

XML: introduzione alla codifica dei testi

Con la codifica dei testi si intende la rappresentazione dei testi stessi su un supporto digitale in un formato utilizzabile dall'elaboratore (**M**achine **R**eadable **F**orm) mediante un opportuno linguaggio formale.

Tale linguaggio descrive il testo in ogni sua parte attraverso delle **marche** o **tag**, ovvero stringhe di caratteri definite dalle due parentesi uncinate < >, al cui interno sono specificati dei comandi.

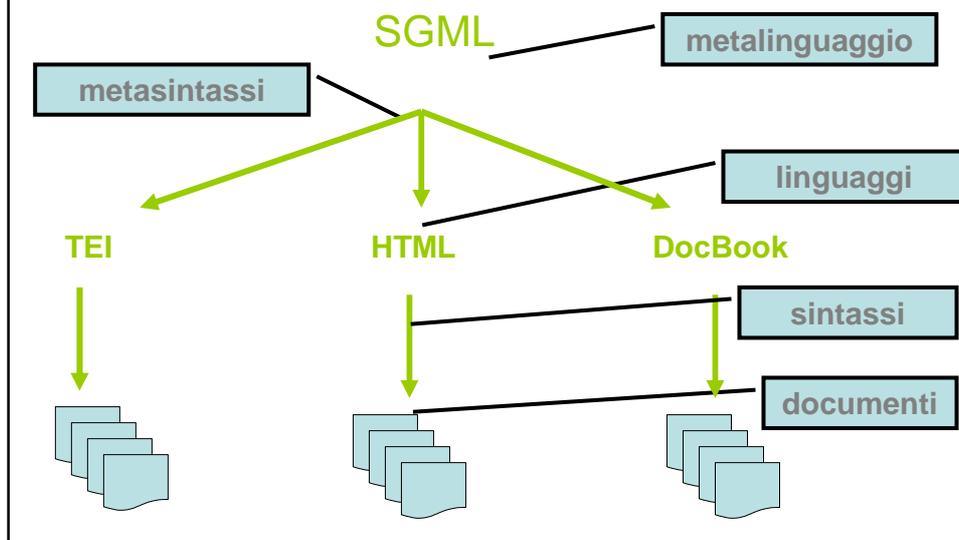
La nascita del linguaggio

Il progetto XML ha avuto inizio alla fine del 1996, nell'ambito della *SGML Activity* del W3C. L'interesse che questo progetto, il più importante dell'organizzazione dopo il World Wide Web, ha suscitato, ha condotto il W3C a creare un gruppo di lavoro, *XML Working Group*, composto da esperti mondiali delle tecnologie SGML, ed una commissione, *XML Editorial Review Board*, deputata alla redazione delle specifiche del progetto.

Nel Febbraio del **1998**, le specifiche sono divenute una raccomandazione ufficiale, con il nome *Extensible Mark-up Language* (XML) versione 1.0.

- L'**XML** non è soltanto uno strumento per il WEB ma è qualcosa di più: è uno strumento che permetterà la condivisione totale dei dati, intesi come pura informazione, a prescindere dal tipo di visualizzazione e di utilizzo che se ne farà in futuro.
- XML è un acronimo per *eXtensible Markup Language*: *Language* perché si tratta di un linguaggio, *Markup* perché è fondato sull'utilizzo dei marcatori ed *eXtensible* perché consente a chi lo utilizza di creare i *tag* di cui ha bisogno.
- **Metalinguaggio**. Non set di marcatori predefinito come HTML (cfr. DTD) ma istruzioni di SINTASSI per la creazione di molteplici linguaggi di codifica.

Il concetto di metalinguaggio di codifica



Le raccomandazioni del W3C sull'XML sono:

- XML dovrà supportare un largo campo di applicazioni (motori per la visualizzazione di contenuti, strumenti di traduzione e applicazioni di database);
- XML dovrà essere compatibile con SGML
- XML dovrà essere facilmente interpretabile in modo da facilitarne la diffusione;
- XML non dovrà avere opzioni perché possono dare problemi di compatibilità;
- XML dovrà essere leggibile dall'uomo anche se questi non ha un Parser XML ma un semplice editor;
- XML dovrà avere una progettazione formale e concisa non come SGML;
- I documenti XML dovranno essere facili da creare anche con un semplice editor.

Concetti di base

- **Il rapporto struttura logica e layout:** XML codifica la STRUTTURA del documento ma non si preoccupa di come apparirà sullo schermo.
- **Documenti VALIDI e BEN FORMATI.** XML consente di codificare documenti usando marcatori a discrezione del codificatore, scelti sulla base della STRUTTURA del documento (BEN FORMATO - well formed) oppure di utilizzare una DTD (esistente o creata ad hoc per il tipo di documento - VALIDO).
- **Lavorare con l'XML:** il documento, la DTD (*Document Type Definition* - opzionale), il foglio di stile (XSL - *eXtensible Style Sheet Language*).

Il Markup: Elementi, attributi, entità

- Markup come sistema di descrizione dei documenti; usa la convenzione dei delimitatori (apertura <nome>; chiusura</nome>) per racchiudere l'informazione sul testo.
- **ELEMENTO**: caratteristiche della partizione logica della sezione di testo
(es.<nomepersona>Francesca</nomepersona>)
- **ATTRIBUTO**: caratteristiche specifiche dell'elemento
(es.<nomepersona tipo="f">Francesca</nomepersona>)
- **ENTITA'**: riferimento ad oggetti "esterni" al documento. Convenzione &nomeentità;

La sintassi dell'XML

- No elementi vuoti , nuova sintassi: oppure
- Case sensitive: ma non
- Valore dell'attributo fra virgolette
- Corretta nidificazione: NO <i></i> MA <i></i>
- i simboli "< >" vanno usati solo per includere i comandi dei *tag*, il simbolo & deve essere usato come riferimento per le entità (á), dove per entità si intende una sequenza arbitraria di byte a cui viene dato un nome mediante la dichiarazione della DTD.

Editor, parser e browser

- Scrivere un documento XML: editor di testo (come Blocco Note di Windows) e software visuali (come Xmetal, TextPad, XML Pro, etc.)
- Analizzare la correttezza del documento: parser validanti (verificano che il documento segua le regole sintattiche del linguaggio e le norme specificate nella DTD) e non validanti (verificano solo l'adeguatezza sintattica del documento)
- La distribuzione sul Web: Internet Explorer 5 (con parser integrato)

Un semplice documento XML

Essendo come l'HTML un formato "solo testo" è possibile realizzare il documento partendo da un editor di testo qualsiasi (come "Blocco Note" di Windows).

```
<catalogo>
  <libro numero="1">
    <autore>Cesare Pavese</autore>
    <titolo>La casa in collina</titolo>
  </libro>
  <br/>
  <libro numero="2">
    <autore>Francesco Petrarca</autore>
    <titolo>Il Canzoniere</titolo>
  </libro>
</catalogo>
```

Il Prologo e l'istanza

- **Prologo** - due parti relative a
 - 1. versione del documento (unica dichiarazione obbligatoria), riferimento alla presenza di una DTD, set di caratteri utilizzato (UTF-8 o ISO-8859-1)
 - 2. eventuale DTD cui fare riferimento:

Dichiarazione XML (case-sensitive!)

```
<?xml version="1.0" standalone="yes" encoding="ISO-8859-1"?>
```

Dichiarazione del tipo di documento (opzionale)

```
<!DOCTYPE nomeelementoradice SYSTEM "nomedtd.dtd">
```

- **Istanza del documento**: Elemento **radice** e serie dei tag in struttura gerarchica.

La DTD

Document Type Definition

La DTD raccoglie l'elenco dei marcatori utilizzabili in fase di codifica, gli attributi eventuali e definisce le relazioni fra gli elementi.

Esistono DTD già fatte ma è possibile anche crearne di nuove.

Usare una DTD di riferimento in XML non è obbligatorio.

Se il documento XML fa riferimento ad una DTD si dice

VALIDO; se il set di marcatori è d'invenzione il documento è solo **BEN FORMATO** cioè rispetta le regole sintattiche del linguaggio.

DTD – dichiarazione dell'elemento

Definisce il tipo di *elemento* ed il *content model* (indica quali elementi possono occorrere all'interno dell'elemento stesso, in che ordine e con quale ricorrenza)

<!ELEMENT testo (intro, corpo, app)>

Content model dell'elemento

In questo caso l'elemento definito è *testo* ed all'interno delle parentesi tonde è espresso il *content model*. Quest'ultimo ha una sintassi molto complessa, composta da **indicatori di occorrenza e connettivi**.

Il foglio di stile

XML codifica la STRUTTURA del documento e delega ad altri linguaggi il compito di definire quale sarà invece il layout, cioè l'aspetto, la formattazione del documento, come cioè il testo codificato apparirà sullo schermo in fase di visualizzazione.

Il foglio di stile contiene dunque le istruzioni di formattazione: ogni porzione di testo che sta tra due marcatori assumerà l'aspetto definito per quel marcatore nel foglio di stile

CSS e XSL

XML per la formattazione usa i CSS (*Cascading Style Sheet*). Ma è anche nato un nuovo linguaggio ad hoc per l'XML che si chiama XSL (eXtensible Mark up Language). Si basa sulle regole di formattazione dell'HTML e su quelle dei CSS.

Prevede l'assegnazione di caratteristiche fisiche per ciascuno dei marcatori utilizzati nel documento XML e quindi definisce come ogni porzione di testo apparirà in fase di visualizzazione.

La visualizzazione del documento - I fogli di stile

Anche il foglio di stile è realizzabile ricorrendo ad un semplice editor di testo. Il file avrà estensione .xsl.

Intestazione:

```
<?xml version="1.0" ?>  
<xsl:stylesheet version="1.0"  
  xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
```

L'istruzione iniziale per applicare tutta la formattazione all'intero documento:

```
<xsl:template match="/">  
  <xsl:apply-templates select="nomedellaroot"/>  
</xsl:template>
```

La formattazione degli elementi con l'HTML

```
<xsl:template match="autore">  
  <center>  
    <font color="red" size="+3" face="Arial">  
      <xsl:apply-templates/>  
    </font>  
  </center>  
</xsl:template>
```

[Per le istruzioni HTML fare riferimento alle slide dedicate]

Inserire immagini e creare link

```
<xsl:template match="elemento">
  
  <a href="link.html">
    <xsl:apply-templates/>
  </a>
</xsl:template>
```

Linkare un'immagine

File XML:

```
<figure name="immagine"/>
```

File XSL:

```
<xsl:template match="figure[@name='immagine']">
  <a href="link-immagine.estensione" target="blank">
    
  <xsl:apply-templates />
</a>
</xsl:template>
```

Altre regole per i fogli di stile

```
<xsl:template match="divisione/paragrafo">
  <xsl:apply-templates/>
</xsl:template>
```

Se ho più elementi "paragrafo" che dipendono però da elementi di livello superiore diverso, posso specificare su quale degli elementi "paragrafo" voglio applicare la formattazione.

Altre regole per i fogli di stile

```
<xsl:template match="elemento[@attributo='valore']">
  <xsl:apply-templates/>
</xsl:template>
```

ESEMPIO:

```
<xsl:template match="paragrafo[@numero='1']">
  <xsl:apply-templates/>
</xsl:template>
```

In questo modo posso applicare la formattazione solo a quegli elementi che hanno un attributo con il valore specificato. Senza specificare il valore posso farlo con tutti gli elementi che hanno quell'attributo, indipendentemente dal valore.

```
<xsl:template match="elemento[@attributo]">
  <xsl:apply-templates/>
</xsl:template>
```

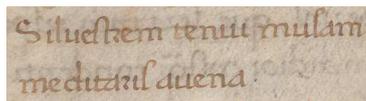
Altre regole per i fogli di stile

```
<xsl:template match="paragrafo">
  <xsl:value-of select="@attributo">
  <xsl:apply-templates/>
</xsl:template>
```

In questo modo posso specificare che voglio visualizzare il valore dell'attributo (associato all'elemento indicato) in fase di layout.

Esempio

siluestrem tenui <sic corr="M">m</sic>usam meditaris a<sic corr="v">u</sic>ena;



Layout con <sic>	<i>S</i> siluestrem tenui <i>musam</i> meditaris aena; /
<pre><xsl:template match="sic"> <i><xsl:apply-templates/></i> </xsl:template></pre>	
Layout con <corr>	siluestrem tenui / <i>usam</i> meditaris aena;
<pre><xsl:template match="sic"> <i><xsl:value-of select="@corr"/></i> </xsl:template></pre>	

Il documento e il foglio di stile

[Dichiarazione XML:]

```
<?xml version="1.0" encoding="ISO-8859-1"
standalone="no"?>
```

[Collegamento al foglio di stile:]

```
<?xml-stylesheet href="nomedelfogliodistile.xml" type="text/xsl"?>
```

<elemento radice>

[testo con marcatori]

</elemento radice>

Altre caratteristiche del documento XML - Entità predefinite

Il segno 'et' (&), le virgolette doppie e singole (' e ") e i segni dei tag (< e >) non possono apparire nel testo del documento. La ragione è che questi cinque caratteri sono riservati per le istruzioni di processo di XML. Volendo utilizzare questi caratteri devono essere sostituiti con i cosiddetti riferimenti di entità (*entity references*) e in questo modo non vengono interpretati come parti del markup. Un riferimento di entità è la combinazione di vari caratteri scritti fra un & e un punto e virgola.

I cinque riferimenti di entità predefiniti in XML sono:

&	=	&
<	=	<
>	=	>
"	=	"
'	=	'

Altre caratteristiche del documento XML - Sezioni CDATA e Commenti

- CDATA sono delle sezioni di testo che il parser XML non cerca di interpretare. Tutte le occorrenze di & in una sezione CDATA, per esempio, verranno letti dal parser come & e un simbolo < non sarà interpretato come istruzione di markup. Le sezioni CDATA s'iniziano con <![CDATA[e terminano in]]>.
- Il commento è una sezione speciale che comincia con <!-- e finisce con --> . Tutti i dati scritti fra questi due tag sono ignorati dal processore XML. I commenti sono per lo più utilizzati per aggiungere brevi note nel documento XML, o per commentare intere sezioni del documento XML.

Dichiarazioni di entità ENTITA' GENERALI INTERNE ANALIZZATE

Si possono dichiarare nella DTD delle entità, cioè degli oggetti che risiedono esternamente al documento, che vengono poi richiamate all'interno del documento nella forma &nome;

```
<!ENTITY nome "stringa da sostituire">
```

Nel documento la sostituzione avverrà tutte le volte che l'elaboratore leggerà &nome;

Esempio:

```
<!ENTITY ft "Francesca Tomasi"> da richiamare nel documento XML con &ft;
```

Dichiarazioni di entità

ENTITA' GENERALI ESTERNE

L'oggetto da sostituire al riferimento di entità è un oggetto che si trova esternamente al documento XML; può essere un oggetto di pubblico dominio (PUBLIC) o un oggetto "privato" creato specificatamente (SYSTEM):

```
<!ENTITY file SYSTEM "http://www.miosito.it/pippo.xml">
```

```
<!ENTITY file SYSTEM "http://html.it/pippo.xml">
```

Il file viene recuperato dal sistema tutte le volte che trova la stringa &file;

Inserire la DTD nel prologo

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
```

```
<!DOCTYPE nomeelementoradice [
```

```
  <!ELEMENT nomeelementoradice (contenuto)>
```

```
]>
```

```
< nomeelementoradice >
```

Questo è l'elemento superiore del documento.

```
</ nomeelementoradice >
```

Una DTD comincia con <!DOCTYPE e termina con]>. Questo segnala al processore XML dove la DTD inizia e finisce.

Direttamente dopo il <!DOCTYPE viene il nome dell'elemento radice seguito da un [.

Richiamare la DTD nel prologo

`<!DOCTYPE nomelementoradice SYSTEM "nomedtd.dtd">`

DTD ad uso privato, generalmente situata nell'hard disk (e in questo esempio deve essere anche nella stessa cartella)

`<!DOCTYPE nomelementoradice PUBLIC nomeconvenzionale "URL/nomedtd.dtd">`

DTD ad uso pubblico. In questo caso alla DTD viene assegnato un nome univoco che l'elaboratore XML cerca per procurarsi la DTD (oppure utilizza l'URL specificato).

Es: `<!DOCTYPE TEI.2 PUBLIC "-//TEI P3//DTD Main Document Type//EN">`

In questo caso indica che il documento è della TEI.2, che la DTD è stata sviluppata dalla TEI P3 e che la lingua utilizzata è l'inglese (EN).