

Sistemi per il recupero delle informazioni

Gabriele Pozzani

A.A. 2015/2016

**Corso di Laurea Magistrale in
Editoria e Giornalismo**

XML per l'editoria elettronica

XML come formato documentale

- Utilizzare XML per:
 - Pagine Web
 - Libri
 - Articoli
 - Manuali di riferimento/istruzione
 - Corsi
 - Libri di testo

3

XML come formato di salvataggio

- I principali formati di salvataggio sono basati su XML:
 - Office Open XML (OOXML, OpenXML): docx, pptx, xlsx
 - Open Document Format (ODF): odt, odp, ods
 - EPUB

4

Esempio, XML di un EPUB



The screenshot shows a software interface with multiple tabs and panes. The main pane displays an XML document structure for an EPUB file. The XML code is as follows:

```
<?xml version="1.0"?><!DOCTYPE ncx PUBLIC "-//NISO//DTD ncx 2005-1-1//EN" "http://www.daisy.org/ns/specifications/niso-ncx-2005-1-1.dtd">
<head>
    <meta name="dtb:uid" content="444cb7d-76b0-4840ab3d7726" />
    <meta name="dtb:depth" content="1" />
    <meta name="dtb:totalPageCount" content="0" />
    <meta name="dtb:maxPageNumber" content="0" />
</head>
<ncx version="1.0" xmlns="http://www.daisy.org/ns/specifications/niso-ncx-2005-1-1">
    <head>
        <meta name="dtb:depth" content="1" />
        <meta name="dtb:totalPageCount" content="0" />
        <meta name="dtb:maxPageNumber" content="0" />
    </head>
    <text>
        <navPoint id="navpoint-1" playOrder="1">
            <navLabel>
                <text>copyright</text>
            </navLabel>
            <content src="iPad572cbat7.xhtml" />
        </navPoint>
        <navPoint id="navpoint-2" playOrder="2">
            <navLabel>
                <text>About</text>
            </navLabel>
            <content src="iPad572cbat8.xhtml" />
        </navPoint>
    </text>
    <listItem id="list-item-1" playOrder="1" type="list-item">
        <navPoint id="navpoint-3" playOrder="1">
            <navLabel>
                <text>Table of Contents</text>
            </navLabel>
            <content src="iPad572cbat5.xhtml" />
        </navPoint>
    </listItem>
    <listItem id="list-item-2" playOrder="2" type="list-item">
        <navPoint id="navpoint-4" playOrder="1">
            <navLabel>
                <text>About</text>
            </navLabel>
            <content src="iPad572cbat8.xhtml" />
        </navPoint>
    </listItem>
</ncx>
```

5

XML per diversi utenti

- XML soddisfa le esigenze di:
 - Programmatori (o utenti esperti) addestrati a lavorare con le strutture rigide tipiche delle applicazioni orientate ai dati.
 - Scrittori ed editori (o utenti normali) che preferiscono la forma libera di un libro o di un articolo.
 - XML soddisfa le esigenze di entrambe le comunità in maniera equa e soddisfacente.

6

Usi dei tag XML

- L'arricchimento di un documento con una struttura tramite l'uso di XML porta a:
 - Rappresentazione sia del contenuto che di metadati
 - Strutturazione del contenuto tramite più suddivisioni
 - Tra loro indipendenti
 - Con semantiche diverse
 - Con i propri metadati
 - Arricchimento semantico delle parti “interessanti” del contenuto tramite appositi marcatori che ne definiscono il significato del contenuto

17

TEI (Text Encoding Initiative)

<http://www.tei-c.org/>

- Applicazione per la marcatura della letteratura classica
- Primo esempio di DTD orientato a documenti di tipo narrativo
- Include elementi per:
 - Le strutture letterarie più comuni (capitoli, scena, stanza, ...)
 - La tipografia
 - Le strutture grammaticali

18

TEI example: Plato, Parmenides (I)

From <http://www.perseus.tufts.edu/>

```
...  
  <revisionDesc>  
    <change>  
      <date>July, 1992</date>  
      <respStmt>  
        <name>WPM</name>  
        <resp>(n/a)</resp>  
      </respStmt>  
      <item>  
        Tagged in conformance with Prose.e dtd.  
      </item>  
    </change>  
  </revisionDesc>  
...  
...
```

19

TEI example: Plato, Parmenides (II)

```
...  
  <text><group>  
    <text n="Parm.">  
    <body>  
      <head>Parmenides</head>  
      <castList>  
        <castItem type="role">  
          <role>Cephalus</role>  
        </castItem>  
        <castItem type="role">  
          <role>Antiphon</role>  
        </castItem>  
        <castItem type="role">  
          <role>Aristoteles</role>  
        </castItem>  
      </castList>  
...
```

20

TEI example (2): la rivista Scandinavian-Canadian Studies

Dall'XML vengono ottenuti in automatico la corrispondente pagina Web e il PDF

The screenshot shows a web page with a grey header containing the text "SCANDINAVIAN-CANADIAN STUDIES/ÉTUDES SCANDINAVES AU CANADA" and "Vol.16 (2006) pp.443-444". Below the header is a large XML code block. To the right of the XML, there is a summary of the book "Ingmar Bergman: A Reference Guide" by Birgitta Steene, published by Brian McIlroy. The summary includes details like the publisher (Amsterdam University Press), date (2005), and ISBN (9053564063). It also mentions that the book is nearly four times the size of the earlier volume and is a strong case for the sheer massiveness of the compilation.

```
<TEI.2 id="mcilroy-1.16">
  <teiHeader>...</teiHeader>
  <text>
    <front>
      <docTitle n="Ingmar Bergman: A Reference Guide">
        <titlePart type="Main">
          <name reg="Steene, Birgitta">Birgitta Steene</name>
          <title level="m">Ingmar Bergman: A Reference Guide</title>
        </titlePart>
        <titlePart type="ReviewedBook">
          <listBibl>
            <biblStruct>
              <monogr>
                <author>
                  <name reg="Steene, Birgitta">Birgitta Steene</name>
                </author>
                <author>
                  <name reg="McIlroy, Brian">Brian McIlroy</name>
                </author>
                <imprint>
                  <publisher>Amsterdam University Press</publisher>
                  <pubPlace>Amsterdam</pubPlace>
                  <date value="2005">2005</date>
                  <biblScope type="pages">1152 pages</biblScope>
                </imprint>
                <note>ISBN 9053564063. (hbk) €62.5.-</note>
              </monogr>
            </listBibl>
            <titlePart>
              <docAuthor>
                <name key="mcilroy_brian" reg="McIlroy, Brian">Brian McIlroy</name>
                <note>teaches film studies at the University of British Columbia. He is title level="m">World Cinema 2: Sweden</title>
                <note>(London: Flicks Books, 1986).</note>
              </docAuthor>
              <titlePart type="short_affil">Dept. of Film, University of British
            </front>
            <div0>
              <p>
                Nearly twenty years ago, Birgitta Steene published a reference guide to Ingmar Bergman's major film artist of G. K. Hall, and it amounted to three hundred pages. It was a series of filmographical entries for current and future Bergman film and theatre scholars. Birgitta Steene, she makes a strong case by the sheer massiveness of this compilation that Ingmar Bergman is not just Sweden's major film artist of the twentieth century, but possibly Europe's.
              </p>
            </div0>
          </listBibl>
        </titlePart>
      </docTitle>
      <ref>View metadata</ref>
    </front>
    <back>
      <ref>View metadata</ref>
    </back>
  </text>

```

DocBook

<http://www.oasis-open.org/docbook/>

- Applicazione per documenti nuovi (non vecchi)
- Formato per la creazione del testo (non di un prodotto finito e pronto per essere presentato al pubblico)
- Utilizzato nella documentazione relativa al campo informatico
- Sintassi molto semplice
- Modulari e quindi utilizzabili solo in parte (porzioni e strutture che servono)

DocBook example (I)

```
<!DOCTYPE article PUBLIC "-//OASIS//DTD DocBook V4.1//EN">
<article>
  <articleinfo>
    <title>An Example Article</title>
    <author>
      <firstname>Your first name</firstname>
      <surname>Your surname</surname>
      <affiliation>
        <address>
          <email>foo@example.com</email>
        </address>
      </affiliation>
    </author>
    <copyright>
      <year>2000</year>
      <holder>Copyright string here</holder>
    </copyright>
```

23

DocBook example (II)

```
<abstract>
  <para>If your article has an abstract
    then it should go here.</para>
</abstract>
</articleinfo>
<sect1>
  <title>My First Section</title>
  <para>This is the first section
    in my article.</para>
<sect2>
  <title>My First Sub-Section</title>
  <para>This is the first sub-section
    in my article.</para>
</sect2>
</sect1>
</article>
```

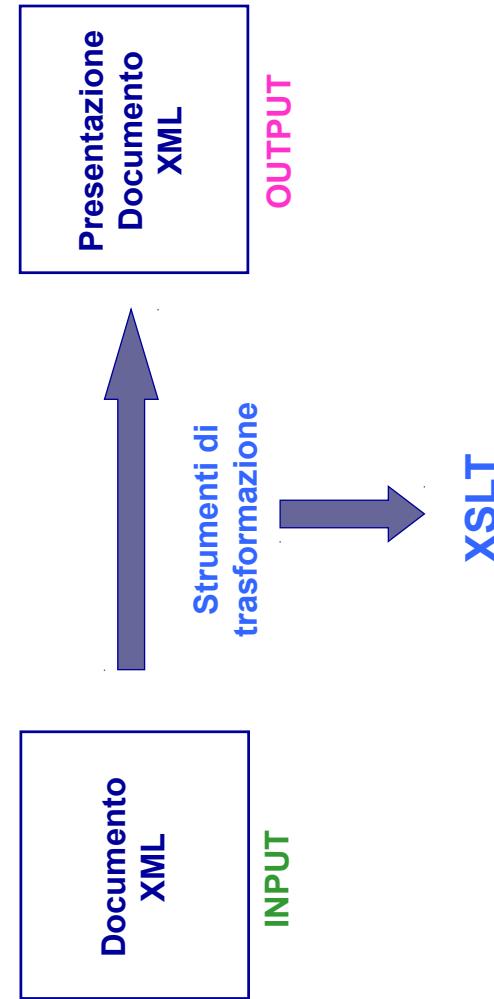
24

Transformazione e presentazione (1)

- La marcatura di un documento XML descrive la struttura del documento e non la sua presentazione.
 - Specifica l'organizzazione.
 - Non specifica come deve apparire.
- Un documento XML può essere letto nel suo formato nativo (marcatori + testo).
- Normalmente viene tradotto in un formato differente adatto alla presentazione.
- Per XML il formato di INPUT non deve necessariamente corrispondere al formato di OUTPUT.
 - Il formato di INPUT serve per agevolare chi scrive.
 - Il formato di OUTPUT serve per agevolare chi legge.

25

Transformazione e presentazione (2)



26

XSLT (eXtensible Stylesheet Language Transformation)

- L'XLST descrive il processo secondo cui un documento XML viene trasformato per produrre il documento finale, pronto per la presentazione
- Un documento XSL contiene un elenco di modelli
 - Definisce dei fogli di stile
- Ogni stile verrà applicato a determinati elementi in modo da renderli “presentabili”

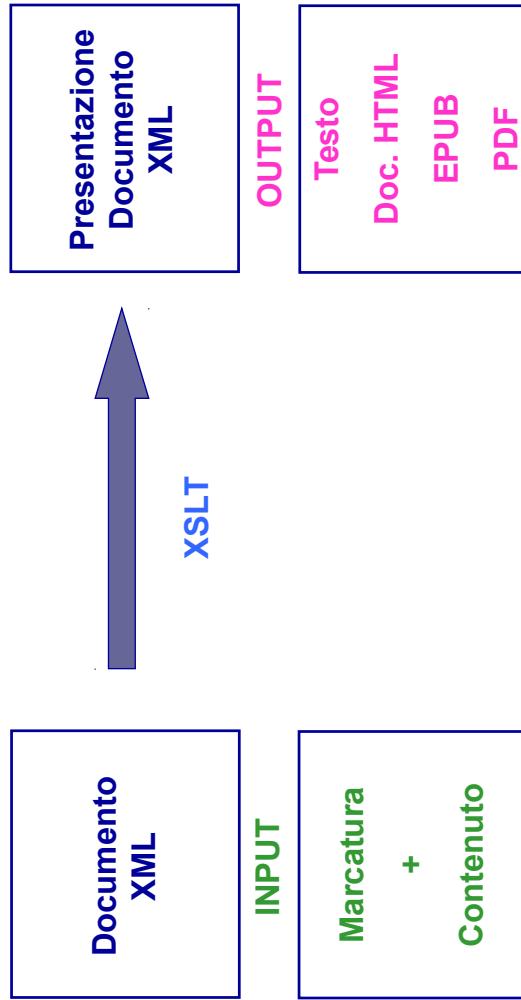
27

XSLT: funzionamento (1)

- Legge documento in INPUT
- Quando incontra qualcosa (del documento in INPUT) che corrisponde a un modello del foglio di stile
 - Produce in OUTPUT il modello con il relativo contenuto
- Il modello rappresenta il modo in cui il contenuto deve venir presentato

28

XSLT: funzionamento (2)



29

XML per l'editoria

- Come si situano tutte queste tecnologie nell'ambito dell'editoria?
- Perché gli editori dovrebbero usare/aver bisogno di XML?
- Vediamo alcune considerazioni
 - effettuate sia da informatici che da editori che possano aiutarci a rispondere a queste domande

31

XML per l'editoria elettronica

- Come si può facilmente comprendere l'XML, in quanto tecnologia informatica, mostra i suoi maggiori vantaggi nell'ambito dell'editoria elettronica
 - L'editoria elettronica è ancora in uno stadio di comprensione e sviluppo
 - Gli editori non hanno ancora ben compreso come sfruttare le nuove possibilità
 - Quale modello di business?
 - Quali sono le nuove funzionalità?
 - L'uso dell'XML è ancor più limitato

32

Utilizzi di XML per l'editoria

- La formattazione di documenti/libri/articoli in XML apre due principali possibilità:
 - 1) Arricchimento del testo con metadati e annotazioni semantiche
 - Permettono una maggiore efficacia e ricchezza nella ricerca di informazioni
 - SRI per dati semistruturati
 - 2) Processo di lavorazione “singolo input, output multipli”
 - Flusso di lavoro XML

33

Singolo input, output multipli (I)

- XML è un formato
 - Neutrale
 - Intermediario
 - Di interscambio
- Esso può essere usato come punto di partenza per un processo di produzione editoriale basato sul riuso
 - Uno stesso doc XML può essere usato per produrre, automaticamente, diversi formati/tipi di pubblicazione

34

Singolo input, output multipli (II)

- Uno stesso XML può essere usato per produrre PDF, vari ebook, pagine web, ...
- Con o senza DRM

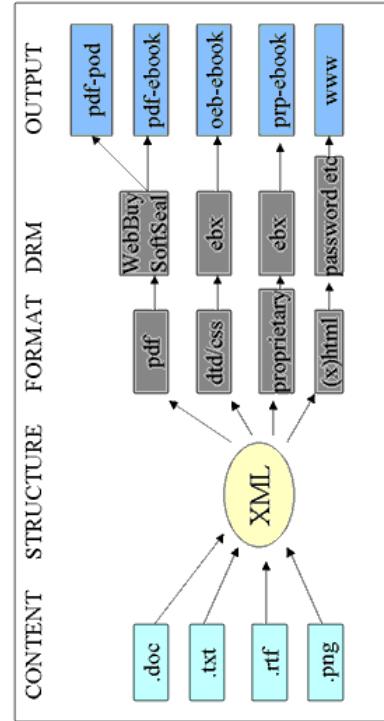


Figure: From raw content and XML to distributed publications

- Per le vecchie produzioni, l'XML può essere ottenuto dai classici formati doc, txt, ...
- Per le nuove produzioni, l'XML può essere prodotto direttamente per poi ottenere, se necessario anche documenti in formato doc, ...

35

XML Workflow

- Più in generale l'XML è alla base dell'XML workflow, o flusso di lavoro XML
 - In contrapposizione alla lavorazione “tradizionale”

Flusso di lavoro tradizionale

- Con la lavorazione tradizionale
 - i contenuti vengono importanti nel programma di impaginazione
 - l'impaginatore deve applicare manualmente uno ad uno gli stili appropriati a ciascun elemento del testo (titoli di capitoli, titoli di sezioni, sottotitoli, ecc...)
- Svantaggi:
 - richiede un certo tempo ed è noioso e ripetitivo
 - è possibile commettere errori
 - Non c'è alcun modo di riutilizzare questo lavoro per un libro che abbia delle caratteristiche simili

Flusso di lavoro XML (I)

- Il flusso di lavoro viene separato in due processi paralleli: il lavoro sui contenuti e il lavoro sulla presentazione
 - Da un lato l'editor organizza i contenuti in modo tale da consentire al wordprocessor di applicare dei tag XML a ciascun blocco di testo
 - Un tag XML, in questo caso, non è altro che un'etichetta che descriva il tipo di blocco di testo
 - Può essere fatto già con MS Office e LibreOffice

Flusso di lavoro XML (II)

- L'impaginatore preparerà la presentazione, cioè crea gli stili appropriati per ciascun tipo di elemento strutturale del testo (titoli di capitolo, sottotitoli, box, immagini, ...)
 - Usando XSL
 - Per ciascun “formato” di libro che si vorrà ottenere

Flusso di lavoro XML (III)

- Quando l'editor e l'impaginatore hanno concluso il loro lavoro
 - Si importa l'XML dei contenuti nel programma di impaginazione
 - Si chiede al programma di impaginazione di applicare i corretti stili a ciascun blocco di testo

Flusso di lavoro XML (IV)

- Vantaggi:
 - L'applicazione degli stili è automatica, dunque si risparmia tempo
 - Non possono esservi errori nell'applicazione di stili.
 - Il template di stili creato dall'impaginatore è riutilizzabile per qualsiasi altro libro che abbia delle caratteristiche analoghe.
 - Dunque è possibile ad esempio produrre il secondo volume di una collana di ricette lavorando solo alla strutturazione dei contenuti e passare direttamente all'importazione dei contenuti nel programma di impaginazione.
 - L'impaginazione ottenuta con questo processo è molto più idonea alla generazione di ebook

Flusso di lavoro XML (V)

- I classici wordprocessor (MS Office e LibreOffice) sfruttano già l'XML
 - I file .docx, .xlsx, .odt, ... non sono altro che uno ZIP di file XML
 - Entrambi permettono di esportare in XML
 - LibreOffice permette di esportare in formato XML valido per DocBook
 - Vi sono poi editor appositi per XML
 - I principali sw di impaginazione basati su XML permettono di importare direttamente i file docx, odt, ...

Flusso di lavoro XML, esempio (I)

1) Preparo i contenuti in LibreOffice

- Applicando gli stili che diventeranno poi tag XML
 - Non guardo a come esse appaiono, quello sarà compito dell'XSL

NoSQL Database

Introduction

NoSQL (commonly interpreted as "not only SQL") is a broad class of database management systems identified by non-adherence to the widely used relational database management system model. NoSQL databases are not built primarily on tables, and generally do not use SQL for data manipulation.

NoSQL database systems are often highly optimized for retrieval and appending operations and flexibility compared to full SQL systems. The reduced run-time performance for certain data models. The data can be structured, but NoSQL is used when what really matters is the ability to store and retrieve great quantities of data, not the relationships between the elements.

This organization is particularly useful for statistical or real-time analyses of growing lists of elements. NoSQL does not use SQL as its query language. NoSQL database systems arose alongside major Internet companies, such as Google, Amazon, and Facebook, which had challenges in dealing with huge quantities of data which conventional RDBMS solutions could not cope with. Note that both

Flusso di lavoro XML, esempio (III)

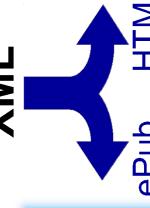
2) Apro l'ODT con l'impaginatore (e.g., <oXygen/>)

- Questo dà accesso direttamente all'xml contenuto nell'odt



Flusso di lavoro XML, esempio (III)

- L'impaginatore contiene già alcuni XSLT di base che possono essere applicati e/o modificati



A screenshot of a web browser displaying the transformed HTML output. The page title is "NoSQL Database - Mozilla Firefox". The content includes a table of contents for a "NoSQL Database" document, with sections like "1. NoSQL Database", "1.1. Introduction", "1.2. History", etc. The browser's address bar shows the URL "file:///media/david/Documenti/Univ/NoSQL/Word/TEP/test.html".

A screenshot of a web browser displaying the transformed PDF output. The page title is "NoSQL Database - Mozilla Firefox". The content is identical to the HTML version, showing the table of contents for the "NoSQL Database" document. The browser's address bar shows the URL "file:///media/david/Documenti/Univ/NoSQL/Word/TEP/test.html".

1. NoSQL Database

1.1. Introduction

NoSQL, commonly interpreted as "not only SQL", is a broad class of database management systems identified by non-adherence to the widely used relational database management system model. NoSQL databases are primarily built on tables, and generally do not use SQL for data manipulation.

NoSQL database systems are often highly optimized for retrieval and updating operations and often offer little functionality beyond record storage (e.g., key-value stores). The performance and scalability of NoSQL databases is compensated by marked gains in processing power for certain data models.

The data can be structured, but NoSQL is used when what really matters is the ability to store and retrieve great quantities of data, not the relationships between the elements.

Obiezioni

- Diversi editori obiettano che
 - il riuso non è così importante
 - I maggiori costi (economici, temporali, ...) nell'uso dell'XML non sono quindi giustificati
- Queste obiezioni nascono da considerazioni principalmente editoriali “di vecchio stampo”
- Ovviamente non sempre è necessario il riuso
 - E con esso le tecnologie che lo supportano
 - Dipende dal contesto lavorativo e dal contenuto
- Il fatto che le tecnologie esistano, non significa che vadano sempre usate, ma innanzitutto vanno comprese

46

Ribattere alle obiezioni

- Certamente l'e-publishing è un ambito in cui il riuso è fondamentale
 - E le tecnologie che lo supportano possono
 - Semplificare
 - Velocizzare
 - Automatizzare
- Il processo lavorativo
 - Il passaggio ad un processo di produzione editoriale basato sull'XML richiede
 - La formazione di editori ed autori
 - La comprensione della separazione tra contenuto e sua presentazione
 - La comprensione che è il costrutto semantico ad essere importante, non la presentazione finale
- Con lo sviluppo dell'e-publishing ci si può aspettare che riuso e XML possano crescere di importanza e in necessità

47

<Oxygen/>

http://en.wikipedia.org/wiki/Comparison_of_XML_editors

- Editor XML
 - Permette l'editing di documenti XML sia in formato testuale che “grafico” (WYSIWYM)
 - Permette la validazione del documento rispetto al suo DTD
 - Supporto particolare all'editing di documenti TEI, DocBook, e altri principali formati XML
 - Inclusi i formati MS Office, LibreOffice e EPUB
 - Trasformazione dei doc XML tramite XSLT
 - Include già alcuni XSLT per la trasformazione in EPUB, PDF, (X)HTML
 - Creazione ed editing degli stili XSLT