Validation failed, followed by the number of errors in the document. There will also throw up an alert along the bottom of the screen detailing what the error is.

❖ You can manually check the validation of the file by clicking the button that looks like a document with a red check mark:
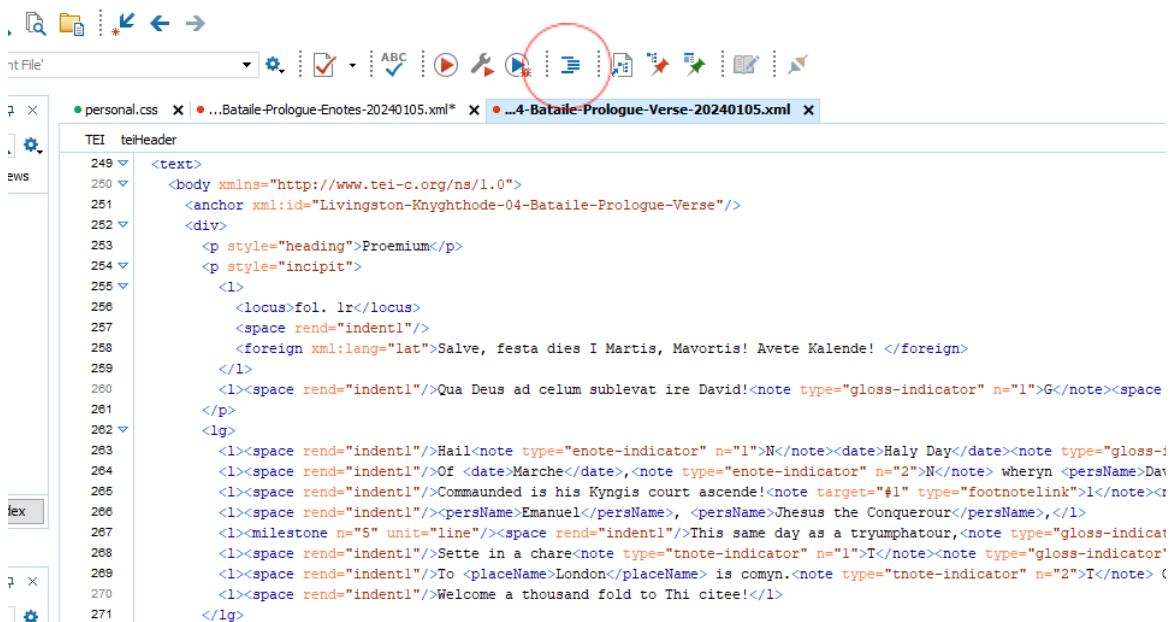


❖ This will pull up a Results bar along the bottom of the window. This will detail what errors there are, as well as where they exist.



Format and Indent

When you first open an XML document, it prioritizes putting each element on its own line. This makes it good for seeing the structure quickly, but it's very difficult for us to navigate and audit. Instead, we'll want to change the way the document looks in oXygen. Click the Format and Indent button – it looks like a list with indents – in the top menu.

Then, the document will look like this:

```
TEI   teiHeader

  1   <?xml version="1.0" encoding="utf-8"?>
  2 ▽ <TEI xmlns="http://www.tei-c.org/ns/1.0" version="3.3.0">
  3 ▶│   <teiHeader> [266 lines]
270 ▽ │ <text>
271 ▽ │   <body xmlns="http://www.tei-c.org/ns/1.0">
272   │     <anchor xml:id="Livingston-Knyghthode-04-Bataile-Prologue-Verse"/>
273 ▽ │     <div>
274   │       <p style="heading">Proemium</p>
275 ▽ │       <p style="incipit">
276 ▽ │         <l>
277   │           <locus>fol. 1r</locus>
278   │           <space rend="indent1"/>
279   │           <foreign xml:lang="lat">Salve, festa dies I Martis, Mavortis! Avete Kalende! </foreign>
280   │         </l>
281 ▽ │         <l><space rend="indent1"/>Qua Deus ad celum sublevat ire David!<note
282   │             type="gloss-indicator" n="1">G</note><space rend="indent1"/><gloss n="1">(See next
283   │             three lines)</gloss></l>
284   │       </p>
285 ▽ │       <lg>
286 ▽ │         <l><space rend="indent1"/>Hail<note type="enote-indicator" n="1">N</note><date>Haly
287   │             Day</date><note type="gloss-indicator" n="2">G</note> devout! Al hail kalende<note
288   │             type="gloss-indicator" n="3">G</note><space rend="indent1"/><gloss n="2">Holy
289   │             Day</gloss><gloss n="3"> kalends (i.e., first day)</gloss></l>
290 ▽ │         <l><space rend="indent1"/>Of <date>Marche</date>,<note type="enote-indicator" n="2"
291   │             >N</note> wheryn <persName>David</persName> the Confessour</l>
292 ▽ │         <l><space rend="indent1"/>Commaunded is his Kyngis court ascende!<note target="#1"
293   │             type="footnotelink">1</note><note type="enote-indicator" n="3">N</note></l>
294 ▽ │         <l><space rend="indent1"/><persName>Emanuel</persName>, <persName>Jhesus the
295   │             Conquerour</persName>,</l>
296 ▽ │         <l><milestone n="5" unit="line"/><space rend="indent1"/>This same day as a
297   │             tryumphatour,<note type="gloss-indicator" n="4">G</note><space rend="indent2"/><gloss
298   │             n="4">supreme victor</gloss></l>
299 ▽ │         <l><space rend="indent1"/>Sette in a chare<note type="tnote-indicator" n="1">T</note><note
300   │             type="gloss-indicator" n="5">G</note> and throne of majestee,<space rend="indent1"
301   │             /><gloss n="5">chair</gloss></l>
302 ▽ │         <l><space rend="indent1"/>To <placeName>London</placeName> is comyn <note
```

### Clicking on tags to see starting and ending elements
A quick way to see where elements open and close is by clicking in one of the tags. When you click into a tag, the entirety of the element will be underlined and a bracket will appear in the left margin, connecting the opening and closing elements.

```
320   │       cultural diffusion from the continent to England, as well as one of the most important
321   │       cities in the territories of the dukes of Burgundy. Here staples and luxury goods of all
322   │       kinds from around the world were bought and sold.<note target="#6" type="footnotelink"
323   │         >6</note> Caxton's office in the Merchant Adventurers required him to negotiate both
324   │       international contracts and local agreements with London, the ports of Flanders, the
```

### Collapsing elements
Once XML files are formatted and indented, they can be very lengthy. It can be easier to navigate the documents by collapsing down larger elements.
   ❖ Along the left margin is a line count for the code. Note the arrows beside some of the numbers. These indicate tags that have closing elements on a different line.

❖ Click the down arrow to collapse an entire section so it's not in your way.
  ➢ We recommend you do this for the &lt;teiHeader&gt;, since it won't need to be audited for this stage.
  ➢ Click into &lt;teiHeader&gt; and notice how it is underlined and a bracket appears, stretching far down the page. You can use the **down arrow** next to line 3 to **collapse** this entire header section so it's not in your way.

```
1   <?xml version="1.0" encoding="utf-8"?>
2 ▽ <TEI xmlns="http://www.tei-c.org/ns/1.0" version="3.3.0">
3 ▽   <teiHeader>
4 ▽     <fileDesc>
5 ▽       <titleStmt>
6           <title>Of Knyghthode and Bataile</title>
```

● Ex: The cursor is in the &lt;lg&gt; on XML line 285, so it appears with a small black bracket beneath it. This line group closes on XML line 305, so the &lt;/lg&gt; is also bracketed there. To the left, these lines are highlighted by another small black bracket stretching from line 285 to 305.

```
284        </p>
285 ▽      <lg>
286 ▽        <l><space rend="indent1"/>Hail<note type="enote-indicator" n="1">N</note><date>Haly
287              Day</date><note type="gloss-indicator" n="2">G</note> devout! Al hail kalende<note
288              type="gloss-indicator" n="3">G</note><space rend="indent1"/><gloss n="2">Holy
289              Day</gloss><gloss n="3"> kalends (i.e., first day)</gloss></l>
290 ▽        <l><space rend="indent1"/>Of <date>Marche</date>,<note type="enote-indicator" n="2"
291              >N</note> wheryn <persName>David</persName> the Confessour</l>
292 ▽        <l><space rend="indent1"/>Commaunded is his Kyngis court ascende!<note target="#1"
293              type="footnotelink">1</note><note type="enote-indicator" n="3">N</note></l>
294 ▽        <l><space rend="indent1"/><persName>Emanuel</persName>, <persName>Jhesus the
295              Conquerour</persName>,</l>
296 ▽        <l><milestone n="5" unit="line"/><space rend="indent1"/>This same day as a
297              tryumphatour,<note type="gloss-indicator" n="4">G</note><space rend="indent2"/><gloss
298              n="4">supreme victor</gloss></l>
299 ▽        <l><space rend="indent1"/>Sette in a chare<note type="tnote-indicator" n="1">T</note><note
300              type="gloss-indicator" n="5">G</note> and throne of majestee,<space rend="indent1"
301              /><gloss n="5">chair</gloss></l>
302 ▽        <l><space rend="indent1"/>To <placeName>London</placeName> is comyn.<note
303              type="tnote-indicator" n="2">T</note> O <persName>Saviour</persName>,</l>
304          <l><space rend="indent1"/>Welcome a thousand fold to Thi citee!</l>
305        </lg>
306 ▽      <lg>
```

● When you collapse this line group, 285 will collapse and the next visible line will be 306. oXygen also helpfully tells you how many lines you've collapsed, so look for the bracketed [x number of lines] in case you fear things have vanished.
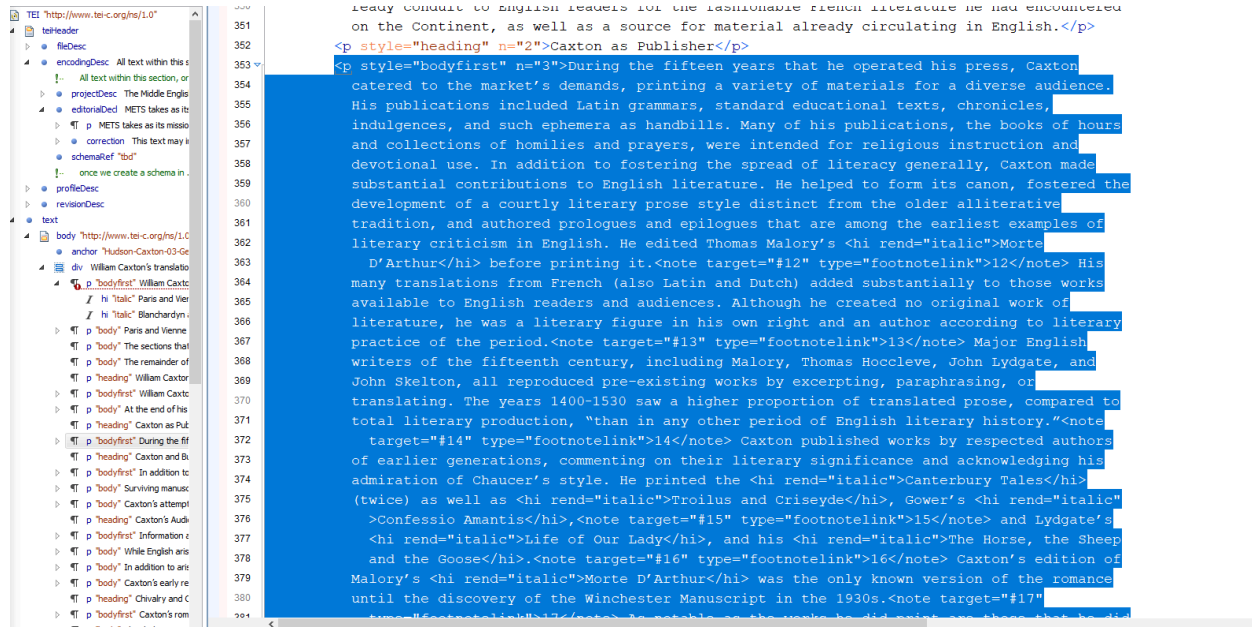
```
284        </p>
285 ▶      <lg> [20 lines]
306 ▽      <lg>
307 ▽        <l><space rend="in
```

**N.B.:** If you are having trouble navigating through the main elements in an XML file, it helps to have the Outline window open.

❖ If it's not already open, then click Window > Show View > Outline
❖ An outline will appear along the left side of your screen.
❖ When you click on an element in the outline view, then the main display of the XML file will jump to the element you have selected, as well as highlight the entirety of the element in blue.



The color coding system

oXygen color-codes information to make it easier to parse the code at a glance. You will typically see the following colors in your XML documents:

❖ Black: any content/data contained by the code
❖ Dark blue: main elements used in TEI-XML
❖ Orange: attributes that help define the elements
❖ Dark red: definitions for the attributes
❖ Green: a comment made within the code. This is essentially a note-to-self that doesn't display.
❖ Light pea-green: an escaped character.

Ex.

```
<p style="bodyfirst" n="1">W
```

In this case, the "p" is the element, the orange "style=" and "n=" are attributes, and "bodyfirst" and "1" are the further definitions for the attributes.

Escaped characters
Certain characters, such as ampersands, are command characters for oXygen, which are used to signal the code to perform some sort of action. In order for these characters to display in the content, we have to use the unicode character code. Some diacritics also display through escaped characters. Anything represented by unicodes will display as a light pea-green color.
Ex:

```
.c">wher-thurgh</hi>. (adv. &amp; conj.), sense 3a(c).
```