LOS LENGUAJES DE MARCADO APLICADOS A LOS REGISTROS BIBLIOGRÁFICOS. XML MARC DTD; XML MARC SCHEMA.

I. INTRODUCCIÓN A LOS LENGUAJES DE MARCADO

Los **lenguajes de marcado**, también denominados lenguajes de marcas, provienen y toman su nombre de la <u>práctica tradicional e histórica de marcar los manuscritos con instrucciones de impresión en los márgenes</u>. Esta tarea, tan habitual en la industria editorial desde la aparición de la imprenta, fue dando lugar a un grupo de **marcas estandarizadas**, cuya <u>esencia</u> ha sido trasladada hoy al <u>mundo de la informática</u>.

Actualmente, se distinguen **dos tipos fundamentales de lenguajes de marcado**, si bien en la práctica pueden combinarse:

- LENGUAJES DE MARCADO DE PROCEDIMIENTO, orientados hacia la presentación
 del texto. Sus símbolos o marcas indican la clase de operaciones tipográficas que
 deben ser aplicadas a cada uno de los elementos del documento electrónico, para dar
 formato al texto. Su misión es, pues, configurar la apariencia física de éste (fuente,
 estilo, tamaño de letra, etc.), tanto en pantalla como impreso.
- LENGUAJES DE MARCADO DESCRIPTIVO, diseñados para identificar las piezas o bloques estructurales que componen el texto. Sus marcas determinan la estructura lógica del documento electrónico y/o la descripción de su contenido, no ya su tipografía ni el formato que presentará cada fragmento en su posterior visualización.

Así pues, en el contexto automatizado actual, los lenguajes de marcado permiten **codificar documentos** <u>intercalando</u>, junto al texto, <u>etiquetas o marcas provistas de información adicional</u> <u>sobre su estructura o presentación</u>. En realidad, más que de lenguajes, podría hablarse de **metalenguajes o conjuntos de reglas** que se encargan de establecer y definir la <u>forma digital</u> de los documentos, bien para controlar su procesamiento, bien para representar su significado.

ORIGEN Y EVOLUCIÓN

A pesar de existir proyectos anteriores, como el **GenCode**, la iniciativa que sentó las bases de los lenguajes de marcas actuales partió en los **años 60** del investigador de IBM **Charles F. Goldfarb**, considerado hoy padre de este tipo de lenguajes por su participación en el desarrollo del *Generalized Markup Language* o **GML**.

La generalización de los lenguajes de marcas: SGML (años 80)

El gran éxito del lenguaje GML propició su extensión e incluso adopción por parte del gobierno de los EE.UU., hecho que acrecentó la necesidad de estandarizarlo. Así, tras un largo proceso, en 1986, se convirtió en norma ISO bajo el nombre de SGML (Standard Generalized Markup Language).

ISO 8879:1986

A pesar de su falta de precisión y dificultad, SGML fue la **piedra angular** o **referente definitivo** de los lenguajes de marcado modernos, ya que, además de proveer una <u>sintaxis para la inclusión de marcas en los textos</u>, introdujo por primera vez una <u>sintaxis para especificar qué etiquetas estaban permitidas y dónde</u>. Éste sería el **punto de partida de dos conceptos clave**, abordados más adelante: **DTD** (*Document Type Definition*) y **Schema**.

La popularización: HTML (años 90)

El Lenguaje de Marcas de Hipertexto o <u>HTML (HyperText Markup Language)</u> fue definido en **1990**, a partir de la sintaxis de SGML, por **Sir Tim Berners-Lee**, creador de la WWW. En la actualidad, constituye la tecnología predominante en la construcción de páginas web, dada su simplicidad para <u>estructurar textos</u> y su capacidad para <u>establecer enlaces</u> con otros archivos. El **Consorcio World Wide Web (W3C)** regula las recomendaciones y versiones normalizadas de este lenguaje, aceptado como **norma ISO** desde el año **2000**.

La <u>estructura HTML clave</u> se basa en el uso de **ELEMENTOS**, compuestos mayoritariamente por **etiquetas**, **atributos** y **contenido**:

- Etiquetas. Están representadas por corchetes angulares (< >) con una instrucción en su interior. Cuando los elementos presentan una etiqueta inicial, texto y una etiqueta final, se denominan <u>llenos</u>. Cuando no presentan contenido ni necesitan cerrarse mediante etiqueta alguna se conocen como <u>elementos vacíos</u>.
- Atributos. Se incluyen en la etiqueta inicial del elemento pertinente, justo detrás de la instrucción (estructura nombre="valor"), y recogen la <u>variable a aplicar</u>.
- Contenido. Es el texto propiamente dicho, consignado entre las etiquetas y posteriormente reconocido como información en la visualización HTML.

Todo documento HTML presenta dos zonas claramente diferenciadas: cabecera (*HEAD*) y cuerpo (*BODY*).

elemento con argumento

La cabecera alberga el <u>título del documento</u> que aparecerá <u>en la parte superior de la ventana del navegador</u> así como otros datos no explícitos a posteriori, pero de interés para el servidor o para los buscadores web.

• El cuerpo contiene la información que el usuario verá en su navegador.

Puesto que HTML fue diseñado, en origen, para <u>intercambiar información en entornos académicos</u>, sus etiquetas están eminentemente pensadas para la **organización lógica del contenido** (titulo, párrafo, etc.) y no tanto para su presentación. Por eso, el W3C tuvo que idear las denominadas **hojas de estilo** (*Style Sheets*).

La madurez: XML (años 2000)

La respuesta a las deficiencias surgidas en el entorno HTML fue la aparición, en **1998**, del **lenguaje XML** (*eXtensible Markup Language*), desarrollado y evolucionado por el **W3C**.

Se trata de un lenguaje **sencillo**, que propicia la <u>compatibilidad</u> entre sistemas, permite incluir <u>enlaces multidireccionales</u> y cuenta a su alrededor con <u>tecnologías complementarias</u>, como el **Lenguaje de Hojas de Estilo Extensible** (**XSL**, *Extensible Style Language*), promovido también por el W3C para la presentación de documentos XML.

Desde el <u>punto de vista estructural</u>, XML guarda ciertas similitudes con HTML, ya que sus documentos cuentan con dos partes claramente delimitadas, **PRÓLOGO** y **CUERPO**, y sus componentes presentan denominaciones parecidas. Si bien el **prólogo** es <u>opcional</u>, el **cuerpo** es completamente <u>obligatorio</u>, pues recoge el contenido en sí del documento, dividido en texto y marcado, con **ELEMENTOS** (llenos o vacíos).

Como en HTML, los <u>componentes clave del marcado</u> son las **etiquetas** y sus **atributos**, aunque también cobran especial importancia los **comentarios** o las denominadas **entidades predefinidas** y **secciones CDATA**, ambas necesarias para representar caracteres propios de las marcas, que, en vez de procesarse como tal, deben figurar en la visualización posterior. Los <u>datos consignados entre las marcas</u> constituyen el **contenido**, perfectamente legible y sin codificación alguna.

Aunque se denomina "**extensible**", porque <u>no limita el número de marcas o etiquetas posibles</u>, es un lenguaje especialmente **estricto** en cuanto a lo que está permitido y lo que no. Así, <u>todo documento XML debe cumplir **dos condiciones**:</u>

- Estar **bien formado** (*well-formed*), es decir, guiarse por lo especificado en la recomendación W3C.

Ser válido (fully validated), es decir, respetar las restricciones establecidas por una definición externa, DTD o XML Schema. Pero, ¿qué son exactamente ambos conceptos?

II. INTRODUCCIÓN A LOS CONCEPTOS DTD Y SCHEMA

DTD (Definición del Tipo de Documento)

Una DTD define los tipos de <u>elementos, atributos y entidades permitidos</u> así como sus posibles <u>combinaciones</u>. Las primeras DTD fueron desarrolladas en los laboratorios de **IBM**, hacia **1978**, cuando aún se experimentaba con el futuro lenguaje SGML. Actualmente, <u>las más comunes</u> son las usadas para entornos HTML y, sobre todo, para XML.

Su <u>función básica</u> consiste en **describir la estructura y sintaxis de los datos**, proporcionando un **formato común** que dé consistencia a todos los documentos regidos por la misma DTD. Aunque ésta puede incluirse <u>dentro del propio documento</u>, normalmente se almacena <u>en un fichero aparte vinculado</u>.

Las **limitaciones** y la **rigidez** de las DTD favorecieron la <u>aparición de otras herramientas de</u> <u>descripción estructural</u>, alternativas y más completas, como los denominados **Schemas**.

SCHEMAS

Un Schema es **similar a una DTD** en el sentido de que define <u>qué elementos</u> están permitidos, <u>cómo deben organizarse</u> y <u>qué tipo de atributos</u> pueden albergar, pero añade **VENTAJAS** como:

- El uso de la sintaxis de XML.
- Una mayor especificación del tipo de datos.
- Su extensibilidad.

El resultado, esto es, **XML Schema**, es un **lenguaje de esquema** utilizado para describir, con total precisión, la <u>estructura y restricciones de contenido de los documentos XML</u>, más allá de las normas sintácticas impuestas por el propio XML. Fue desarrollado por el **W3C** y alcanzó el nivel de **recomendación** en mayo de **2001**.

III. LENGUAJES DE MARCADO APLICADOS A LOS REGISTROS BIBLIOGRÁFICOS MARC

Los **formatos de marcas** constituyen desde hace tiempo una poderosa <u>alternativa</u>, <u>o complemento</u>, <u>a los sistemas tradicionales de codificación de datos</u> para describir recursos electrónicos. Los proyectos más importantes en este ámbito, liderados desde los años 90 por la *Library of Congress*, siguen, en la actualidad, **dos líneas fundamentales de trabajo** relacionadas con los lenguajes de marcado:

- Diseño de modelos de descripción alternativos a MARC, como el Dublin Core (DC).
- Adaptación de los modelos tradicionales MARC (norma ISO 2709) a los nuevos formatos de Internet.

Para muchos especialistas, el <u>formato MARC debe seguir primando frente a otros sistemas</u> de metadatos propuestos, como el DC, ya que:

- Sus más de 30 años como soporte de la comunidad bibliotecaria y de la industria de software especializado lo avalan. La mayoría de opciones restantes se encuentran aún en fase de desarrollo, no pudiéndose garantizar su permanencia en un futuro.
- Posee mayor capacidad expresiva, además de una semántica claramente definida, aceptada y adaptada a las distintas realidades nacionales.
- Da cabida a una amplia gama de registros (datos bibliográficos, autoridades, fondos y localizaciones, etc.).
- Se ha convertido en la base de los SIGB actuales.

Por éstas y otras razones, se han realizado diversos <u>esfuerzos para ligar MARC e Internet</u>, empezando por la inclusión del **campo 856** (*electronic location and access*) o la **adaptación de MARC a SGML** durante los años 90. Para ello, se crearon **dos DTD** capaces de convertir registros <u>MARC a SGML y viceversa</u>, sin <u>pérdida de información</u>, que pronto migrarían a XML para adecuarse a las nuevas necesidades tecnológicas.

MARC DTD

Oficina de Desarrollo de Redes y Normas MARC

En **2002**, la *Network Development and MARC Standards Office* (NDMSO) de la **LC** publicó un **esquema XML**, yendo <u>más allá de un mero mecanismo de conversión</u> y facilitando la <u>representación de registros MARC en formato XML</u>, para eliminar complejidades innecesarias y evitar que MARC quedara relegado frente a otras propuestas, en el marco de la biblioteca electrónica.

IV. MARC XML DTD

MARC SCHEMA

Las MARC XML DTD, derivadas de las mencionadas DTD para SGML iniciales, se dividen en:

- 1. XML DTD para registros bibliográficos, fondos e información a la comunidad.
- 2. XML DTD para registros de autoridad y clasificación.

Surgieron en 2001, con los objetivos relatados, de acuerdo con los siguientes PRINCIPIOS:

- 1º Generalidad, por su independencia respecto a aplicaciones específicas basadas en MARC.
- 2º **Reversibilidad**, por su <u>capacidad para convertir datos de una estructura a otra</u> y poder volver al formato original, sin pérdida de contenido intelectual o de elementos semánticos esenciales.
- 3º Flexibilidad, ya que ofrece diversas posibilidades, en vez de imponerlas directamente.
- 4º Amigabilidad para el usuario (del inglés user-friendly), gracias a sus jerarquías lógicas.
- 5º **Relación con TEI** (*Text Encoding Initiative*), otra iniciativa sólida para el <u>intercambio de</u> información textual y para la descripción normalizada de documentos electrónicos.

V. MARC XML SCHEMA

Con este nuevo **avance** de la NDMSO de la **LC** se ha logrado una <u>conversión más simple y</u> <u>flexible de los registros MARC</u>, como norma ISO 2709, <u>al lenguaje XML</u>, adaptado a la WWW. En concreto, la utilización de la tecnología **XML Schema**, para definir la estructura de los registros MARC actuales, reporta **VENTAJAS** como:

- Soporta todos los datos codificados con MARC, <u>independientemente de la variante</u> <u>utilizada</u>, y los transforma para que sean perfectamente legibles en un entorno web.
- La **conversión** de MARC a XML y la **recomposición** de XML a MARC se realiza <u>sin</u> <u>pérdida de información</u> alguna.
- Posee una arquitectura extensible para describir recursos originales en sintaxis XML, permitiendo realizar modificaciones individuales o por lotes.
- Posibilita la transformación de MARC XML a otros formatos de metadatos.

La **gran mejora** de esta tecnología con respecto a las XML DTD radica en que, cuando se traduce un registro MARC con XML Schema, <u>se mantienen todos y cada uno de los campos, indicadores y subcampos MARC</u>, tal cual. Las DTD, por el contrario, obligan a redactar de nuevo toda esa información en formato XML, lo cual acaba por crear una DTD muy grande y con muchas más líneas.

VI. CONCLUSIÓN

En definitiva, los lenguajes de marcado son herramientas indispensables para la <u>estructuración</u> <u>y presentación de los recursos bibliográficos en línea</u>. Hoy, además, se postulan como elementos clave en el diseño de la **Web Semántica**, aquella que no sólo permite acceder a la <u>información</u>, sino también definir su <u>significado</u>, para facilitar su procesamiento automático.

XML, a través de sus **DTD** y **Schemas**, ha hecho de **MARC** un poderoso estándar descriptivo, también en el ámbito electrónico. No obstante, no faltan las **voces disonantes** que advierten de su <u>complejidad</u>, <u>lentitud</u> y <u>coste</u> ante grandes cantidades documentales. Por eso, han surgido otras alternativas, como el **DC** o el lenguaje **RDF**, que ofrecen mayor rapidez descriptiva, aunque en menor detalle.

PUNTOS TRATADOS EN EL TEMA

LOS	LENGUAJES	DE MARCADO	APLICADOS	A LOS	REGISTROS	BIBLIOGRÁFICOS.
XML	MARC DTD; X	ML MARC SCHE	EMA.			

VI. CONCLUSIÓN

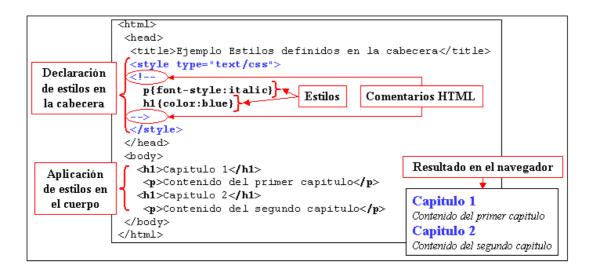
ANEXO

Reglas de la hoja de estilo aplicadas a un elemento en HTML:

```
Mediante el atributo style se puede alterar el
aspecto del elemento al que se le aplica.
RESULTADO (aplicado al elemento párrafo):

Mediante el atributo style se puede alterar el aspecto del elemento al que
se le aplica.
```

Reglas de la hoja de estilo aplicadas a la cabecera de un documento en HTML:



Reglas de la hoja de estilo vinculadas (a uno o varios documentos) como fichero externo:

