

# Text Encoding Initiative Workshop: Intro to Text Encoding

Michelle Dalmau & John Walsh, Indiana University  
Catapult / Scholars' Commons Production

# Slides, Exercises & Additional Primers

---

□ <http://dcl.slis.indiana.edu/teiworkshop/>

# Overview: Introduction to TEI

---

- Introduction to text encoding
  - the what and why
- Introduction to the Text Encoding Initiative / TEI
  - the how

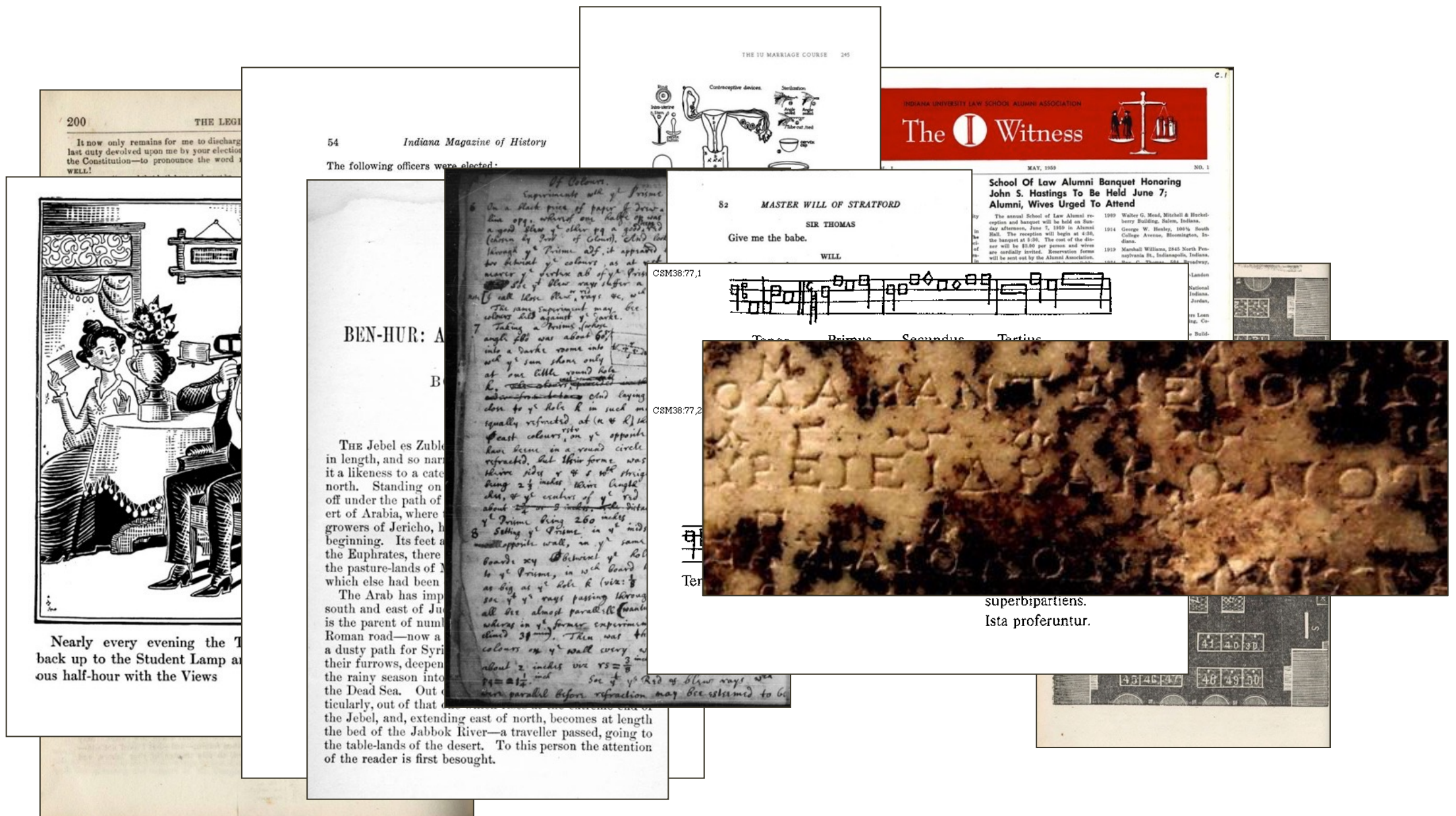
# Motivations for Text Encoding

- Access & preservation
- Discovery & dissemination
  - Searching/browsing
  - Interoperability & portability: harvesting/repurposing
- Analysis
  - Linguistic analysis
  - Concordances
  - Topic models
- Visualization
  - Interactive timelines (see [VWWP](#))
  - Map-based interfaces (see [Swinburne Project](#))

# Representing the Text with Encoding

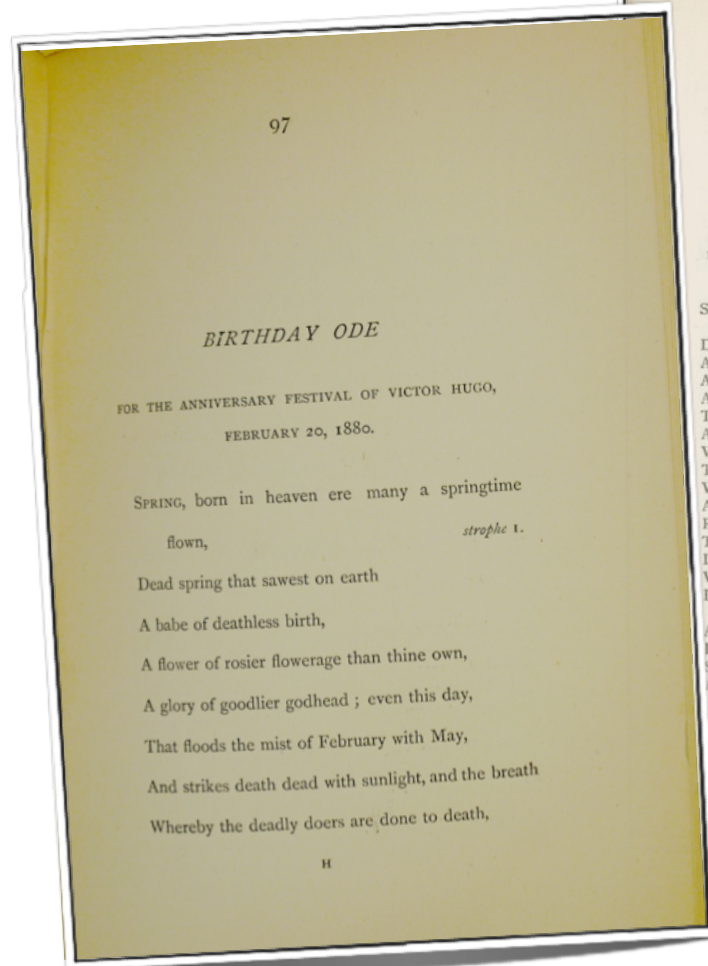
- Structural features
  - Text divisions (chapters, sections, etc.), paragraphs, lists, tables, line groups, lines, etc.
- Content & context
  - Metadata for the electronic and for the source document
  - References to people, places, events, organizations, etc. within the text (phrase-level)
  - Thematic and interpretive annotation
- Formatting & design
  - Bold, italics, small case, indentations, color, dimensions, binding, watermarks, and other features of the material document

# What is a Textual Document?

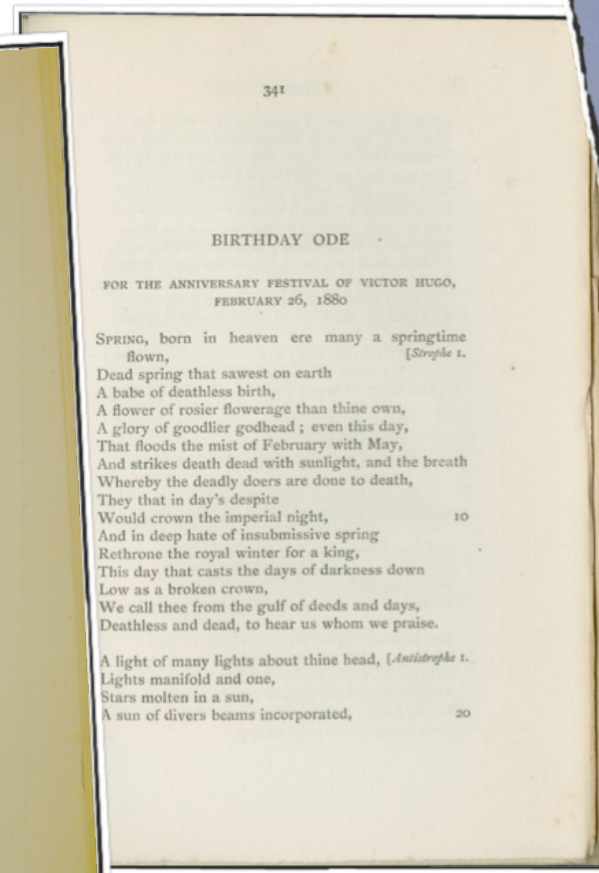




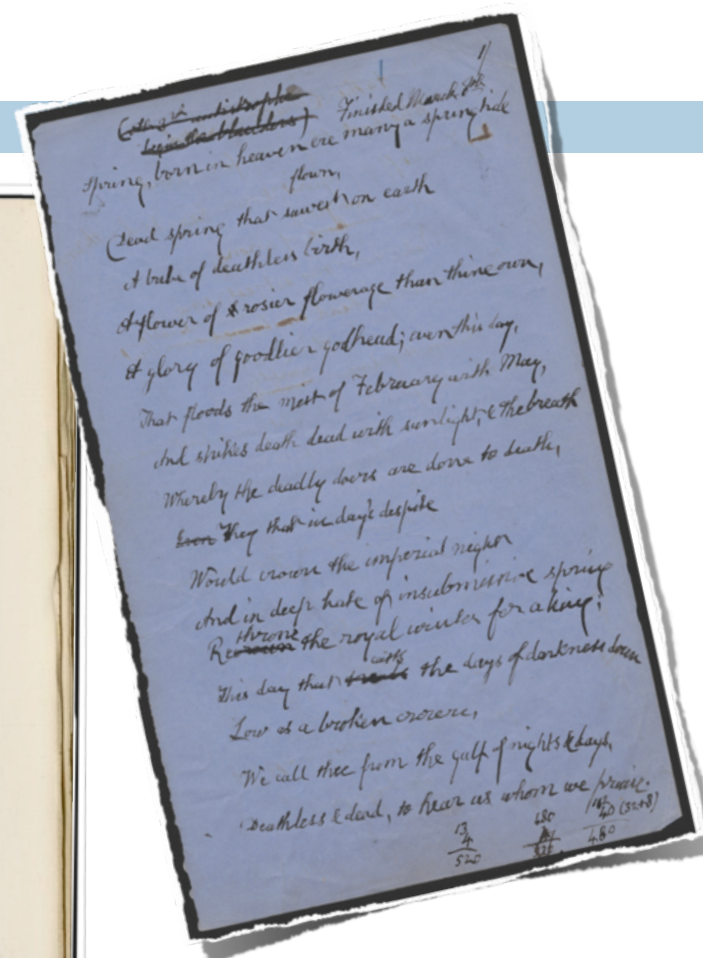
# Variants



Swinburne's Songs of the Springtides  
(1880)

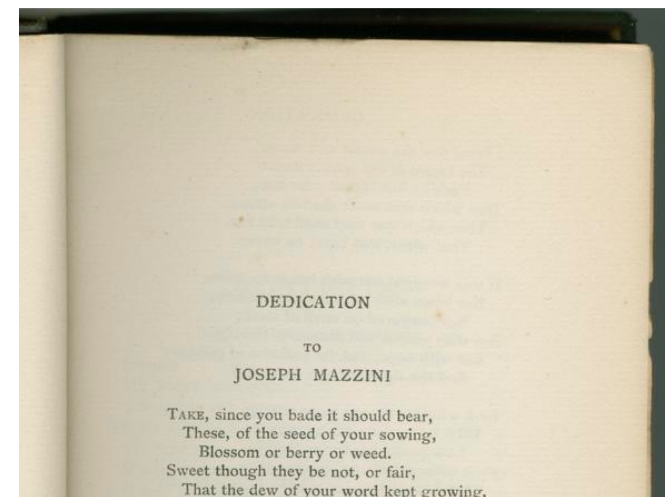
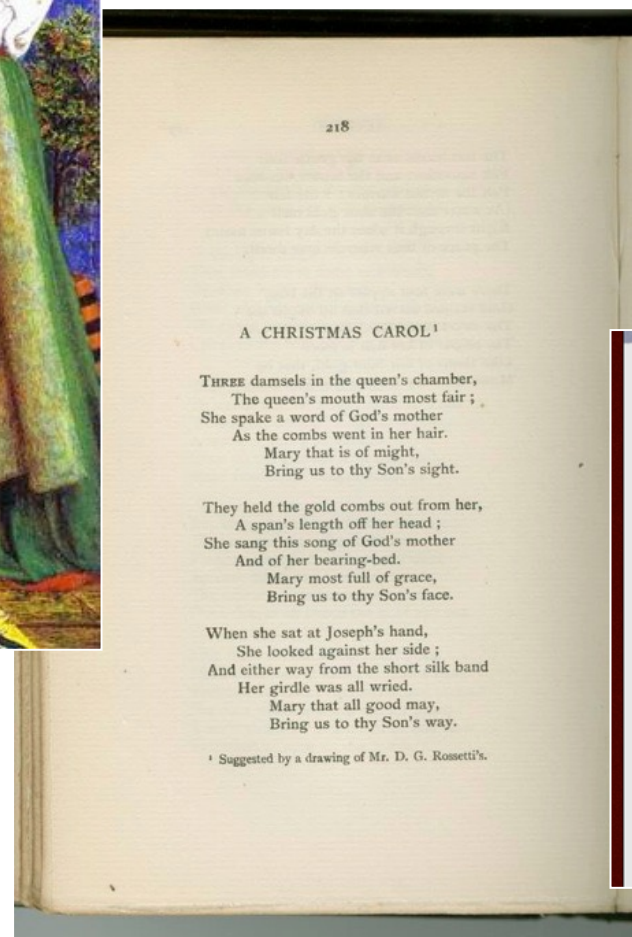


Swinburne's Poems (1904)



MS. Special Collections  
Research Center.  
Syracuse University  
Library

# Intertextual and Contextual Information



<< Back to Search Results

Search within this document:

1 occurrence of mazzini [Clear Hits]

Algernon Charles Swinburne: The Poems of Algernon Charles Swinburne

- 1 Dedication to Joseph Mazzini
- 2 Songs Before Sunrise
- 3 Songs of Two Nations

[show document information]

Dedication to Joseph Mazzini  
Algernon Charles Swinburne

page: [v]

DEDICATION  
TO  
JOSEPH MAZZINI<sup>1</sup>

TAKE, since you bade it should bear,  
These, of the seed of your sowing,  
Blossom or berry or weed.  
Sweet though they be not, or fair,  
That the dew of your word kept growing,  
Sweet at least was the seed.

Men bring you love-offerings of tears,  
And sorrow the kiss that assuages,  
And slaves the hate-offering of wrongs,  
And time the thanksgiving of years,  
And years the thanksgiving of ages;  
I bring you my handful of songs.

If a perfume be left, if a bloom,  
Let it live till Italia be risen,  
To be strewn in the dust of her car  
When her voice shall awake from the tomb  
England, and France from her prison,

Mazzini, Giuseppe (1805-1872)  
Italian patriot, idolized by Swinburne.

Resources:  
<http://en.wikipedia.org/wiki/Mazzini>

10



# Advantages of Text Encoding

- Re-use and flexibility: build once, use many
- Presentation and output of text controlled by style sheets.
  - Generate different views of the same text and different formats: PDF, HTML, ePub (ebooks), plain text (for text analysis), etc.
- The document *and* the markup can serve as an object of analysis and increased discoverability

# Interpretative & Project-specific Encoding

- Text encoding is not necessarily simple data entry/capture; it is not objective but interpretive. Every encoded text is a “reading” of the text.
- There are often many ways to apply a particular markup language to a particular text.
- Individual projects typically require project-specific guidelines and documentation in addition to the general markup language specification or guidelines.

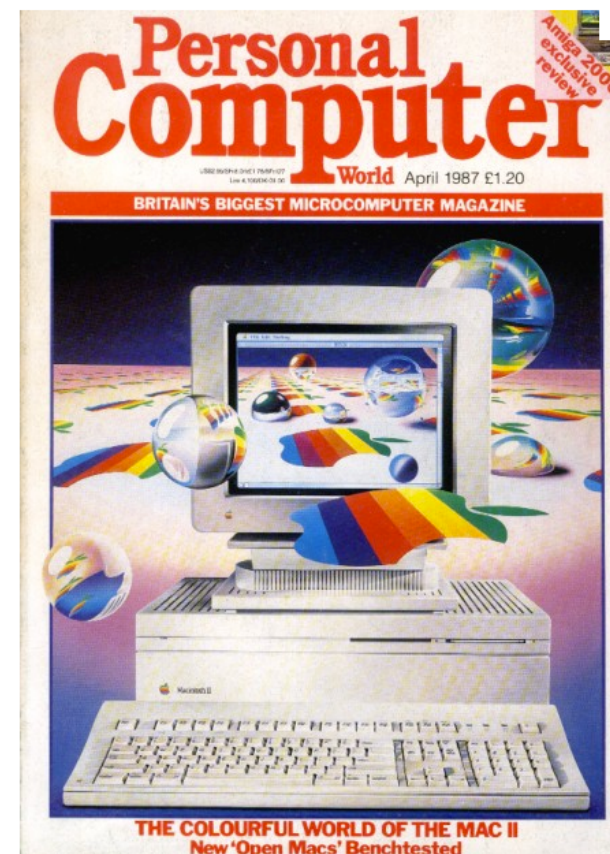
# Intro to the Text Encoding Initiative (TEI)

- TEI is:
  - a formally constituted organization, the TEI Consortium;
  - a scholarly community—with an annual conference, open-access journal, and active email discussion list.
  - a text encoding standard produced by that organization, TEI's Guidelines for Electronic Text Encoding and Interchange.
- For our purposes, TEI refers to the technical text encoding standard.

# History of TEI

Vanhoutte, Edward. "An Introduction to the TEI and the TEI Consortium." *Literary and Linguistic Computing* 19.1 (2004): 9-16.

*It all started in 1987 ...*





# Questions?



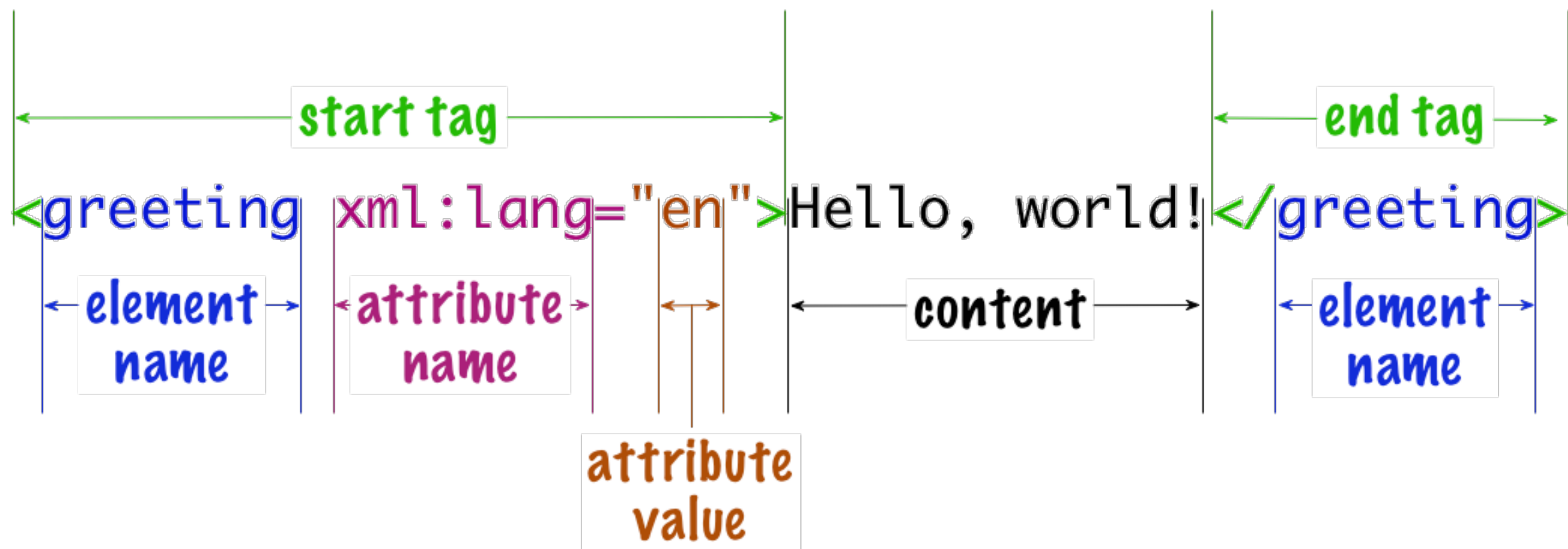
# Quick Introduction to XML

- XML, or eXtensible Markup Language, is a **non-proprietary meta language** for creating markup languages suited for different tasks, domains, and disciplines.
- An XML markup language consists of "tags" used to define the structure and other features of a text.
- HTML:
  - `<p>(paragraph of text)</p>`
  - ``
  - `<a href="http://www.indiana.edu">Indiana University</a>`
- TEI:
  - `<sp who="#rosamond"> (speech) </sp>`
  - `<lg> (line group, stanza) </lg>`
  - `<salute>Dear Fred,</salute>`

# XML Key Terms

- **Elements** are the basic, named structural units of an XML document (**nouns of encoding**)
  - `<title>The Odyssey</title>`
- **Attributes** are name/value pairs (name="value") associated with elements (**adjectives of encoding**)
  - `<creator type="author">Homer</creator>`
  - An element may have multiple attributes
- **DTDs** (Document Type Definitions) and Schemas define the rules that govern a particular type of XML document. They declare elements and attributes and the allowable content for those elements and attributes (**grammar rules**)

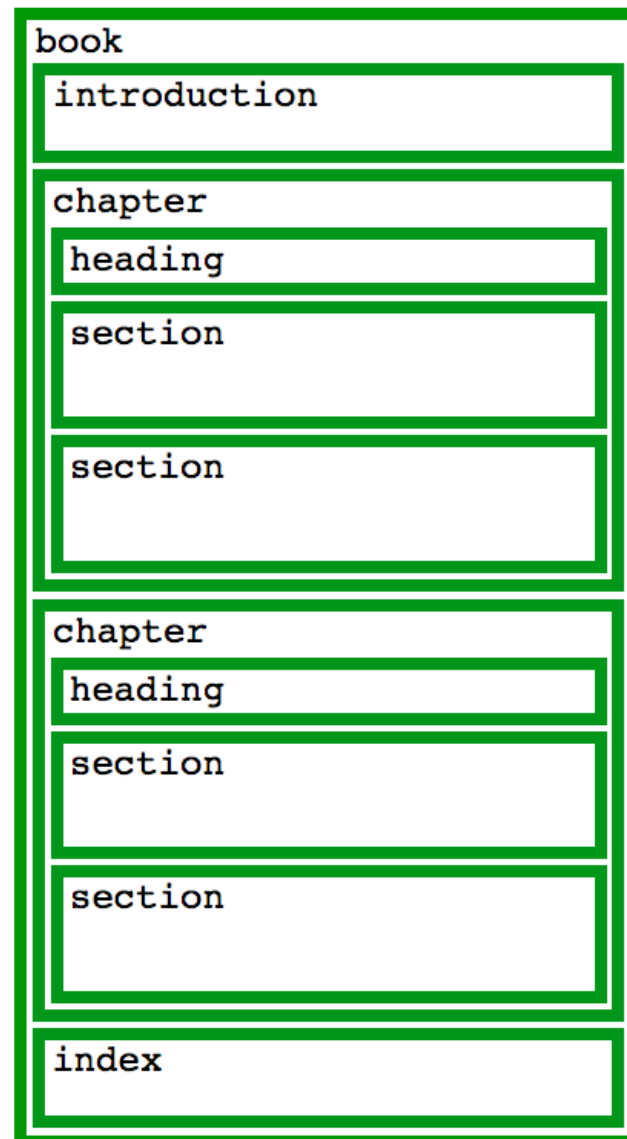
# XML: Anatomy of an Element



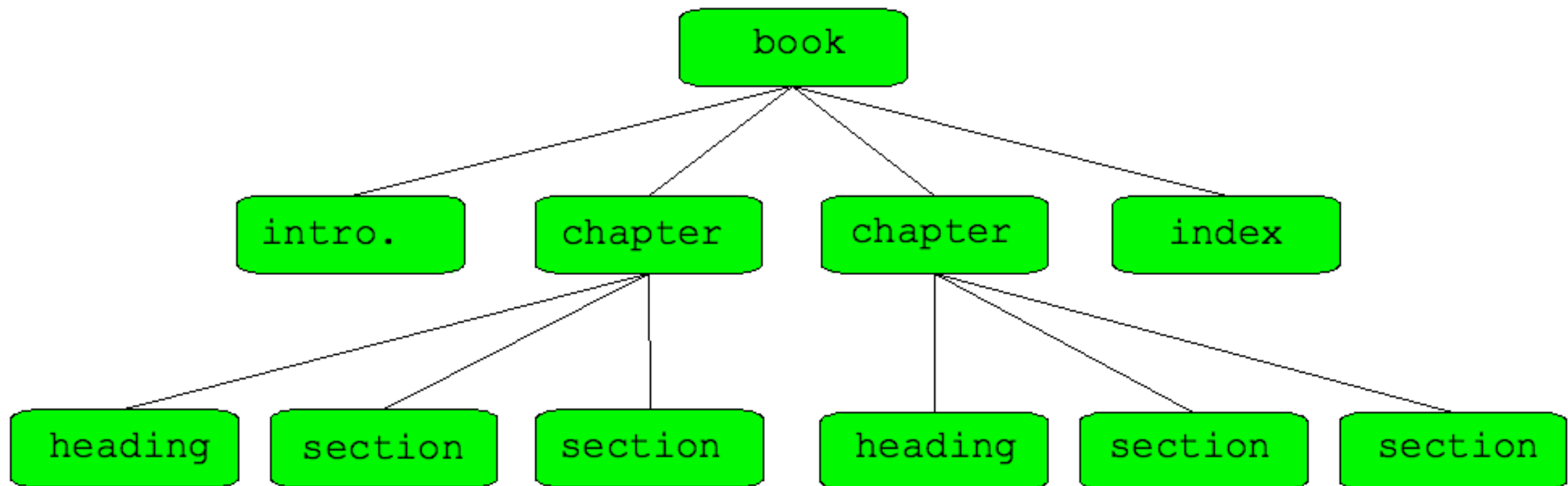
Empty tags or milestone elements: `<lb />` = `<lb></lb>`



# XML Representation: Boxes



# XML Representation: Tree



# XML Representation: Markup

```
<?xml version="1.0" encoding="UTF-8"?>
<book>
  <introduction>Blah blah blah ... </introduction>
  <chapter>
    <heading>Wines</heading>
    <section>White wines ... </section>
    <section>Red wines ... </section>
  </chapter>
  <chapter>
    <heading>Beers</heading>
    <section>Ales ... </section>
    <section>Lagers ... </section>
  </chapter>
  <index> stuff ... </index>
</book>
```

# XML: Well-Formed and Valid

- All XML documents need to be well-formed according to some basic rules:
  - Open and close all tags/elements
  - Tags/elements may not overlap
  - Attribute values must be quoted
- XML documents should be valid according to a DTD or Schema:
  - Use the appropriate elements & attributes
  - Adhere to the “grammar rules” (e.g., allowable attributes for elements)
- Software programs help reinforce these principles
  - XML Editors like Oxygen



# Questions?

---

- John Walsh's XML-Primer:
- <https://www.youtube.com/watch?v=JhhKyyP0e18>

# Intro to the TEI Guidelines and Tag Set

---

- TEI Guidelines: Quick Overview
- TEI P5 Guidelines
- TEI Basic Components
- Basic Markup: Prose
- Basic Markup: Verse
- Basic Markup: Drama
- Basic Markup: Letters

# TEI Guidelines: Quick Overview

- Text Encoding Initiative (TEI) / *Guidelines for Electronic Text Encoding and Interchange (TEI)*
- The TEI *Guidelines* "are addressed to anyone who works with any text in electronic form. They provide means of representing those features of a text which need to be identified explicitly in order to facilitate processing of the text by computer programs" (Sperberg-McQueen).
- TEI provides elements, attributes, and other mechanisms for encoding prose, poetry, drama, dictionaries, critical apparatus, linguistic corpora, and other scholarly and non-scholarly texts.

# TEI Guidelines: Quick Overview

- The TEI Guidelines:
  - Can be applied lightly or heavily
  - Are designed as a set of modules/mechanisms that can be selected as needed:
    - core: Elements common to all TEI documents
    - figures: Tables, formulae, music notation, and figures
    - linking: Linking, segmentation, and alignment
    - msdescription: Manuscript description
    - namesdates: Names and dates
  - Can adapt to local conditions



# TEI P5 Guidelines

- P5 Guidelines:

<http://www.tei-c.org/release/doc/tei-p5-doc/en/html/index.html>

- Prose documentation with examples

- P5 Tag/Element Set:

<http://www.tei-c.org/release/doc/tei-p5-doc/en/html/REF-ELEMENTS.html>

- Listing of the tag set with examples and relevant links to prose documentation

# TEI P5: Basic Components

- **<TEI>**: The root element of a TEI document
  - **<teiHeader>**: The metadata header for a TEI document. Includes bibliographic, technical, administrative, and other metadata about the digital file and the analog source, if one exists.
  - **<text>**: The text itself, e.g., the title page and chapters of a novel, the acts and scenes of a drama, the books or cantos of a long poem. The **<text>** element is further subdivided into:
    - **<front>**: Front matter, e.g, the title page(s), table of contents, potentially a preface or dedication
    - **<body>**: The main body of a document, excluding front and back matter
    - **<back>**: Back matter, e.g., indices, appendices

# TEI P5: Basic Markup: Prose

- **<div>**: (division) is used for basic structural divisions of a text, e.g, volumes, chapters, sections, cantos, tables of contents, indices, appendices, etc. The @type attribute may be used to designate the type of the division.
  - `<div type="chapter">...</div>`
  - `<div type="section">...</div>`
  - `<div type="contents">...</div>`
  - `<div type="canto">...</div>`
- **<head>**: (heading) contains any type of heading, for example the title of a section, or the heading of a list, figure, table, etc.
- **<p>**: (paragraph)
- **<pb>**: (page break) marks the boundary between one page of a text and the next

# TEI P5: Basic Markup: Prose

## Chapter 1: The Manor House

Charles hadn't visited the manor house since Easter, 1955, and now he remembered why. "Hullo", he called out as he walked up the drive, and then, as if to himself, "To be or not to be?, to walk or not to walk...oh, **hang** it all!" His meditation on Hamlet was interrupted as he collided with a peacock. "Sacré bleu!" he exclaimed with irritation, his sang-froid completely deserting him. It was going to be a long week. His catalog of irritations included:

1. The weather
2. The peacocks
3. His meager grasp of French

# TEI P5: Basic Markup: Prose

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <div type="chapter">
3   <head>Chapter 1: The Manor House</head>
4   <p>Charles hadn't visited the manor house since
5     Easter, 1955, and now he remembered why.</p>
6   <p><said>Hullo</said>, he called out as he walked up the
7     drive, and then, as if to himself, <said>To be or
8     not to be?, to walk or not to walk...oh,
9     <emph rendition="#b">hang</emph> it all!</said>
10    His meditation on Hamlet was interrupted as he
11    collided with a peacock. <said xml:lang="fr">Sacré
12    bleu!</said> he exclaimed with irritation, his
13    <foreign xml:lang="fr">sang-froid</foreign> completely deserting him.
14    It was going to be a long week. His catalog of irritations included:
15      <list type="ordered">
16        <item>The weather</item>
17        <item>The peacocks</item>
18        <item>His meager grasp of French</item>
19      </list>
20    </p>
21 </div>
```

# TEI P5: Basic Markup: Verse/Poetry

- `<lg>`: (line group) contains a group of verse lines functioning as a formal unit, e.g. a stanza, refrain, verse paragraph, etc. The `@type` and `@subtype` attributes may be used to classify the type of line group
- `<l>`: (line) contains a line of verse

# TEI P5: Basic Markup: Poetry/Verse

## THE ROUNDEL

A ROUNDEL is wrought as a ring or a starbright sphere,  
With craft of delight and with cunning of sound unsought,  
That the heart of the hearer may smile if to pleasure his ear  
    A roundel is wrought.

Its jewel of music is carven of all or of aught—  
Love, laughter, or mourning—remembrance of rapture or fear—  
That fancy may fashion to hang in the ear of thought.

As a bird's quick song runs round, and the hearts in us hear  
Pause answer to pause, and again the same strain caught,  
So moves the device whence, round as a pearl or tear,  
    A roundel is wrought.



# TEI P5: Basic Markup: Poetry/Verse

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <div type="poem">
3   <head rendition="#center #uppercase">The Roundel</head>
4   <lg>
5     <l><hi rendition="#small-caps">A roundel</hi> is wrought as a ring or a starbright sphere,</l>
6     <l>With craft of delight and with cunning of sound unsought,</l>
7     <l>That the heart of the hearer may smile if to pleasure his ear</l>
8     <l rendition="#l-indent-03">A roundel is wrought.</l>
9   </lg>
10  <lg>
11    <l>Its jewel of music is carven of all or of aught-</l>
12    <l>Love, laughter, or mourning-remembrance of rapture or fear-</l>
13    <l>That fancy may fashion to hang in the ear of thought.</l>
14  </lg>
15  <lg>
16    <l>As a bird's quick song runs round, and the hearts in us hear</l>
17    <l>Pause answer to pause, and again the same strain caught,</l>
18    <l>So moves the device whence, round as a pearl or tear,</l>
19    <l rendition="#l-indent-03">A roundel is wrought.</l>
20  </lg>
21 </div>
```

# TEI P5: Basic Markup: Drama

- **<sp>**: (speech) contains individual speech in a performance text, or a passage presented as such in a prose or verse text.
- **<speaker>**: contains a specialized form of heading or label, giving the name of one or more speakers in a dramatic text or fragment.
- **<stage>**: (stage direction) contains any kind of stage direction within a dramatic text or fragment.

# TEI P5: Basic Markup: Drama

## Scene 1

*Enter Fay*

Fay:

I say, Dinah, has anyone seen my gloves?

*Enter Dinah*

Dinah:

No, miss, perhaps the parakeet has got them again?

*Exit Fay and Dinah*

# TEI P5: Basic Markup: Drama

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <div type="scene">
3     <head rendition="#center">Scene 1</head>
4     <stage rendition="#i">Enter Fay</stage>
5     <sp>
6         <speaker>Fay:</speaker>
7         <p>I say, Dinah, has anyone seen my gloves?</p>
8     </sp>
9     <stage rendition="#i">Enter Dinah</stage>
10    <sp>
11        <speaker>Dinah:</speaker>
12        <p>No, miss, perhaps the parakeet has got them again?</p>
13    </sp>
14    <stage rendition="#i">Exit Fay and Dinah</stage>
15 </div>
```

# TEI P5: Basic Markup: Letters

- **<opener>**: groups together dateline, byline, salutation, and similar phrases appearing as a preliminary group at the start of a division, especially of a letter.
- **<closer>**: groups together dateline, byline, salutation, and similar phrases appearing as a final group at the end of a division, especially of a letter.
  - **<dateline>**: contains a brief description of the place, date, time, etc. of production of a letter, prefixed or suffixed to it as a kind of heading or trailer.
  - **<salute>**: (salutation) contains a salutation or greeting in the closing of a letter, preface, etc.
  - **<signed>**: (signature) contains the closing salutation

# TEI P5: Basic Markup: Letters

---

1906 August the 5<sup>th</sup>

Cape Cod

My dear Becky

How lovely the oysters are this evening!

Yours very truly

Maria





# Hands-on Exercises: Basic Genres

---

- <http://dcl.slis.indiana.edu/teiworkshop/>
- Launch Oxygen
- Complete exercises one at a time: Prose, Verse, Drama and Letters
- Encode a short document of your own choosing