

(4日) (個) (目) (目) 目) のQC





What's in a text? Digital Texts, XML, and TEI Upon Julia's Clothes WHEN as in silks my Julia goes, Then, then (me thinks) how sweetly flowes That liquefaction of her clothes. Next, when I cast mine eyes and see That brave Vibration each way free; O how that glittering taketh me! Upon Julia's Clothes When as in silks my Julia goes, Then, then (me thinks) how sweetly flowes That liquefaction of her clothes. Next, when I cast mine eyes, and see That brave Vibration each way free; O how that glittering taketh mel

ヨト イヨト ヨー のへで





・ コット (雪) (小田) (コット 日)







◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの



What does markup capture?

Digital Texts, XML, and TEI

Compare

<head>Upon Julia's Clothes</head> <lg><l>Whenas in silks my <hi>Julia</hi> goes,</l> <l>Then, then (me thinks) how sweetly flowes</l> <l>That liquefaction of her clothes.</l> </lg>

and

Likewise..

Digital Texts, XML, and TEI

Compare

<hi rend="dropcap">H</hi>&WYN;ET WE GARDE <lb/>na in gear-dagum þeod-cyninga <lb/>prym gefrunon, hu ða æþelingas <lb/>ellen fremedon. oft scyld scefing sceaþe<add>na</add> <lb/>preatum, moneg<expan>um</expan> mægpum meodo-setl<add>a</add>

<lb/>of<damage desc="blot"/>teah egsode <sic>eorl</sic>
syððan ærest wear<add>p</add>

```
<lb/>fea sceaft funden...
```

and

```
<lg>
<lg>
<lp><l>Hwæt! we Gar-dena in gear-dagum</l>
<lp><l>hwæt! we Gar-dena in gear-dagum</l>
<lp><l>peod-cyninga þrym gefrunon, </l>
<lp><l>hu ða æþelingas ellen fremedon, </l>
</lg>
<lg><lg><lp><lsOft Scyld Scefing sceaþena þreatum, </l>
<ls><l>monegum mægþum meodo-setla ofteah; </l>
<ls><l>ester ester ester
```





	Some alphabet soup	
Digital Texts, XML, and TEI	SGML HTML W3C XML DTD CSS Xpath XSLT RelaxNG	Standard Generalized Markup Language Hypertext Markup Language World Wide Web Consortium eXtensible Markup Language Document Type Definition (or Declaration) Cascading Style Sheet XML Path Language eXtensible Stylesheet Language - Transformations Regular Expression Language for XML (New Generation)

Oh, and then there's also TEI, the *Text Encoding Initiative*

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ○ □ ● ○ ○ ○ ○



An example XML document <?xml version="1.0" encoding="utf-8" ?> Digital Texts, XML, and TEL <cookBook> <recipe n="1"> <head>Nail Soup</head> <ingredientList> <ingredient>an onion</ingredient> <ingredient>two carrots</ingredient> <ingredient>water</ingredient> <ingredient>a nail</ingredient> <ingredient>some gullible peasants</ingredient> </ingredientList> <procedure> <step>put the water on to boil</step> <step>take out the nail and serve</step> </procedure> </recipe> <recipe n="2"> <!-- contents of second recipe here --> </recipe> <!-- hic desunt multa --> </cookBook>







Representing an XML tree

Digital Texts, XML, and TEI

- An XML document is encoded as a linear string of characters
- It begins with a special processing instruction
- Element occurrences are marked by start- and end-tags
- The characters < and & are Magic and must always be "escaped"
- Comments are delimited by <!- and ->
- CDATA sections are delimited by <![CDATA[and]]>
- Attribute name/value pairs are supplied on the start-tag and may be given in any order
- Entity references are delimited by & and ;



	Splot the mistake	
Digital Texts, XML, and TEI		
	<pre><greeting>Hello world!</greeting> <greeting>Hello world!</greeting> <greeting><grunt>Ho</grunt> world!</greeting> <grunt>Ho <greeting>world!</greeting></grunt> <greeting><grunt>Ho world!</grunt></greeting> <grunt type="loud">Ho</grunt> <grunt type="loud"> <grunt type="loud"> <grunt type="loud"> </grunt> </grunt> </grunt> </pre>	

▲□▶▲圖▶▲≣▶▲≣▶ ■ のへの

Defining the rules

Digital Texts, XML, and TEI

A valid XML document conforms to rules which are stated in an external schema of some sort. A schema specifies:

- the name of the root element
- names for all elements used
- names and datatypes and (occasionally) default values for their attributes
- rules about how elements can nest
- and a few other things, depending on the schema language

n.b. A schema does *not* specify anything about what elements "mean"





The root element of the document itself

▲□▶▲□▶▲□▶▲□▶ □ のQ@



Namespace declarations

Digital Texts, XML, and TEI

- An XML document can use elements declared in different name spaces.
 - a namespace declaration associates a namespace prefix with an external identifier (which looks like an URL)
 - the default namespace *may* be declared using a special xmlns attribute
 - other name spaces must all use a special prefix, which is also declared

```
<TEI xmlns="http://www.tei-c.org/ns/1.0"> ... </TEI>
```

There is a special xml namespace, used by the TEI for global attributes xml:id and xml:lang

The Doctype Declaration

Digital Texts, XML, and TEI

In DTD world, an optional "Document Type" declaration may appear:

```
<?xml version="1.0" ?>
<!DOCTYPE hello [<!ELEMENT hello (#PCDATA)>]>
<hello xmlns="http://www.greetings.org">
hello world
</hello>
```

- The DTD is one way of associating the document with its schema (but is not used by W3C or Relax NG for this purpose)
- The DTD subset is used to provide declarations additional to those in the schema
- The DTD subset may be internal, external, or both

















◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの