

El Metalenguaje XML y el Esquema de Tipo de Elemento

JUAN VOUTSSÁS MÁRQUES

Centro Universitario de Investigaciones Bibliotecológicas
de la UNAM, 04510, México D.F., Tel: 56-23-03-29
E-Mail: voutssás@servidor.unam.mx

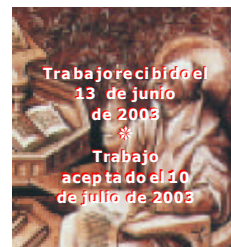
RESUMEN

En las últimas décadas se ha ido desarrollando una herramienta para el procesamiento de información electrónica conocida como los metalenguajes de marcado. Desde el formato "MARC" de los sesentas hasta el "XML" es tab le ci do en nues tros días ha ha bi do sin duda una evo lu ción en este sen ti do bas tan te no ta ble. El pre sen te do cu men to tra ta de es ta ble cer una evo lu ción his tó ri ca y fun cional sobre los len gua jes de mar ca do, sus ca rac te rís ti cas, ven ta jas, li mi ta ciones y sus ex pec ta ti vas.

Dentro de las especificaciones "XML" se ha desarrollado recientemente una técnica particular de declaración y manejo de la información documental conocida como "esquema de tipo de elemento" o simplemente es que mas, la cual va más allá de la simple y ya co no ci da de fi ni ción del do cu men to a tra vés de un des cri p tor de tipo de do cu men to (DTD). Los es que mas nos ofre cen una po si bi li dad mu cho más am plia para de fi nir cla ses de do cu men tos. Este tra ba jo in ten ta es ta ble cer esas po si bi li da des para po der en ten der la po ten cia li dad de los es que mas, su es ta do del arte, sus ten den cias y fi nal men te la im por tan cia que el de sa rrol lo de es que mas es pe cia les para cier tas cla ses de do cu men tos pue de te ner en nues tro me di o bi bli o te ca ri o me xi ca no para un me jo res ta ble ci mien to y ex plo ta ción de las co lec ciones di gi ta les.

Pa la bras Cla ve: Me ta len gua jes, Len gua jes de mar ca do do cu men tal, XML, HTML, SGML, MARC, DTD, Es que ma de tipo de ele men to, Es que mas, Bi bli o te ca s di gi ta les, De sa rrol lo de co lec ciones di gi ta les, Tec no lo gías de in for ma ción

Artículo



*METALANGUAGE XML AND THE TYPE
OF ELEMENT SCHEME*
JUAN VOUTSSÁS-MÁRQUES

ABSTRACT

During the last decades a tool for electronic processing known as markup languages has been and is still being developed. From the MARC Format of the 60's to the XML established in our days, a very clear evolution has, no doubt, taken place. This paper intends to settle the issue of the historical and functional evolution of the mark-up languages and of their characteristics, advantages, limitations and expectations.

Within the specifications of XML a particular technique for the declaration and handling of the documentary information has been developed recently. The technique is known as Type of Element Scheme, or simply Schemes, but it goes well beyond the simple and very well known definition of document through a Descriptor of Type of Document (DTD). The schemes offer a much wider possibility for defining classes of documents. This paper attempts to establish the possibilities in order to understand the potentiality of Schemes, its state of art, its tendencies and also the importance that the development of special schemes may have for certain classes of documents in our Mexican library environment for a better implementation and exploitation of digital collections.

Key Words: Metalanguages; Document Markup Languages; XML; HTML; SGML; MARC; DTD; Esquemas; Digital Library; Digital Collection Development; Information Technology

ANTECEDENTES

Desde la década de los sesenta, el personal técnico relacionado con el intercambio de documentos por medio de una computadora comenzó a preocuparse por estructurar esa información en una forma normalizada que facilitara el intercambio y la manipulación de esos materiales. En el ámbito de las bibliotecas surgió el proyecto del formato “MARC”, (*Machine Readable Cataloguing*), pionero en esta temática. Durante 1965 la Biblioteca del Congreso de los Estados Unidos (LC) estableció un proyecto piloto denominado “MARC I” cuyo objetivo era definir una metodología para crear registros catalográficos en un formato legible **por una computadora**. Un proyecto similar apareció en el Reino Unido en esa misma época ordenado por el Consejo de la Biblioteca Británica (British Library) con el fin de establecer un formato para registrar y explorar los registros que se rían usados en la Bibliografía Británica. A este proyecto se le llamó “BNB MARC” (*British National Bibliography with Machine Readable Cataloguing*). Ambos proyectos derivaron en poco tiempo en una cooperación en el ambiente de las bibliotecas anglo-sajonas y dieron origen, en 1968, al proyecto “MARC II”. El resultado fue el formato ampliamente conocido para intercambiar registros catalográficos de monografías vía electrónica, el cual fue complementado poco después con los consecuentes formatos para publicaciones periódicas, mapas, discos, etcétera, hasta completar prácticamente todo los tipos de documentos que existen en las bibliotecas.

Con el tiempo MARC fue tan exitoso y tan ampliamente utilizado en todo el mundo que empezaron a usarse múltiples variantes del mismo, hasta llegar a unas veinte que fueron adaptadas a las nuevas necesidades de otras comunidades en el planeta. Para 1977 se definió el formato universal “UNIMARC” como una especie de “espejito” para MARC con objeto de que los distintos “MARC” pudieran intercambiar registros entre sí usando “UNIMARC” como intermediario. Además el formato también

empezó a ser de finido para otras en tida des de da tos más allá de la simple pro ducción de ca tá los, como los regis tros de au to ri dad, pues su di se ño per mi tía que fue ra usa do en prác ti ca men te cual quier pla ta for ma de có m pu to del mun do. De bi do a su am plia acep ta ción MARC se con vir tió en nor ma AN SI para la unión ame ri ca na en 1971 (Z39.2 - 1971) y en nor ma ISO in ter na cional en 1973 (ISO 2709: 1973 E). Está de más ahon dar sobre lo que MARC re pre sen tó para el regis tro elec tró ni co de in for ma ción. [*Library of Congress*, 2003]

En esen cia, MARC no ha sido nun ca un ca tá lo go ni un mé to do para ca ta lo gar, como mu chos han pen sa do a lo lar go del tiem po. La **gran a por ta ción de MARC al mane jo documen tal**, y que ha sido se gui da por mu chos for ma tos poste rior men te has ta nues tros días, fue la de aso ciar una “eti que ta” a cada uno de los ele men tos que con for man una fi cha ca ta lo grá fi ca, con ob je to de que las com pu ta do ras pue dan iden ti fi car y ca ta lo gar sus par tes y poste rior men te efec tuar múl ti ples ac ciones con ellas en be ne fi cio de las per so nas. Este prin ci pio tan sim ple y tan ob vio en nues tra épo ca mar có un hi to en la his to ria del mane jo documen tal por me dio de com pu ta do ras. Hoy se ha so fis ti ca do y di ver si fi ca do en or me men te, pero el prin ci pio se man tie ne has ta nues tros días. [A general in tro duc tion to marc] [If lanet. *What is marc*]

De he cho, este con cepto de aso ciar cada ele men to de una fi cha, y pos te rior men te de cual quier otro do cu men to, con una “eti que ta” (lla ma da en in glés *la bel, tag* o *token*) fue con o cián do se como “mar ca do” del do cu men to, y ha sido de tal tran scen den cia que to dos los for ma tos que han apa re ci do pos te rior men te para el regis tro de do cu men tos se ba san en el con cepto del mar ca do, y así lo os ten tan en su nom bre (*markup* o *mark up*.) Abun da ré más ade lan te en este con cepto. Por lo an te rior, MARC pue de ser con si de ra do el “abuelo” de to dos los len gua jes de mar ca do documen tal para com pu ta do ra que exis ten hoy.

Fue ra de las bi bli o te cas, a ni vel ge neral y por esa mis ma épo ca de fi nes de los se sen tas, la em pre sa IBM creó “GML” (*Generalized Markup Language*) o “Len gua je de mar ca do ge neral i za do” para con ten der con las ne ce si da des de sus pro pios sis te mas in ter nos de pu bli ca ción. Ellos usaron GML para pro du cir de mane ra nor ma li za da, li bros, re por tes, ma nua les y otros ti pos de do cu men tos a par tir de un ú ni co con jun to de ar chivos ori gi na les en una com pu ta do ra. Mu chas otras so lu ciones fue ron in tro ducidas por otras or ga ni za ciones para es truc tu rar in for ma ción documen tal, pero na da real men te que abar ca ra este pro pósito a gran es ca la.

La pri me ra tec no lo gía de es truc tu ra ción de in for ma ción documen tal nor ma li za da de cier ta sig ni fi can cia fue el “SGML” (*Standard Generalized Markup Language*) o “Len gua je de mar ca do ge neral i za do es tán dar”, en lo su ce si vo “SGML”, el cual tam bién pro vi no de la em pre sa IBM a prin ci pios de los och en tas. Este len gua je fue cre a do con el fin de for ma te ar y or ga ni zar la documen ta ción legal den tro de esa em pre sa, pero poste rior men te fue ex pan di do y adap ta do para ser usa do en una am plia va ri e dad de em pre sas como un es tán dar para mane jar todo ti po de in for ma ción y para 1986 se con vir tió en una nor ma ISO [ISO 8879, 1986].

Si bien SGML es extremadamente poderoso, es igualmente complejo y requiere de una considerable cantidad de programación adicional para procesarlo. Debido a su complejidad y los recursos extra que demanda, SGML no era una opción viable para representar hipertextos en las primeras épocas de la Internet, dadas las dimensiones y limitaciones de los equipos de esa época.

En 1989, Tim Berners-Lee y Anders Berglund, dos investigadores del Laboratorio Europeo de Partículas Físicas (CERN), crearon un lenguaje basado en etiquetas para marcar documentos técnicos y ser éstos compartidos en Internet. Este lenguaje fue expandido en 1990 a una versión simplificada del SGML llamada HTML (*Hyper-Text Markup Language*) “Lenguaje de marcado de hipertextos”, y ha llegado a ser el formato estándar para el manejo de información en la “Web”.

HTML podía representar perfectamente información *estática* en una página Web; es decir, textos previamente establecidos, imágenes, botones, ligas, etcétera, pero pronto hubo necesidad de ligar las páginas a bases de datos con el fin de traer a la pantalla catálogos, estados de cuenta, textos completos, etcétera; es decir, información *dinámica* o *cambian te*, para lo cual HTML no fue diseñado. Algunos añadidos han surgido para HTML en los últimos años, para resolver esta problemática, algunos de ellos muy exitosos, como los derivados del lenguaje “Java”, y otros han caído ya en el olvido. Al final, HTML comienza a estar demasiado “recomendado” y resulta ya complejo para poder satis hacer los requerimientos del manejo de información de la actualidad. [Teasdale, 1995]

Como resultado de esto muchos técnicos han empezado a volver sus ojos nuevamente hacia SGML, ya que su complejidad era de origen y no resultado de adiciones, y siempre han tenido capacidad para manejar documentos complejos de varios tipos. Otros grupos de personas empezaron a reescribir una versión simplificada de SGML pero ca paz de con tener con las carencias de HTML; de ahí surge un nuevo lenguaje llamado XML: el (*eXtensible Markup Language*) o “Lenguaje de marcado extensible”, en lo sucesivo “XML”. En 1996, el “Consortio para el desarrollo de la World Wide Web”, conocido como el W3C, sentó las bases para estos desarrollos. Se establecieron las ventaj as propias del SGML: estructura, extensibilidad y validación, y se creó un grupo de trabajo que estableció las bases para un nuevo lenguaje de marcado que conserva las ventaj as cen trales del SGML y que tu vie ra la sim pli ci dad del HTML; esto es, un SGML “aero di ná mi co” para la Web. El resultado fue que en 1998 se publi ca ron las especificaciones de la primera versión del XML [Extensible Markup Language 1.0, 2000]. En cuanto a dimensiones, la especificación XML resultó ser de menos de una décima parte de la de SGML, con lo cual puede estimarse el grado de compactación. [MSDN on Line, 2000]. Desde ese año a la fecha, el con sor cio W3C se ha cons titui do como la instancia oficial a nivel mundial para definiciones del estándar XML.

CONCEPTOS BÁSICOS

Metalinguaje

Si se observa la literatura al respecto, puede notarse que HTML, SGML, XML y otros entes del manejo de información documental son denominados como “lenguajes”, pero más propiamente dicho, se trata de “metalinguajes”. ¿Qué significa realmente este término? es importante aclararlo para poder comprender el contexto de estas herramientas.

El concepto de “metalinguaje” no es simple y pretende establecer las reglas por las cuales un lenguaje dado puede ser de finido o examinado **formalmente**, pero existen muchas clases de lenguajes: los lenguajes que los humanos usamos para comunicarnos cotidianamente, español, inglés, francés, etcétera; lenguajes para programar una computadora, Fortran, Algol, “C”, Pascal, Cobol, Java, etcétera; el álgebra, etcétera.

Así pues existen muchos metalinguajes; más de los lenguajes que pueden ser creados. Incluso un lenguaje cualquiera que hablamos las personas incluye varios “sublenguajes” o “argots” propios de una profesión o clase social. Distinguiamos por ejemplo el lenguaje propio de los médicos dentro de un lenguaje dado, o el del personal técnico de cómputo, etcétera. Los enfoques pueden ser entonces filológicos, matemáticos, computacionales, etcétera. Para fines de este trabajo nos centraremos en los lenguajes “documentales”, que son los de nuestro interés.

De esta forma podemos establecer un metalinguaje para definir un lenguaje dado y por tanto, **definiremos entonces de una forma simple “metalinguaje” como un medio para definir formalmente un lenguaje dado**. Si bien entran también aquí consideraciones de sintaxis, contexto, objetos, etcétera, concentrándonos en el tema que nos ocupa; es decir los documentos, resulta que los metalinguajes (MARC, SGML, HTML, XML) son un medio para describir un lenguaje de marcado; dicho de otra forma, son metalinguajes para definir formalmente un lenguaje de marcado de documentos.

Las reglas para definir el marcado de documentos, su sintaxis, sus limitaciones y características están dentro de cada uno de esos entes (MARC, SGML, HTML, XML); es decir, ellos internamente van con formando todo un lenguaje para marcar o codificar documentos. El lenguaje está por lo tanto intrínsecamente contenido en cada uno de ellos, pero la manera de presentar el conjunto de las reglas y definiciones de ese lenguaje es un metalinguaje. Por eso decimos que estas herramientas son metalinguajes más propiamente dicho que lenguajes.

Históricamente, la palabra marcado (*markup*) viene de las instrucciones que se le daban a un impresor o tipógrafo sobre cómo debería imprimirse un cierto pasaje de un texto: las negritas, los subrayados, la tipografía, el tipo y el tamaño de letra, los símbolos especiales, etcétera, se señalaban con anotaciones o “marcas” al margen del texto; de ahí el nombre. Con los textos se fueron automatizando, el término se fue extendiendo hasta cubrir toda clase de códigos de marcado insertados en esos

textos electrónicos con objeto de incrustar una serie de características propias de cada sección o parte del texto.

De este concepto se desprende que el marcado es una codificación del texto, gracias a la cual cada porción del mismo se va haciendo explícita y toma forma en un contexto. De hecho, y hablando formalmente, todos los textos han tenido desde hace largo tiempo una manera de hacerse explícitos y entendibles al lector. Los signos de puntuación, los acentos, los espacios entre palabras, el uso de mayúsculas, los párrafos, los capítulos, los incisos, etcétera, se han agregado a los textos desde hace mucho tiempo para que sus lectores puedan leer explícitamente un conjunto de letras y dárles un sentido. Imagine se el lector un texto en el cual todas las letras están con tilde, sin espacios ni puntuaciones ni nada adicional, y podrá darse cuenta de la enorme ayuda que estos elementos proporcionan para hacer explícitas estas palabras. Si el lector no se lo imagina analice entonces este texto:

adanyrazaazarynada

Como puede observarse, el texto anterior difícilmente puede ser interpretado al primer golpe de vista; es necesario leerlo y releerlo varias veces para que el cerebro humano, tan avezado en este tipo de actividades, pueda establecer patrones, cortes, analogías, y describir su significado y su contexto. Puesto de esta forma:

palíndromo:

Adán y raza, azar y nada

es mucho más fácil de interpretar. Ciertas convenciones de espacios y puntuación lo hacen explícito. Incluso, gracias al encabezado los lectores se cercioraron de que el texto es un palíndromo, es decir, algo que puede ser leído igualmente a la derecha que de izquierda a derecha. Algunos no se habrán percatado de ello hasta que el encabezado, es decir, cierto *marcado*, lo hizo explícito.

Del mismo modo, y en ello radica su importancia, un lenguaje de marcado debe permitir hacer explícito un texto **para los sistemas y las máquinas** que los operan electrónicamente: además de dónde empieza y termina cada palabra, lo cual puede lograrse con espacios, nos indica dónde empieza una parte de un texto en un documento y dónde acaba, qué partes tiene ese documento y cómo se llaman, etcétera. El lenguaje de marcado es entonces un conjunto de convenciones preestablecidas que se usan de manera conjunta y ordenada para marcar o codificar un texto de modo que éste pueda ser entendido explícitamente por máquinas para que éstas, a su vez, se lo puedan hacer explícito a las personas.

Un lenguaje de marcado debe especificar entonces cuáles marcas son permitidas, cuáles son obligatorias u opcionales, cómo deben diferenciarse esas marcas del texto propiamente dicho y, por supuesto, qué significa cada marca.

Trataré de ilustrar con un ejemplo los temas anteriormente expuestos.

Usted quiere describir una parte de un documento que es la dirección postal de una persona dada, la cual, como sabemos, está formada a su vez por varios elementos. ¿Cómo describiríamos esta entidad en lenguaje llano pero a la vez **formalmente**?

Tal vez sería algo como esto:

Una dirección postal consiste de: una parte-que-tiene-el-nombre, seguida de una parte-que-tiene-la-dirección, seguida de una parte-que-tiene-el-código-postal.

La parte-que-tiene-el-nombre consiste de: un nombre-de-pila, seguido de un apellido-paterno, seguido de un apellido-materno, si se tiene.

La parte-que-tiene-la-dirección consiste de: un nombre-de-calle, seguido de un número-exterior, seguido (si lo hubiese) de un número-interior, seguido de un nombre-de-colonia.

La parte-que-tiene-el-código-postal consiste de: un nombre-de-ciudad, seguido de un nombre-de-entidad-federativa, seguido de un código-postal

Un código-postal consiste de un número de cinco dígitos.

Usemos un metalenguaje para describir esto. Por simplicidad, usaremos la “forma de Backus-Naur” (Backus-Naur Form o BNF), un metalenguaje muy utilizado para definir lenguajes de programación desde 1959 y que toma el nombre de sus creadores John Backus y Peter Naur. [*Free Dictionary*...2000].

Las bases de construcción del metalenguaje son las siguientes; primero los símbolos:

::=: significa “se define como”

|: significa “o”

<>: corchetes angulares usados para encerrar nombres de categorías o “entes”.

{ }: corchetes de llave usados para encerrar nombres de categorías opcionales.

“ ”: Texto entre comillas que se integra literalmente tal como está escrito.

<EOL>: fin de la línea de definición

La manera de usar la forma de Backus-Naur es la siguiente: se pone del lado izquierdo la entidad a definir y del lado derecho la definición, separadas por el signo de “se define como”. Se trata de pasar de lo general a lo particular y se evitará a toda costa la tautología; es decir que la parte a definir sea igual que la definida. No obstante, **una parte** de lo que está a la derecha como definida puede ser parte de lo que se define. Por ejemplo, en el álgebra, la expresión $n = n + 1$ no puede existir, ya que de acuerdo a sus reglas, no hay ningún número que sea igual a él mismo más uno. En notación BNF $n ::= n + 1$ sí puede existir y significa que n se define como él mismo incrementado en una unidad. Perfectamente a manera de establecer un contador. Si se observa el ejemplo anterior, se ve que n forma parte de lo que se define y de lo definido, pero es correcto bajo estas reglas y no es una tautología.

Si en algún momento una de finición puede ser de finida a su vez en componentes más simples, esa nueva definición se escribe a continuación.

De acuerdo con lo anterior, podemos definir más formalmente una dirección postal según la notación de Backus-Naur como:

```

<dirección-postal> ::= <parte-que-tiene-el-nombre> <parte-que-tiene-la-dirección> <parte-que-tiene-el-código-postal>
<parte-que-tiene-el-nombre> ::= <nombre-de-pila> <apellido-paterno> {apellido-materno}
<nombre-de-pila> ::= <nombre> | <inicial> "." <EOL>
| <nombre> <nombre-de-pila>
<parte-que-tiene-la-dirección> ::= <nombre-de-calle> <número-exterior> {número-interior} <nombre-de-colonia>
<parte-que-tiene-el-código-postal> ::= <nombre-de-ciudad> <nombre-de-entidad-federativa> <código-postal>
<código-postal> ::= <dígito> <dígito> <dígito> <dígito> <dígito>
<dígito> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
    
```

De acuerdo con este ejemplo, y utilizando la notación BNF podemos observar lo siguiente:

- ❖ Hemos definido dirección-postal como una categoría formada por tres partes (nombre, dirección y código) y éstas de ben existir siempre y en ese orden.
- ❖ Hemos definido a su vez al nombre como una categoría formada por el nombre-de-pila, seguido por el apellido paterno y opcionalmente, el apellido materno, y en ese orden.
- ❖ Hemos definido también al nombre de pila como una categoría que puede estar formada por dos posibilidades: una de ellas es un nombre o una inicial más un punto obligatorio; la otra posibilidad es compuesta, y puede estar formada por un nombre seguido de lo que habíamos ya utilizado como nombre de pila; es decir, el concepto se vuelve recursivo. Esto es, puedo usar un nombre más otro nombre, o un nombre más una inicial, o una inicial más un nombre, et cétera. Dicho de manera práctica, un nombre de pila puede ser *José* o *J.* o *José Luis* o *J. Luis* o *José L.* o *J. L.* . Al usar el concepto de recursividad una entidad se define como un nombre seguido de lo que se había utilizado como nombre con anterioridad; es decir, se pueden ir acumulando. Todos los nombres mencionados en el ejemplo caben dentro de esta definición. El símbolo <EOL> indica que lo del primer renglón es aparte de lo del segundo y no deben mezclarse.
- ❖ Hemos definido también que la dirección está formada por el nombre de una calle, seguido de un número exterior y, opcionalmente, un número interior, seguido de un nombre de colonia, en ese orden.
- ❖ Hemos definido que el código postal está formado por el nombre de una ciudad, seguido por el nombre de una entidad federativa y el número de un código postal, en ese orden.

- ❖ Hemos definido que un código postal está formado por cinco dígitos consecutivos.
- ❖ Hemos definido a un dígito como cualquiera de los guarismos del cero al 9.

Si bien no hemos definido una serie de validaciones y restricciones con respecto a estos posibles datos, ya que se trata de un ejemplo muy sencillo como puede observarse, la definición de un entellamado de dirección postal es posible, preciso y poco ambiguo cuando usamos un metalenguaje; en este caso la notación BNF. Esto es más fácil que tratar de explicarlo en lenguaje corriente, ya que habría que estar haciendo muchas explicaciones sobre cada elemento, y además el resultado es una definición que tiene mayor formalidad. Además estas definiciones pueden ser interpretadas por un programa de computador. Puede hacerse en tonces que ese programa es criba las categorías en sitios precisos dentro de una pantalla, hoja o etiqueta de papel: primero el nombre, luego la dirección y luego el código. Puede hacerse un listado sólo de los nombres de las personas y omitir el de más; puede hacerse un listado que ordene los nombres alfabéticamente por apellidos, etcétera. Todas estas acciones serían muy difíciles de efectuar si el nombre y la dirección de una persona fue en una sola entidad escribida corrido, con lo cual una máquina se vería imposibilitada de distinguir cada una de esas partes.

Cabe hacer notar que esta notación de Backus-Naur es sólo un ejemplo de un metalenguaje sencillo para ilustrar conceptos básicos sobre este tópico. No se usa en la práctica para describir textos ya que carece de un buen número de elementos necesarios para realizar un correcto encodificado de ellos mismos. Para ello se han creado metalenguajes propios para ese efecto.

HTML (HYPER-TEXT MARKUP LANGUAGE)

Como ya se ha mencionado, a principios de los noventa se creó el *World Wide Web*, *WWW*, o simplemente *Web*, como un mecanismo para editar y acceder a información a nivel mundial con base en las tecnologías de Internet o red global de computadores. En un inicio la información que viajaba en el Internet era sólo de tipo texto, y las instrucciones debían ser tecladas a modo de comandos. El advenimiento de la Web le agrega interfaces gráficas a la distribución de información; es decir que una página de la web contiene ahora además de textos, imágenes (fotografías, dibujos, diagramas, etcétera), así como otros elementos de acceso más cómodos para el usuario: botones para seleccionar opciones, menús, etcétera. Agrega también un elemento muy importante: el hipervínculo; es decir una “*liga*” con otro texto, que puede estar en esa misma página o en cualquier otra, incluso en otra computadora, sistema o hasta otro país. Gracias al hipervínculo, el usuario puede, al posicionar el sobre un punto en la pantalla, dirigir se hacia otro conjunto de información de su interés sobre la Web. Posteriormente se agregaron

al medio sonidos e imágenes en movimiento, como animaciones o video, para hacer de la edición en la Web un ambiente realmente multimedia.

Precisamente para poder efectuar esta edición de información por medio de la Web se desarrolla el metalenguaje HTML que como ya se ha visto nace como un subconjunto o aplicación muy específica de SGML, del cual sólo se desearían algunas de sus características que permitieran describir las partes de una página pensando sobre todo en cómo iban a verse sobre la pantalla. Lo principal era la distribución física de la información sobre la pantalla, es decir el *lay out*: dónde deben quedar los encabezados o los títulos, su tipo de letra, tamaño y colores; el diseño y colores de los fondos; ubicación, color, tamaño, etcétera de párrafos sobre la pantalla; la ubicación, tamaño, etcétera, de las imágenes sobre la pantalla; los botones para seleccionar alguna opción de un menú; las ligas con otras páginas, etcétera.

Como puede verse HTML podía definir un cierto número de entidades en una página a la cual podía arribar un usuario, y brindarle acceso a la información. Pero adolece de un serio defecto: los textos en su interior son manejados solamente como eso, como “textos”; es decir, como un conjunto de palabras contiguas, sin ningún discernimiento de su contexto, su significado, su ponderación, sus relaciones, etcétera. Por esta razón al multiplicarse el número de páginas en la red se hizo necesario un mecanismo que nos permitiera saber en dónde podría aparecer algo sobre un tópico de nuestro interés. Nacen así los “buscadores”, con el fin de poder buscar y encontrar algo en la red. El problema es que los textos sólo es tan marcados como textos, cada una de las palabras sin ninguna ponderación sobre su contexto o sus partes. Las primeras búsquedas sólo hacen un barrido basándose en la “fuerza bruta” de la computadora que trata de identificar las palabras con tenidas dentro de esos textos. Todos los que hemos navegado y utilizado la red sabemos el resultado de estos procesos: docenas, cientos y hasta miles de referencias inútiles llenas de “ruido” e irrelevantes para nuestro propósito. Quienes conocen algo de catalogación bibliográfica saben que es muy diferente buscar por ejemplo, a Octavio Paz como autor, como título o como tema en una obra. Peor aún si sólo buscamos por “paz” fuera de todo contexto, pues encontraremos a otras personas apellidadas o llamadas “Paz”, pero también noticias sobre la paz, o la falta de ella, documentos, tratados, estudios, etcétera, que nos hablan de este estado social que nada tienen que ver con nuestro Premio Nobel.

Peor aún fue cuando las páginas comenzaron a ser pueras de entrada para enormes repositorios de información, digamos un catálogo de biblioteca; cómo encontrar la información que está dentro de esa base de datos, pero que físicamente no está en la página Web? Las palabras para recuperar documentos usadas como llaves están asociadas a la página en sí misma, no a los bancos de datos ligados a ellas. Esa información está perdida para la recuperación.

Los buscadores tuvieron que irse sofisticando para poder contender con este problema y hoy varios mecanismos se han introducido con este fin, como los “meta datos”, los “buscadores inteligentes” y algunos otros aditamentos como: frecuencia de la palabra en una página, proximidad o contigüidad entre ellas, etcétera. Bastante se

ha avanzado en este sentido, y debemos reconocer que en la actualidad hay páginas con mejores índices y buscadores realmente más adecuados para encontrar información **pertinente**. Pero el problema viene de origen, desde la manera en que la información es preparada y descrita en HTML. Este formato tiene algunas otras deficiencias, entre las cuales las más relevantes son que HTML:

- ❖ No es extensible; es decir, el conjunto de etiquetas es cerrado y nadie puede definir sus propias etiquetas para requerimientos específicos. Por ejemplo, en el mundo bibliográfico querríamos tener etiquetas tales como <autor>, <título> o <ISBN> en lugar de la etiqueta genérica <p> de párrafo.
- ❖ No permite representar las especificaciones de las estructuras de los datos, como se requiere en la creación y uso de bases de datos.
- ❖ No provee soporte para validar los datos.

HTML, por lo tanto, **no maneja los aspectos de contenido** y no puede ser utilizado, por ejemplo, para calificar el contenido de páginas Web. HTML sólo se enfoca en la presentación de la información. Más allá de esto cualquier cosa que requiere de un considerable esfuerzo de programación adicional, ya sea en forma de *applets* o de aplicaciones o programas.

No obstante lo anterior, sería un grave error menospreciar el HTML. Este estándar ha logrado resultados inusitados en la expansión de la divulgación documental electrónica en los últimos años. Gracias a él la Web ha llegado a ser lo que es hoy en día. HTML reforzó el hecho de que la plataforma de cómputo fue irrelevante para el intercambio de información, y creó las bases del mecanismo de transporte para mover documentos a lo largo de distintas redes (lo que llegó a ser el protocolo HTTP, *Hyper-Text Transport Protocol*), y también el esquema de direccionamiento de documentos tanto locales como remotos (lo que llegaría a ser el direccionamiento URL *Universal Resource Locator*). Basante se hizo con HTML a pesar de su belleza y sus limitaciones, simplemente resulta que las necesidades de manejo documental siguen creciendo y empieza a verse limitado para contener con ellas. Como muchas otras grandes herramientas el HTML está siendo rebasado por el mismo dinamismo de la información de la red. El consorcio W3C ha considerado otras opciones que puedan contender con esas demandas y es por ello que SGML es vuelto a tomar en cuenta y nacieron nuevos formatos como XML o XHTML. Sin embargo por un buen tiempo seguramente seguirán viendo muchas aplicaciones desarrolladas en HTML. La aparición de XML y su uso seguramente se implantarán por medio de un cambio gradual de los desarrollos actuales y no mediante una sustitución inmediata.

Si bien los esfuerzos del grupo de trabajo del W3C se concentraron durante un tiempo en la definición del XML, como lo hemos ya mencionado, recientemente se liberaron las especificaciones de una primera versión de un metalenguaje denominado XHTML (eXtensible Hyper-Text Markup Language) compatible con XML, cuyas especificaciones han sido liberadas con anterioridad por el mismo consorcio [XHTML

1.0]. En lo personal no me queda muy claro por qué existen dos sucesos de HTML provenientes de la misma fuente. Como ya se ha explicado, HTML se veía cercado por una serie de limitaciones que eran resueltas con base en aplicaciones adicionales, y por ello un grupo de trabajo desarrolla un nuevo metalenguaje que, con tener el anterior (HTML), pueda con tener con las necesidades actuales, y por esto se desarrolla y libera así la primera versión de XML; hasta ahí todo suena muy lógico. Por qué el mismo con sorcio libera después una nueva versión de TML y la llama XTML y es compatible con ML no está claro y crea mucho mayor confusión en los desarrolladores al no saber si sus nuevas aplicaciones provenientes de HTML deberán hacerse en XHTML o en ML. Pienso que el W3C quiso condescender con diversas facciones de grupos de desarrolladores y aceptó esta nueva versión del HTML convertida a XML con objeto de dar gusto a todos, pero me parece que la versión con más futuro será a la larga la denominada XML.

SGML (Standard Generalized Markup Language ISO 8879:1986)

Como ya se ha comentado, el SGML sirvió de base al HTML. El primero, mucho más completo que su sucesor, además de los aspectos propios de la presentación del documento puede contener una o más definiciones de tipos de documentos (DTD, o “*Document Type Definition*”), las cuales son descripciones formales de la sintaxis de los documentos que se les asignan a cada tipo de ellos que se desea de finir. Por lo tanto, se requiere de una definición DTD para poder interpretar y verificar un documento SGML.

SGML tiene entre sus principales características ventajas; primero, el ser un estándar no propietario, es decir, que no está atado a una marca, patente o compañía en particular, y que es apoyado por un gran número de proveedores de *software*. Por ello, un documento que cumpla con los estándares SGML tendrá una vida más larga que uno basado en un estándar propietario. En segundo lugar, los documentos codificados bajo SGML aún son bastante legibles para las personas y por tanto a la vez legibles para los programas de computadora. En tercer lugar, los documentos bajo SGML describen la estructura de los datos y su semántica, y no sólo la manera en que van a ser presentados en la pantalla.

Entre sus desventajas podemos mencionar: primero, el hecho de que es muy generalizado, que incluye especificaciones particulares para todo tipo de documentos y que se vuelve muy complejo: sus especificaciones se extienden por más de 500 páginas. En segundo lugar, muchas de esas especificaciones son irrelevantes para el uso del documento en la Web, lo cual las hace superfluas cuando éste es el uso que se pretende darle al documento. Dado que hay muchas opciones dentro de éste, la interoperabilidad entre diversas empresas se reduce sensiblemente.

Por lo anterior, los principales usos de SGML se han dado en ambientes “cerrados” donde es más o menos fácil controlar la estandarización de documentos, tales

como los sectores militares, de inteligencia, de manufactura de aeronaves, empresas de publicación, grandes sistemas de archivo, etcétera. [SGML ISO 8879, 1986].

XML (eXtensible Markup Language)

En los anteriores capítulos hemos visto ya que las ventajas y desventajas de HTML y SGML se contraponen: HTML es práctico y esbelto, pero ha sido rebasado por las necesidades del momento en la Web, y GML, siendo mucho más completo tiene muchos elementos imprácticos para su uso en la red. Hemos visto también que el W3C, tratando de contener con esta problemática, favoreció el desarrollo de un nuevo estándar que tratase de tomar “lo mejor de ambos mundos” y así como nace el nuevo estándar XML, que pretende combinar las mayores capacidades descriptivas y semánticas del SGML con lo práctico, sencillo y universal del HTML. Podría afirmarse que se trata de **“la versión corregida y aumentada de HTML” o “un SGML aerodinámico para la Web”**.

En efecto, XML presenta una serie de ventajas prácticas que lo hacen ideal para ir sustituyendo al HTML en la web, he las aquí:

- ❖ Sus especificaciones vienen contenidas en sólo 26 páginas.
- ❖ Los programas analizadores sintácticos de XML, que se construyen en una computadora dada, no requieren del mencionado descriptor del tipo de documento (DTD) para poder separar un documento en sus componentes.
- ❖ No permite ninguna desviación de la sintaxis estándar. Esto quiere decir que todos los documentos codificados en XML pueden ser editados, guardados y enviados sin importar el programa analizador de XML del receptor.
- ❖ Los documentos XML pueden proveer *hojas de estilo* que le permiten a los navegadores (Netscape, Explorer, etcétera) convertir los documentos para presentarlos en HTML en la pantalla.
- ❖ Los analizadores XML pueden analizar incluso documentos HTML bien formados. Esto permite que la transición de documentos ya existentes en HTML en una empresa que esté evolucionando hacia XML pueda hacerse de manera gradual y programada.

Como ya se ha establecido, XML es un metalenguaje que permite establecer un conjunto de reglas para definir una sintaxis específica, la cual será aplicada solamente en un escenario preestablecido. Los documentos que pueden especificarse entonces bajo XML pueden ser de muy diversas naturalezas, por ejemplo las propias del mundo de las finanzas, al especificar datos de acciones, cotizaciones, tipos de cambio de divisas, etcétera; **los datos de tipo bibliográfico contenidos en un documento dado**; o bien datos de compuestos químicos dentro de una fórmula; o planos y datos sobre la construcción de un avión, etcétera.

Se han presentado una serie de postulados acerca de las principales diferencias entre el estándar HTML y el XML. Tal vez la mejor manera de comprender esto plenamente

sea con un pequeño ejemplo que lo ilustra. Tomemos un mínimo trozo del guión dramático *Don Juan Tenorio*, tal como aparece el texto impreso en un libro (sin lenguaje de marcado electrónico):

ACTO I

Escena 1. Don Juan, con antifaz, sentado a una mesa escribiendo; Ciutti y Buttarelli a un lado esperan do. Al levantarse el telón se ven pasar por la puerta del fondo máscaras, estudiantess y pueblo con hachones, música, etcétera.

Don Juan: ¡Cuán gritan esos malditos!

¡Pero mal rayo me parta
Si en concluyendolacarta
No pagan caro sus gritos!

Buttarelli: A Ciutti

Buen carnaval

Ciutti: A Buttarelli

Buen agosto
Para rellenarlaarquilla.

Para leer e interpretar este texto de Zorrilla, una persona no requiere marcado adicional. Es capaz de reconocer el significado de cada una de las partes que lo forman basándose en su experiencia y el simple “*marcado*” tipográfico.

La versión HTML de este texto se vería así:

```
<H1> ACTO I </H1>
<P><I>Escena 1. Don Juan, con antifaz, sentado a una mesa escribiendo; Ciutti y Buttarelli a un lado esperan do. Al levantarse el telón se ven pasar por la puerta del fondo máscaras, estudiantess y pueblo con hachones, música, etcétera.</I></P>
    <P><B>Don Juan:</B> ¡Cuán gritan esos malditos!
    ¡Pero mal rayo me parta
    Si en concluyendolacarta
    No pagan caro sus gritos! </P>
<P><B>Buttarelli:</B>    <I>A Ciutti</I>
    Buen carnaval </P>
<P><B>Ciutti:</B>    <I>A Buttarelli</I>
    Buen agosto
    Para rellenarlaarquilla.</P>
```

La nomenclatura de marcado HTML es la siguiente:

- <> encierra el principio de una etiqueta de marcado.
- </> encierra el fin de una etiqueta de marcado. Siempre van por pares: <> comienza marcado; </> termina marcado para esa etiqueta.

- H significa "Header" o encabezado. Puede haber más de uno y por tanto se numeran: H1 es el "Header 1". <H1> marca de comienzo del encabezado 1. </H1> marca del fin del encabezado 1.
- P significa "Paragraph" o párrafo. <P> y </P> indican el principio y el fin de un párrafo.
- I significa "Itálicas". "B" significa "Bold", negritas. Puestos en pares para indicar dónde comienza y termina esa característica del texto, con ayuda de la diagonal "/".

La versión XML de este mismo texto podría ser:

```
<ACT> <TITLE>ACTO I</TITLE>
<SCENE><TITLE>Escena 1. Don Juan, con antifaz, se ta do a una mesa es cri bien do; Ciutti y But ta re lli a un lado es pe ran do. Al le van tar se el tel ón se ven pa sar por la puerta del fondo máscaras,estudiantes y pueblo con hachones, música,etcé tera.</TITLE>
<SPEECH>
<SPEAKER>      Don Juan:</SPEAKER>
<LINE>
           ¡Cuán gritan esos malditos!
           ¡Pero mal rayo me parta
           Si en concluyendo la carta
           No pa gan caro sus gri tos!
</LINE>
</SPEECH>
<SPEECH>
<SPEAKER>Buttarelli:</SPEAKER>
           <LISTENER>A Ciutti</LISTENER>
           <LINE>Buen car na na val </LINE>
</SPEECH>
<SPEECH>
<SPEAKER>Ciutti:</SPEAKER>
           <LISTENER>A Buttarelli</LISTENER>
           <LINE>Buen agos to
           para rellenarla arquilla.</LINE>
</SPEECH>
.....
```

Las terminaciones de los marcadores </SCENE> y </ACT> estarán has ta donde se encuentre el final de la escena y el final del acto, por tanto no se escriben aquí.

Los marcadores <> y </> (abrir y cerrar etiqueta) se manejan igual que en HTML y por pares.

Si se observan estos ejemplos marcados se puede notar que en el primer caso, con HTML, los marcados nos indican encabezados, párrafos, itálicas y negritas; es decir, **indicaciones para el despliegue del texto, pero nada sobre las partes o el contexto**

que van componiendo la obra. En el segundo ejemplo, el de XML, puede observarse que se ha marcado, además del encabezado o título, dónde empieza el acto I; dónde empieza la escena; dónde empieza y termina el título; dónde comienza y termina cada diálogo (*speech*); quién es el actor del diálogo y quién es el que escucha, y cuál es la línea del diálogo que le pertenece a cada quien. De esta forma, el documento ha quedado marcado con una serie mucho más completa de elementos, la cual es susceptible de ser analizada no sólo por los seres humanos sino también por una computadora.

Supongamos que queremos que la computadora reproduzca por medio de un programa sintetizador de voz los diálogos a través de las bocinas de la misma; por supuesto con voces apropiadas para cada personaje. Con el primer marcado, el HTML, sería muy difícil lograrlo; con el segundo, la máquina puede, al ir “leyendo” el texto, identificar fácilmente quién es el actor que habla <SPEAKER>, y utilizar siempre la misma voz para ese personaje, y reproducir sólo el texto que se encuentra en <LINE>, ignorando para el parlamento las indicaciones de escena, título, nombre del actor, quiénes u interlocutor, etcétera.

A propósito hemos denominado las etiquetas en inglés: *title, act, speech, listener, line*, etcétera. Podríamos haberlas denominado título, acto, diálogo, escucha, línea, etcétera. Sería válido y funcionaría en nuestro ámbito. Pero ¿Qué pasaría si ese diálogo quiere ser pues to en la Internet para su uso universal? ¿Cómo sabrán las máquinas de otras partes del mundo que <diálogo> es eso, un diálogo? Si bien no hay un estándar establecido en cómo debe llamarse a una etiqueta, es muy probable que en poco tiempo, los dramaturgos que quieran poner sus textos en esta forma se inclinen por un nombre universal para la marca, y lo más probable es que sea un nombre en inglés y que, por convención, todo el mundo acepte en poco tiempo que “<speech>”, por ejemplo, es siempre un diálogo en una obra de teatro.

Otro ejemplo de aplicación de este marcado es que una computadora, bajo pedido, puede hacer fácilmente una lista de los personajes de la obra, <SPEAKER> y ponerla en su directo rio de palabras para recuperación, con objeto de que un usuario en la red bus que y en cuenta a tal o cual personaje que aparece en una obra. Ha cer esto con el primer marcado sería sumamente difícil.

Cabe hacer no tara qué en ambos ejemplos no hemos utilizado todos los elementos posibles del marcado de cada metalenguaje. Ni con mucho hemos agotado todas las etiquetas utilizables. Por supuesto, en XML pueden agregarse marcas para tipografía, pero nosotros sólo hemos seleccionado unas cuantas, simplificando el ejemplo, para ilustrar el concepto. No es el propósito de este documento presentar todos los elementos de marcado de cada forma to, por que esto sería inmen so, pero sí existen textos que hacen esto. La idea es resaltar las características de critas a lo largo de este texto.

Resumiendo, podríamos afirmar que XML es:

Simple - La especificación completa mide menos de 30 cuartillas. XML ha sido diseñado para facilitar aún más la escritura de programas con respecto a HTML o SGML.

- Extensible** - Cada quien puede inventar sus propias etiquetas para marcar cualquier tipo de documento, y ser éstas compartidas. De hecho, XML es un metalinguaje que le permite al usuario definir su propio lenguaje de marcado.
- Un estándar abierto** - XML es SGML. Ello significa que no es necesario saber programar; existen muchas herramientas eficientes que permiten ya sea crearlo, manejarlo o implantarlo en una computadora y distribuirlo.
- Eficiente** - XML tiene entes interconstruidos para reutilizar fragmentos de documentos, así, estos sólo tienen que ser transmitidos una vez.
- Basado en la experiencia** - XML ha sido diseñado por personas que tienen amplia experiencia en los lenguajes de marcado y han capitalizado las enseñanzas que el uso de ellos les ha suministrado a lo largo de los años.
- Consensuado** - El diseño de XML incluye los puntos de vista de los organismos coordinadores de HTML y SGML, así como los de personas que han desarrollado importantes aplicaciones con estos estándares.
- Libre** - Nadie tiene la propiedad o patente de XML, ni podrá tenerla, ya que tanto SGML como XML han sido definidos como estándares internacionales. Por lo mismo su uso o desarrollo no implica el pago de ninguna regalía.
- Internacional** - XML tiene interconstruido un soporte para textos en prácticamente todos los alfabetos del mundo, incluyendo técnicas para consignar el lenguaje y/o código del alfabeto utilizado.
- Listo para ser usado** - Los "browsers" o navegadores del web, en sus últimas versiones, son capaces de leer especificaciones XML. Los hipervínculos, textos y multimedia pueden ser vistos tal como se hace ahora con un documento HTML.
- Manejable** - XML incluye métodos para declarar y reforzar las estructuras documentales usadas actualmente, como las de bases de datos.
- Validable** - XML tiene técnicas que permiten la validación de los documentos involucrados, así que uno puede estar seguro de que los documentos registrados con él son creados correctamente.

ANÁLISIS DE LA PROBLEMÁTICA

Para poder proceder a un análisis crítico de la situación es necesario primero acotar el entorno en el cual deseamos realizarlo.

Sin duda una de las características del establecimiento de colecciones y bibliotecas digitales es que éstas no se limitan sólo a materiales de referencia, sino, cada vez más, a textos completos. Las bibliotecas digitales tratan ahora de establecer los mejores mecanismos para almacenar estos textos completos electrónicos de tal forma que sean totalmente recuperables y explotables por las comunidades académicas, más allá de hacer una simple búsqueda de palabras en el texto y operadores booleanos. Ha quedado demostrado que XML tiene un gran número de ventajas nativas para este

propósito, y que por ello mismo será el lenguaje de marcado más utilizado en esta década hasta que sea sustituido por algo mejor. No obstante, su misma universalidad y extensibilidad, y las mayores ventajas de XML, introducen una serie de problemas en el entorno en que éste puede ser utilizado.

En efecto, debemos empezar a atacar el problema en nuestro medio; como se ha mencionado XML puede ser utilizado para describir toda clase de documentos: desde los rollos del Mar Muerto, pasando por un expediente médico o un libro, hasta las especificaciones de un avión supersónico. Puede ser usado entonces en los sectores educativo, de la construcción, médico, gubernamental, etcétera. Cada sector puede deberle empezar a sentar las bases para definir sus documentos más relevantes. De hecho, puede existir más de una definición para documentos de un mismo sector o tipo, lo que haría muy difícil su manejo por parte de usuarios distintos a aquellos para los cuales se haya diseñado la información.

Por este motivo en la actualidad se están definiendo descripciones de tipo de documento (DTDs) por grupos sectoriales con intereses afines, de forma que están surgiendo estándares avalados por organismos que garantizan que cualquier usuario que los adopte como suyos, trabaje con las mismas etiquetas e idénticas normativas, como se hace con el actual HTML. Como ejemplos de esto tenemos CML, *Chemical Markup Language* para el sector químico, MathML, *Mathematical Markup Language* para definir datos matemáticos, SMIL, *Synchronized Multimedia Integration Language*, para definir presentaciones en recursos multimedia, etcétera. A estos tipos de variantes de documentos se les conoce como *clases* y es un término muy importante.

De hecho para un mismo tipo de documento puede existir diferencias en sus necesidades. Tomemos el caso del documento “libro”. La definición de este documento obedecerá a los intereses particulares de determinado sector. Por ejemplo, para alguien que comercia con libros en la Web, su definición de documento sólo incluiría elementos mínimos de identificación: autor, título, año, editorial, ISBN; pero también datos tales como precio, descuento, peso, costo de envío, disponibilidad y/o tiempo de entrega, condiciones de pago, etcétera. Poco que ver con el ambiente de una biblioteca, en donde faltan muchos otros datos con fines de referencia de los documentos, mientras que algunos otros de nuestro ejemplo salen sobrando. Si el documento por definir se es un texto completo, deben agregarse mucho más datos. Por ello las definiciones hechas para un sector serían inútiles para otro.

En el medio académico, y en particular el de las bibliotecas, es necesario entonces comenzar a establecer las definiciones tipo para cada clase de documento, con objeto de satisfacer cabalmente las necesidades de consulta por parte de sus comunidades de usuarios. Cabe recordar que existen además distintos tipos de bibliotecas, y por ende, deberá haber ciertos “matices” en estas definiciones. Estos estándares deberían, por supuesto, estar acordes con nuestra realidad mexicana y con las características de nuestra producción editorial, entorno histórico y cultural, idioma, etcétera.

Debemos reflexionar entonces si queremos sentarnos a esperar que estas definiciones sean hechas por otros en entornos extranjeros, o por sectores ajenos a la academia y a

la biblioteca, o por personal no profesional en el registro documental; o si es el medio bibliotecario académico el que debe comenzar a sentar las bases para establecer estándares en estas definiciones, como ya lo han comenzado a hacer otros sectores en otras partes del mundo. **El sector bibliotecario es el que debería hacer las especificaciones de los documentos para el medio académico**, y sería ideal que fuera el sector bibliotecario mexicano el que lo hiciera para el medio académico mexicano.

Hemos visto ya el despertar y el auge que están tomando las bibliotecas digitales en todo el orbe, y el que comienza a tener en nuestro país. Empezamos a ver ya algunas colecciones de documentos digitales. Sin embargo, cabe observar que casi todas estas colecciones están presentadas en formato HTML y que casi ninguna de estas colecciones de bibliotecas digitales tiene ya sus definiciones XML, aun que hay que destacar que algunas ya lo están considerando y preparándose para ello. Se venden o distribuyen ya en el medio algunos programas o aplicaciones que se anuncian como “generadores de bibliotecas digitales”, tanto de origen nacional como extranjero. Empieza a surgir ya las primeras bibliotecas digitales con especificaciones XML. La pregunta crucial es ¿quién crea esas definiciones para los tipos de documentos que manejan esas bibliotecas? Un examen más de talleado nos enseña que por lo general no han sido hechas por personal profesional en registro documental, y por lo tanto cuando mucho representan el punto de vista de un sólo profesional o sector. Las “etiquetas” de marcado muchas veces son sólo réplicas de las establecidas por MARC o AACR2, que son buenas para material de referencia pero no para textos completos.

Cabe resaltar aquí que estas deficiencias no se hacen no tarde un modo peyorativo alguno. Las personas u organizaciones que están detrás de ello han hecho su mejor esfuerzo con los elementos disponibles hasta este momento y sin ningún estándar preestablecido y por territorios no explorados. Simplemente deseamos resaltar el hecho y el riesgo que se corre de crear una *Babel* al rededor del tema, hasta que los mejores estándares desarrollados se vayan filtrando en el medio y la situación se estabilice. Pero ello puede tomar mucho tiempo y costar caro, además de consumir recursos que en nuestro país no sobran y deben ser utilizados racionalmente para el mejor desarrollo de nuestras bibliotecas y colecciones digitales y, por ende, de las comunidades académicas a las que éstas sirven.

Es altamente necesario, entonces, comenzar a establecer los estándares mínimos que estas definiciones deben tener para cada tipo de documento orientado hacia el sector académico en general y bibliotecario en particular, considerando no sólo las características de referencia sino los textos completos, y enfocándolo a nuestro medio mexicano, a nuestras bibliotecas y a nuestro entorno histórico, social, económico, etcétera. Estas definiciones deben ser efectuadas por personal experto, profesional y multidisciplinario en tareas de registro, recuperación y distribución documental, con el fin de que sean avaladas por las bibliotecas digitales que se están creando y adoptadas por ellas como suyas. La idea es crear un ambiente de homogeneización y calidad en los registros y las colecciones que se formen.

No se trata de ponerle una camisa de fuerza a las especificaciones de los documentos, sino de crear una definición de calidad sobresaliente como base que pueda ser aceptada con confianza por las bibliotecas digitales de nuevo medio, para que a partir de la cada biblioteca pueda seguir construyendo sus propias especificaciones.

¿Qué es lo que habría que definir entonces?

ALCANCES DEL ESTUDIO

Típicamente un documento XML tiene dos partes: una de definición del lenguaje que se usa, y el contenido del documento en sí mismo. La definición puede ser escrita usando una DTD (de definición o descripción del tipo de documento), que es una sintaxis para describir en forma de HTML. Esta DTD puede estar inmersa en el propio documento o ser externa a él. Al ente descrito se le llama **objeto**. Esta técnica ha sido usada hasta este momento con gran intensidad y éxito, dada la reciente liberación del estándar XML. De inicio, se antojaría que el proyecto debería tener de desarrollar las especificaciones detalladas mencionadas anteriormente en forma de DTDs que pueden ser distribuidas para lograr el propósito establecido.

Sin embargo en los últimos meses se ha venido desarrollando una forma más eficaz para definir documentos bajo el estándar XML, conocida como **esquemas de tipo de elemento** o simplemente **esquemas**, que definen las características de las clases de objetos. Estos esquemas pueden definirse como un DTD que permite además su propia ampliación mediante un metalenguaje propio creado al efecto, que define a su vez al DTD. Dicho de otra forma, el esquema representa el metadato para un documento XML asociado, y describe qué puede y qué no puede ser incluido en un documento XML [Campbell, 2003]. Con esta técnica pueden desarrollarse vocabularios especiales para definir y documentar clases de objetos. Esta técnica promete mayores y mejores alcances para definir documentos y se perfila como más interesante como proyecto de investigación de frontera. El presente trabajo pretende explicarlo que es un esquema XML y de describir las ventajas más allá de un DTD. [Young, 2002]

Vamos a tratar de ilustrarlo, para empezar ¿cómo se ve un esquema de tipo de elemento? aparentemente no difiere mucho de un descriptor DTD. Lo mejor será ilustrarlo con un pequeño ejemplo. Todas las declaraciones de un esquema están contenidas en un **elemento de esquema** como en el siguiente ejemplo:

ESQUEMA DE TIPO DE ELEMENTO

```
<?XML version='1.0' ?>
<s:schema id='EjemploEsquema'>
<!-- el esquema va aquí. -->
.....
.....
</s:schema>
```

Declaración del Tipo de Elemento

El núcleo de un esquema de dato XML es la declaración de **Tipo de Elemento**, la cual define una clase de un objeto. El atributo **id** tiene la doble característica de identificar la definición y denominar una clase específica. Los descriptors básicos de nomenclatura ya están definidos en XML y están en inglés. A esto se le denomina el dominio **namespaces** en XML y es una tabla ya definida, si bien en constante actualización. No deben confundirse los descriptors básicos de nomenclatura (*string*, *integer*, *max:occur*, etcétera) ya establecidos, con los nombres posibles de las etiquetas. Ésa es precisamente la ventaja de XML sobre HTML o MARC: los descriptors **no** están establecidos y uno puede bautizarlos como uno considere pertinente. Si bien podríamos bautizar a todos y cada uno de los elementos en los descriptors en castellano: autor, título, tema, etcétera, en una sana defensa de nuestra lengua, ello no será conveniente en la práctica, porque esta definición que daría sólo para uso local en los países de habla hispana, y como lo que se desea es que estas definiciones sean de índole internacional, tal vez lo mejor sea hacerlo en inglés. No obstante lo anterior, es posible establecer “alias” de los nombres, y estos alias pueden ser las nomenclaturas en castellano, con lo cual pueden utilizarse ambas nomenclaturas agregando unas cuantas definiciones más y cuidando que existan simultáneamente. En los ejemplos subsecuentes utilizaremos descriptors en español para resaltar la facilidad con que uno puede denominarlos.

Para entender cómo se construye un esquema de elementos en primer lugar cómo se establecen las definiciones de los elementos. Hagamos unos ejercicios sencillos previos de cómo se define en XML. Utilicemos para el ejemplo la definición de un elemento “autor” dentro de un documento “libro”:

```
<ElementType id="autor"/>
```

Dentro del tipo de elemento, podemos incluir un subelemento **descripción** (*description*), para que las personas puedan describir cuál es el contenido de ese elemento. De hecho podemos agregar múltiples subelementos a un elemento con objeto de dejarlo perfectamente definido. [Thompson, 1997]

```
<ElementType id="autor">
```

```
  <description>La persona o institución que escribió la obra.</description>
```

```
</ElementType>
```

Propiedades y Contenidos

Dentro del **ElementType** podemos crear subelementos que definan las características de lo ahí contenido.

```
<ElementType id="autor">
```

```
  <string/>
```

```
</ElementType>
```

```
<ElementType id="libro">
```

```
  <ElementType="#autor" occurs="ONEORMORE"/>
```

```
</ElementType>
```

Este ejemplo de fine dos elementos, “autor” y “libro”. Indica que un elemento “libro” puede tener uno o varios elementos “autor”. El elemento autor puede contener una cadena (*string*) de caracteres. Dentro de un tipo de elemento pueden estar contenidos varios subelementos (*element, group, any, empty, string*, etcétera). Éstos indican cuáles pueden ser incluidos en ese tipo de elemento e incluso su secuencia. Los tipos de elementos tienen también un atributo **ocurre** (*occurs*), el cual puede tener cuatro valores: requerido “*required*”, opcional “*optional*”, cero o más “*zeroormore*”, y uno o más “*oneormore*”.

Siguiendo estas reglas, podemos continuar construyendo esta descripción:

```
<ElementType="#titulo" occurs="OPTIONAL"/>
```

Podemos agregar también unos subelementos en “grupo”; digamos “prólogo” e “introducción”, los cuales pueden venir en el documento libro. Cabe hacer aquí la aclaración de que estamos definiendo un documento para ser ingresado en texto completo, y no sólo los elementos de su ficha catalográfica.

Un registro más elaborado (sin pretender ser exhaustivo) y sólo a manera de ilustración, podría verse de esta forma:

```
<ElementTypeid="libro">
  <ElementType="#autor" occurs="OPTIONAL"/>
  <ElementType="#titulo" occurs="ONEORMORE"/>
  <ElementType="#fecha-de-copyright" occurs="OPTIONAL"/>
  <groupoccurs="OPTIONAL">
    <elementType="#prefacio"/>
    <elementType="#introduccion"/>
  </group>
  <ElementType="#ISBN" occurs="OPTIONAL"/>
  <ElementType="#texto-del-libro" occurs="ONEORMORE"/>
  <ElementType="#para-edades" occurs="OPTIONAL"/>
</ElementType>
```

Un documento marcado acorde con esta estructura podría verse así:

```
<libro>
  <autor> Miguel de Cervantes Saavedra </autor>
  <titulo> El IngeniosoHidalgo Don Quijote de la Mancha
  </titulo>
  <prefacio> Aquí iría un prólogo escrito..... </prefacio>
  <introduccion> Aquí iría una introducción.... </introduccion>
  <isbn> 123-456-789-0 </isbn>
  <texto-del-libro> En un lugar de la Mancha, de cuyo nombre no quisiera acordarme.... (sigue todo el texto del libro)
  </texto-del-libro>
  <para-edades> todo público </para-edades>
</libro>
```


Uno de los aspectos de mayor importancia en la definición del esquema, consiste en que, como ya habíamos mencionado, pueden establecerse “alias” de los descriptores. Un “alias” consiste, como su nombre lo indica, en un segundo elemento de identificación de un nombre de elemento. **Su importancia radica en que estos alias pueden ser usados para tener una equivalencia en castellano de cualquier nombre definido previa o posteriormente.** El descriptor se llama *elementTypeEquivalent*.

```
<elementTypeEquivalent id="libro" type="#book"/>
<elementTypeEquivalent id="autor" type="#author"/>
<elementTypeEquivalent id="titulo" type="#title"/>
<elementTypeEquivalent id="personaje" type="#personaje"/>
```

...etcétera.

De esta forma, dentro de la definición del documento, “libro” y “book” **se vuelven sinónimos y pueden ser usados indistintamente dentro de la definición.** Así, con este tipo de elementos, podemos mantener la internacionalidad de la definición y seguir manejando los términos en nuestra lengua.

Ahora bien, este fue sólo un pequeño ejemplo relacionado con los atributos de las primeras versiones de XML para facilitar su comprensión, pero su uso puede irse expandiendo y sofisticando hasta dimensiones inimaginadas. Hemos definido un elemento “autor” formado de una sola cadena de caracteres, pero ello no tiene por qué limitarse a eso. Podemos definir autor como una serie de elementos, digamos, nombre, apellidos, título o cargo, fechas de nacimiento y muerte, seudónimos, etcétera. Suena más parecido a lo que estamos acostumbrados a ver en las reglas angloamericanas de catalogación AACR2. Igualmente podemos subdividir título y subtítulo, y pie de imprenta en ciudad, editorial y año, etcétera.

Pero la definición no se queda en los elementos de la ficha catalográfica. Como se ha visto podemos describir el texto completo de un libro. Podemos subdividirlo en capítulos, secciones, apartados, etcétera. Podemos señalarlo a una página, citas bibliográficas, palabras a ser referidas en un glosario, direcciones de otras páginas web, etcétera. Podemos marcar personajes, escenas, lugares, diálogos expresados o pensamientos del personaje, etcétera, como ya lo hemos mostrado en otros apartados previos.

Podemos también, agregando algunas técnicas llamadas API (*Application Programming Interface*), marcar diferentes valores para algún elemento; es decir, podemos pre-determinar si ese elemento es una cadena de caracteres, un número entero o fraccionario, e incluso con rangos posibles; una fecha en algún formato dado (ISO proporcióna varios estándares), un booleano, un elemento que debe existir de acuerdo con una tabla o catálogo preestablecidos, etcétera.

CONSTRUYENDO UN PRIMER ESQUEMA

Una vez que hemos hecho una pequeña presentación de un esquema podemos pasar a construir un esquema que me verí de tal como puede ser construido en la vida real. Para poder iniciar esta construcción partiré de un ejemplo real de un descriptor de tipo de documento (DTD) para un libro y de ahí con tinuare con la construcción de un esquema. Los nombres de los atributos de los *tags* o etiquetas es tanto ma dos del XML *namespaces* para fines de uso ya normalizado, de acuerdo con lo especificado por XML. Para iniciar el ejemplo utilizaré la obra del Tenorio ya presentada con anterioridad según el orden de Van der Vlist [Van der Vlist, 2001]; de acuerdo con un DTD ya establecido podría verse así:

```
<?xml version="1.0" encoding="UTF-7,5"?>
<libro>
<titulo> Don Juan Tenorio </titulo>
<autor> José Zorrilla </autor>
<isbn="8420639028">
<serie="El libro de Bolsillo. Literatura">
<personaje>
  <nombre> Don Juan Tenorio </nombre>
  <representa-a> el galán de la obra </representa-a>
  <desde > 1844 </desde >
  <calificacion>
    galán extrovertido calavera y pendenciero
  </calificacion>
</personaje>
<personaje >
  <nombre> Doña Inés de Ulloa </nombre>
  <representa-a > la dama joven de la obra pretendida por Don
  Juan </representa-a>
  <desde> 1844 </desde>
  <calificacion>
    cándida y angelical monja pretendida por Don Juan
  </calificacion>
</personaje>
<personaje>
  <nombre> Marcos Ciutti </nombre>
  <representa-a > El criado de Don Juan </representa-a>
  <desde> 1844 </desde>
  <calificacion> hábil, pícaro, diligente y fiel sirviente de Don Juan
  </calificacion>
</personaje>
</libro >
```

Para escribir un esquema que describa los documentos marcados en esta forma, simplemente seguiremos su estructura y elementos conforme los vamos encontrando; para empezar, abrimos el elemento `xs:schema`:

```
<?xml version="1.0" encoding="UTF-7,5"?>
<xs:schema
  xmlns:xs="http://www.unam.mx/xmlimaginario/XMLSchema">
  .../...
</xs:schema>
```

La dirección `http` indicada es donde debe residir el esquema; en este caso es ficticia. Sirve para abrir nuevos esquemas y también guarda la definición del *namespace* destino y varias opciones que se asignan por omisión, algunas de las cuales se rán explicadas más adelante.

Para cuadrar la etiqueta inicial para el elemento “libro”, definimos un elemento con ese nombre. Este elemento tiene atributos e “hijos” no-texto, por lo que lo consideramos como de tipo complejo *complexType*, dado que los otros *dataTypes* como *simpleType* es tan reservados para tipos de datos que contienen sólo valores predefinidos y no elementos o sub-nodos de atributo. La lista de “hijos” del elemento “libro” es descrita por un elemento de secuencia *sequence*:

```
<xs:element name="libro">
  <xs:complexType>
    <xs:sequence>
      .../...
    </xs:sequence>
  </xs:complexType>
</xs:element>
```

La secuencia es una “guía” que define una secuencia ordenada de subelementos. Veremos otras guías, *choice* and *all* en las siguientes secciones.

Podemos definir ahora los elementos de autor y título como tipos simples, no tienen atributos o hijos no-texto y pueden ser descritos directamente como elementos. El tipo `xs:string` (cadena de caracteres) los define de esa forma:

```
<xs:element name="titulo" type="xs:string"/>
<xs:element name="autor" type="xs:string"/>
```

Ahora establecemos el elemento “personaje”, el cual es de tipo complejo y describe a los diferentes personajes de la obra. Nótese cómo se define su cardinalidad, es decir, el número de veces que puede ocurrir:

```

<xs:element name="personaje" minOccurs="0"
  maxOccurs="unbounded">
  <xs:complexType>
  <xs:sequence>
  .../...
  </xs:sequence>
</xs:complexType>
</xs:element>

```

Especificamos la lista de todos sus “hijos” de la misma forma:

```

<xs:element name="nombre" type="xs:string"/>
<xs:element name="representa-a" type="xs:string"
  minOccurs="0" maxOccurs="unbounded"/>
<xs:element name="desde" type="xs:date"/>
<xs:element name="cualificacion" type="xs:string"/>

```

Terminamos su descripción cerrando el *complexType*, y los elementos *element* y *sequence*. Podemos ahora declarar otros atributos de los elementos del documento:

```

<xs:attribute name="isbn" type="xs:string"/>
<xs:attribute name="serie" type="xs:string"/>

```

Y se cierran los elementos restantes.

Este diseño, conocido en el medio como de tipo *matrix* *sbk* por las muñecas rusas de madera que se anidan una dentro de la otra, si se acerca la estructura del documento muestra. Una de las características clave de tal diseño es la de definir cada elemento y atributo dentro de su contexto permitiendo múltiples ocurrencias del mismo elemento para obtener diferentes definiciones.

Lista completa del primer ejemplo de esquema:

```

<?xml version="1.0" encoding="utf-8"?>
<xs:schema
  xmlns:xs="http://www.unam.mx/xmlimaginario/XMLSchema">
  <xs:element name="libro">
  <xs:complexType>
  <xs:sequence>
  <xs:element name="titulo" type="xs:string"/>
  <xs:element name="autor" type="xs:string"/>
  <xs:element name="personaje" minOccurs="0"
    maxOccurs="unbounded">
  <xs:complexType>
  <xs:sequence>

```

```

<xs:element name="nombre" type="xs:string"/>
<xs:element name="representa-a" type="xs:string"
  minOccurs="0"
  maxOccurs="unbounded"/>
<xs:element name="des de" type="xs:date"/>
<xs:element name="cualificacion"
  type="xs:string"/>
</xs:sequence>
</xs:complexType>
</xs:element>
</xs:sequence>
<xs:attribute name="isbn" type="xs:string"/>
</xs:complexType>
</xs:element>
</xs:schema>

```

Nótese cómo el primer documento descrito con marcadores XML corresponde a un documento, la obra del Tenorio de Zorrilla en particular. Este último molisado corresponde al esquema de ese documento; es decir, con este es que podemos describir todos los documentos u obras que hayan sido marcadas con el estilo del documento presentado. Podemos entonces procesar los documentos descritos en esa forma, pero si anexamos el esquema a los documentos los sistemas pueden saber cómo fueron construidos esos documentos. [Morrison.2002]

DIVIDIENDO EL ESQUEMA

Si bien la estructura presentada anteriormente es muy simple, puede llevar a muchos anidamientos de las definiciones inmersas, haciendo la muy difícil de leer y difícil de mantener cuando los documentos son muy complejos. Tiene también la desventaja de ser muy diferente a una estructura de un DTD, lo cual es un obstáculo tanto para las personas como para los programas que de sean transformar un DTD en un esquema XML.

Un segundo diseño que presento está basado en una “tabla plana” de todos los elementos disponibles en un documento, y para cada uno de ellos, listas de los elementos “hijos” y sus atributos. Este efecto se logra a través del uso de referencias hacia los elementos y definiciones de atributos que requieren estar dentro del ámbito del referencista, y se ven en forma de una tabla plana como ya mencionamos.

```

<?xml version="1.0" encoding="UTF-7,5"?>
<xs:schema xmlns:xs=
  "http://www.unam.mx/xmlimaginario/XMLSchema">
<!--definición de los elementos de tipo simple-->
<xs:element name="autor" type="xs:string"/>

```

```

<xs:element name="titulo" type="xs:string"/>
<xs:element name="nombre" type="xs:string"/>
<xs:element name="representa-a" type="xs:string"/>
<xs:element name="desde" type="xs:date"/>
<xs:element name="cualificacion" type="xs:string"/>
<!-- definición de atributos -->
<xs:attribute name="isbn" type="xs:string"/>
<!-- definición de elementos de tipo complejo -->
<xs:element name="personaje">
  <xs:complexType>
    <xs:sequence>
      <xs:element ref="nombre"/>
      <xs:element ref="representa-a" minOccurs="0"
        maxOccurs="unbounded"/>
      <xs:element ref="desde"/>
      <xs:element ref="cualificacion"/>
      <!-- Los elementos de tipo simple son referenciados usando el atributo
        "ref" -->
      <!-- La definición de la cardinalidad se hace cuando el elemento es re
        ferenciado -->
    </xs:sequence>
  </xs:complexType>
</xs:element>
<xs:element name="libro">
  <xs:complexType>
    <xs:sequence>
      <xs:element ref="titulo"/>
      <xs:element ref="autor"/>
      <xs:element ref="personaje" minOccurs="0"
        maxOccurs="unbounded"/>
    </xs:sequence>
    <xs:attribute ref="isbn"/>
  </xs:complexType>
</xs:element>
</xs:schema>

```

Usar una referencia hacia un elemento o atributo es equivalente de “clonar” a un objeto. El elemento o atributo es definido primero y luego puede ser duplicado en otro lugar de la estructura del documento por el mecanismo de referencia. Los dos elementos o atributos se vuelven de esta forma dos instancias de la misma clase.

Hemos establecido ya que podemos definir elementos y atributos conforme los vamos necesitando (técnica *matriushka*), o crearlos primero y referirlos después (tabla plana). El documento de es que mas XML desarrollado por el consorcio W3C nos proporciona un tercer mecanismo usado para definir *dataTypes* (tipos de datos) ya sea simples o de tipo complejo, y utilizar esos tipos para definir atributos y elementos.

Esto se lo gra dán do les un nom bre a los ele men tos *simpleType* y *complexType* y ubi cán do los fue ra de la de fini ción de ele men tos o atri bu tos. Más aún, po de mos de ri var un tipo de da tos *dataType* ha cia otro de finien do una res tric ción so bre los va lo res de este tipo de da tos.

Por ejemplo, para de fi nir un *dataType* llama do “tipo-de-nombre” el cual es tá for ma do por una ca de na de has ta 32 ca rac te res, lo es ta ble ce ría mos así:

```
<xs:simpleType name="tipo-de-nombre">
  <xs:restriction base="xs:string">
    <xs:maxLength value="32"/>
  </xs:restriction>
</xs:simpleType>
```

El ele men to *simpleType* guar da el nom bre del nue vo *dataType*. El ele men to *restriction* ex pre sa el he cho de que el *dataType* se de ri va del *dataType* “string” del es pa cio de nom bres de W3C (*namespace*), pe ro se le apli ca una res tric ción a sus po si bles va lo res. El atri bu to *maxLength*, tam bién de nomi na do *faceta* en los do cu men tos del con sor cio, es pe ci fi ca cuál es esa res tric ción, con una lon gi tud má xi ma de 32 ca rac te res.

Otra “faceta” po de ro sa es el ele men to *pattern* (pa trón), el cual de fi ne una pro po si ción re gular que de be ser sa tis fe cha. Por ejemplo, si no nos in te re sa de fi nir el ele men to isbn con sus gui ones, po dría mos de fi nir el *dataType* isbn co mo:

```
<xs:simpleType name="tipo-isbn">
  <xs:restriction base="xs:string">
    <xs:pattern value="[0-9]{X}"/>
  </xs:restriction>
</xs:simpleType>
```

Es decir, lo de fi ni mos co mo una ca de na de ca rac te res cuyo pa trón po si ble son los nú me ros 0-9 (ce ro al nue ve) y el “nú me ro” X (diez), ú ni ca men te. Los gui ones que dan por en de ex clui dos. Si no nos in te re sa ra agre gar el gui ón co mo po si ble en el pa trón, si mplemente agre ga mos el va lor del gui ón {-} des pués del va lor del nú me ro X; esto es: “[0-9]{X}{-}”

De fi nir y uti li zar *dataTypes* con nom bres es com pa ra ble a de fi nir una cla se y uti li zar la para crear un ob je to. Un *dataType* es un ente abs trac to que pue de ser uti li za do para de fi nir un ele men to o un atri bu to. El *dataType* reali za en ton ces con un ele men to o un atri bu to la mis ma fun ción que una cla se reali za con un ob je to.

```
<?xml version="1.0" encoding="UTF-7,5"?>
<xs:schema xmlns:xs=
  "http://www.unam.mx/xmlimaginario/XMLSchema">
  <!-- definición de tipos simples -->
  <xs:simpleType name="tipo-nombre">
```



```

<xs:restriction base="xs:string">
  <xs:maxLength value="32"/>
</xs:restriction>
</xs:simpleType>
<xs:simpleType name="tipo-desde">
  <xs:restriction base="xs:date"/>
</xs:simpleType>
<xs:simpleType name="tipo-descr">
  <xs:restriction base="xs:string"/>
</xs:simpleType>
<xs:simpleType name="tipo-isbn">
  <xs:restriction base="xs:string">
    <xs:pattern value="[0-9]{X}"/>
  </xs:restriction>
</xs:simpleType>
<!-- definición de tipos complejos -->
<xs:complexType name="tipo-personaje">
  <xs:sequence>
    <xs:element name="nombre" type="tipo-nombre"/>
    <xs:element name="representa-a" type="tipo-nombre"
      minOccurs="0"
      maxOccurs="unbounded"/>
    <xs:element name="desde" type="tipo-desde"/>
    <xs:element name="cualificacion" type="tipo-descr"/>
  </xs:sequence>
</xs:complexType>
<xs:complexType name="tipo-libro">
  <xs:sequence>
    <xs:element name="titulo" type="tipo-nombre"/>
    <xs:element name="autor" type="tipo-nombre"/>
    <xs:element name="personaje" type="tipo-personaje"
      minOccurs="0"/>
    <!-- la definición del elemento "personaje" se hace usando el tipo
      complejo "tipo-personaje" -->
  </xs:sequence>
    <xs:attribute name="isbn" type="tipo-isbn"
      use="required"/>
  </xs:complexType>
  <!-- Se hace referencia a "tipo-libro" para definir el elemento "libro" -->
  <xs:element name="libro" type="tipo-libro"/>
</xs:schema>

```

Grupos, compositores y grupos derivados

Los esquemas de acuerdo con el W3C también permiten la definición de grupos de elementos y atributos:

```
<!-- definición de un grupo de elementos -->
<xs:group name="elementos-principales-de-libro">
  <xs:sequence>
    <xs:element name="titulo" type="tipo-nombre"/>
    <xs:element name="autor" type="tipo-nombre"/>
  </xs:sequence>
</xs:group>
<!-- definición de un grupo de atributos -->
<xs:attributeGroup name="atributos-de-libro">
  <xs:attribute name="isbn" type="tipo-isbn"
  use="required"/>
  <xs:attribute name="available" type="xs:string"/>
</xs:attributeGroup>
```

Estos grupos pueden ser usados en la definición de tipos complejos, como se muestra:

```
<xs:complexType name="tipo-libro">
  <xs:sequence>
    <xs:group ref="elementos-principales-de-libro"/>
    <xs:element name="personaje" type="tipo-personaje"
      minOccurs="0" maxOccurs="unbounded"/>
  </xs:sequence>
  <xs:attributeGroup ref="atributos-de-libro"/>
</xs:complexType>
```

Estos grupos no son *dataTypes* sino contenedores de un conjunto de elementos o atributos que pueden ser usados para describir tipos complejos.

Compositores

Hasta ahora hemos visto el compositor de secuencia `xs:compositor` que define grupos ordenados de elementos o partículas, los cuales pueden ser a su vez grupos de otros compositores. El XML del W3C soporta además dos compositores adicionales que pueden combinarse para permitir varias mezclas. Cada uno de esos compositores puede tener un mínimo y un máximo de ocurrencias (`minOccurs` y `maxOccurs`) para definir su cardinalidad. *Unbounded* significa sin límite.

El compositor *xs:choice* describe la elección entre varios posibles elementos o grupos de ellos. El siguiente grupo (los compositores pueden aparecer dentro de grupos tipos complejos u otros compositores) acepta bien un simple nombre de elemento o una secuencia nombre-de-pila, un apellido-paterno o un segundo-apellido opcional:

```
<xs:group name="tipo-nombres">
  <xs:choice>
    <xs:element name="nombre" type="xs:string"/>
    <xs:sequence>
      <xs:element name="nombre-de-pila" type="
        xs:string"/>
      <xs:element name="apellido-paterno" type="
        xs:string"/>
      <xs:element name="segundo-apellido" type="xs:string"
        minOccurs="0"/>
    </xs:sequence>
  </xs:choice>
</xs:group>
```

El compositor *xs:all* define un conjunto no-ordenado de elementos. La siguiente definición de tipos complejos permite a sus elementos con tenidos aparecer en cualquier orden:

```
<xs:complexType name="tipo-libro">
  <xs:all>
    <xs:element name="titulo" type="xs:string"/>
    <xs:element name="autor" type="xs:string"/>
    <xs:element name="personaje" type="tipo-personaje"
      minOccurs="0" maxOccurs="unbounded"/>
  </xs:all>
  <xs:attribute name="isbn" type="tipo-isbn"
    use="required"/>
</xs:complexType>
```

Para evitar combinaciones que pudieran resultar ambiguas o demasiado complejas para ser resueltas por las herramientas de *es:XML* de W3C, deben agregarse un conjunto de restricciones a la declaración *xs:all*.

- ❖ Sólo pueden aparecer como un hijo por una vez al principio del paquete
- ❖ Sus hijos pueden ser sólo definiciones *xs:element* o referencias y no pueden tener una cardinalidad mayor que uno.

Derivados de tipos simples

Los *dataTypes* simples pueden ser de finidos por derivados de otros *dataTypes*, bien predefinidos e identificados por el namespace del esquema del W3C o definidos en otro lado en el esquema.

Hemos visto ya ejemplos de tipos simples derivados por restricción (usando *xs:restriction*). Las variedades diferentes de restricciones que pueden ser aplicadas a un *dataType* se llaman facetas, como ya hemos establecido, y son particularmente útiles para establecer la longitud de un elemento, los posibles valores que puede tomar, el número mínimo y máximo de ocurrencias, etcétera.

Existen otros dos métodos de derivación que permiten definir espacios en blanco, listas separadas para ser consultadas, uniones de *dataTypes*, etcétera. La siguiente definición usa *xs:union* para extender nuestro tipo-*isbn* y que acepte los valores TBD y NA:

```
<xs:simpleType name="tipo-isbn">
  <xs:union>
    <xs:simpleType>
      <xs:restriction base="xs:string">
        <xs:pattern value="[0-9]{X}"/>
      </xs:restriction>
    </xs:simpleType>
    <xs:simpleType>
      <xs:restriction base="xs:NMTOKEN">
        <xs:enumeration value="TBD"/>
        <xs:enumeration value="NA"/>
      </xs:restriction>
    </xs:simpleType>
  </xs:union>
</xs:simpleType>
```

La “unión” ha sido aplicada a dos tipos distintos de valores para permitir la existencia de dos conjuntos de valores; el nuevo *dataType* acepta ahora los valores enumerados TBD y NA.

El siguiente ejemplo “tipo-*isbn*” usa *xs:list* para definir una lista de varios ISBN separados por un espacio en blanco. Deriva también un tipo “tipo-*isbn10*” usando *xs:restriction* para aceptar entre 1 y 10 valores de números de ISBN separados por un espacio.

```
<xs:simpleType name="tipoisbn">
  <xs:list itemType="tipo-isbn"/9
</xs:simpleType>
<xs:simpleType name="tipo-isbn10">
  <xs:restriction base="tipoisbn">
    <xs:minLength value="1"/>
    <xs:maxLength value="9"/>
  </xs:restriction>
</xs:simpleType>
```

Mezclas en los contenidos

La recomendación del W3C para eso que nos permite también mezclar partes de documentos marcados con otras secciones sin marcar usando la restricción *xs:complexType*, obsérvese el siguiente texto:

```
<isbn libro="8420639028">
  Obra dramática de José Zorrilla
  Su título, Don Juan tenorio, lo dice todo!
</isbn libro>
```

Usando la restricción *xs:complexType*:

```
<xs:element name="libro">
  <xs:complexType mixed="true">
    <xs:all>
      <xs:element name="titulo" type="xs:string"/>
      <xs:element name="autor" type="xs:string"/>
    </xs:all>
    <xs:attribute name="isbn" type="xs:string"/>
  </xs:complexType>
</xs:element>
```

Lo cual validará un elemento XML como el siguiente:

```
<isbn libro="8420639028">
  Divertida obra por <autor>José Zorrilla</autor>.
  Su título(<titulo> Don Juan tenorio</titulo>) lo dice todo!
</isbn libro>
```

La construcción *xs:keyref* es sumamente útil en este tipo de obras ya que permite hacer ligas entre algunos elementos y otros. Por ejemplo, de sea mos asociar a los distintos elementos “personajes” dentro de esta obra; recordemos la definición establecida en un principio:

```
<personaje>
  <nombre>Don Juan Tenorio</nombre>
  <representa-a>el galán de la obra </representa-a>
  <desde >1844</desde >
  <calificacion>
    galán extrovertido calavera y pendenciero
  </calificacion>
</personaje>
```

Y queremos indicar que hay una asociación entre este personaje y otro u otros dentro de la obra; por ejemplo, para que cuando uno sea con sul ta do, se asocie siempre al otro, el cual debe ser referenciado en el mismo libro:

```
<xs:keyref name="referencia-personaje" refer=
  "nombre-personaje">
  <xs:selector xpath="personaje"/>
  <xs:field xpath="representa-a"/>
</xs:keyref>
```

Nos presentará asociaciones con otros personajes de la obra; por ejemplo:

```
<personaje>
  <nombre> Marcos Ciutti</nombre>
  <representa-a > El criado de Don Juan</representa-a>
  <desde>1844</desde>
  <cualificacion>hábil, pícaro, diligente y fiel sirviente de Don Juan
  </cualificacion>
</personaje>
```

Existen todavía algunas otras definiciones para la completa construcción de un esquema, pero no es el caso de este documento con solidarse como un tutorial de esquemas, sino el haber ilustrado las considerablemente mayores posibilidades que un esquema permite lograr mucho más allá de la simple definición DTD de un documento. [W3Cschools, 2003]

LENGUAJES PARA ESQUEMAS

Por todo lo anteriormente expuesto sobre la estructura de esquemas, es claro que nos permites establecer una enorme cantidad de características para el registro, el control y el contexto de un documento dado, pero la estructura puede llegar a complicarse tanto que la construcción de un esquema paso a paso parecería una enorme tarea, como aquí hemos podido observar. Da la impresión de que debe uno memorizar enormes cantidades de instrucciones, símbolos, variables, estructuras, etcétera, lo cual a primera vista parece tedioso y memorizante. Pero no es así. En la práctica los esquemas no tienden a ser escritos por las personas sino por las máquinas.

Como hemos podido observar, la especificación 1.0 de XML de fin los conceptos para una correcta formabilidad y validación. A través de ella es muy fácil verificar un documento certificando que tenga buena formabilidad; no obstante, la validación requiere de más trabajo y le permite al usuario definir restricciones más potentes para la estructura del documento. La validación XML requiere que el documento

siga las restricciones establecidas en el DTD, el cual provee el equivalente de lo que sería en un contexto libre una sintaxis para un tipo de documento [Eito, 2001].

En algunos usos los sistemas pueden requerir de definiciones de marcado más informativas, o restricciones más estrictas o más flexibles en la estructura del documento, o simplemente diferentes de las expresadas en definiciones DTDs según las define XML 1.0. Hay también una tendencia a permitir marcados y restricciones que pueden ser especificados en una sintaxis totalmente XML, para que de esta forma las diversas herramientas electrónicas para los documentos XML puedan ser usadas en las mismas especificaciones. Esto puede lograrse gracias a la creación de **lenguajes de programación para la creación e interpretación de esquemas de marcado**.

De la información proveniente de la página del consorcio W3C para los requerimientos de un esquema [Malhotra, 1999] podemos diferir que el propósito del lenguaje de esquemas XML es el de proveer un inventario de las construcciones de marcado XML con las cuales se pueden escribir esquemas. A su vez, el propósito de un esquema es de definir una clase de documentos XML y usar esas construcciones para delimitar y documentar el significado, uso y relaciones de sus partes constituyentes: tipos de datos, elementos y sus contenidos, atributos y sus valores, entidades y sus contenidos y notaciones, valores implícitos como por ejemplo los de omisión (*default*), etcétera. Los esquemas documentan su propio significado, uso y función. Por ello, el lenguaje de esquemas XML puede ser utilizado para definir, describir y catalogar vocabularios XML para diversas clases de documentos XML. Todo programa XML puede usar el formalismo inherente a un esquema para expresar delimitaciones de valor, estructurales o sintácticas, aplicables a ciertos documentos.

Algunos de los posibles ambientes, además del de las bibliotecas digitales, que podrían verse beneficiados con el uso de la programación basada en estos esquemas pueden ser: el editorial, donde la distribución de información puede involucrar documentos que guardan complejas relaciones entre sí (en estas relaciones los esquemas estructurales describirían las relaciones entre encabezados, nuevos artículos, referencias cruzadas, fotografías, etcétera); el comercio electrónico; la adquisición masiva de datos; la transferencia de información en bases de datos; el intercambio de metadatos, etcétera.

De acuerdo con el grupo de trabajo de esquemas del consorcio W3C, los principios de diseño para un lenguaje de esquemas se basan en las siguientes características:

El lenguaje para esquemas XML debe:

- ❖ Ser más expresivo y detallado que los DTDs.
- ❖ Poder expresarse totalmente en XML.
- ❖ Ser autodescriptivo.
- ❖ Ser utilizable por una amplia variedad de programas que usen XML.
- ❖ Poder ser ampliamente utilizable en el Internet.
- ❖ Estar optimizado para interoperabilidad.

- ❖ Ser suficientemente simple para poder implementarse con recursos moderados.
- ❖ Estar coordinado con otras especificaciones W3C (ligas, apuntadores, etcétera).

En cuanto a sus requerimientos estructurales, el lenguaje para esquemas XML debe poder definir los mecanismos para:

- ❖ Delimitar la estructura de los documentos (espacios, elementos, atributos) así como de los contenidos (tipos de datos, entidades, notaciones).
- ❖ Habilitar la herencia para definir elementos, atributos y tipos de datos.
- ❖ Habilitar documentación autocontenida dentro de otra.
- ❖ Habilitar descripciones y delimitadores específicos para un programa en especial.
- ❖ Contender con la integración de esquemas estructurales con tipos de datos primitivos.
- ❖ Definir las relaciones entre el esquema y los documentos XML.
- ❖ Definir las relaciones entre esquemas y DTDs.

De lo anterior podemos concluir que en los próximos años habrá una tendencia marcada para construir y desarrollar lenguajes de programación que permitan la creación y explotación de esquemas en una forma mucho más sencilla de la que hemos ido estableciendo aquí en este documento. No parece factible que las personas se dediquen a escribir líneas y líneas con todas esas especificaciones y restricciones. Lo que sí es necesario es que se entiendan los conceptos básicos que permiten definir, restringir, correlacionar, etcétera, a las diferentes entidades de un documento; sería imposible hacerlo de otra forma. Alguien tiene que establecer las características y estructuras de un documento.

Sin embargo la generación sí del código del esquema será facilitada a través de lenguajes para esquemas. Algunas pantallas amigables al usuario nos permitirán transmitir de manera sencilla la idea que nos hemos hecho del esquema en múltiples instrucciones que serán leídas por otros programas para interpretar y procesar.

Lo importante de todo esto, en suma, es aprender a crear la estructura mental de orden y contexto que un documento conlleva, de acuerdo con la filosofía de los esquemas que, como hemos podido observar, permite contextualizar los documentos mucho más allá de lo que hemos estado acostumbrados. Desde mi punto de vista, el aprendizaje de desarrollar esquemas no es triba en la habilidad de escribir códigos, sino en la habilidad de plantearse un documento en la mente bajo una gran estructura organizativa, y creo que eso es lo más relevante del aprendizaje de los esquemas, sobre todo para los bibliotecarios que nos hemos acostumbrado a disectar, catalogar y acomodar los múltiples elementos de una ficha. Me parece que el esquema nos permite aplicar las estructuras de orden y sistema no tan solo a los elementos catalográficos, sino al contenido mismo del documento. De ahí su gran valor como educadores de la mente.

CONCLUSIONES

Como ha podido observarse, el marcado XML, aun en un simple DTD, permite establecer definiciones mucho más complejas de un documento que las que podíamos establecer en una ficha catalográfica o un documento HTML para despliegue. Es claro el potencial que el marcado XML nos permite de desarrollar en el marcado y la recuperación documental, aun a nivel de un DTD.

Pero, como ha podido comprobarse, las estructuras planeadas por XML pueden irse sofisticando más y más. La propuesta de creación de esquemas relacionados con un tipo dado de documento nos ofrece un potencial aún mucho mayor de registro y explotación de la información, sin perder la calidad de registro catalográfico, como estamos acostumbrados, sino extendiéndolo a niveles más profundos del contenido.

Los esquemas pueden viajar en la red junto con los documentos marcados bajo su estructura. Esto quiere decir que el nivel y calidad de nuestro marcado no está restringido a nuestro sistema de cómputo. Podemos entregar e intercambiar documentos con estructuras muy complejas y precisas que vayan llevando a niveles más profundos, profesionales y satisfactorios de almacenamiento y recuperación de información.

Como puede deducirse es importante ir creando proyectos que desarrollen esquemas específicos para distintos tipos de documentos dentro de las bibliotecas mexicanas. Ello tendría un gran valor a futuro para las nacientes bibliotecas digitales de nuestro país, ya que podría ofrecerse un estándar mexicano y, sin embargo, perfectamente internacional, hecho por profesionales, y con el cual pueden esas bibliotecas empesarse “marcar” sus documentos de acuerdo con la sugerencia de cada esquema. En primer lugar, podría resolverse de una vez la discusión teórica sobre si usar DTDs o esquemas; o DTDs y esquemas. Con ese resultado, y dependiendo del tipo de material, cada biblioteca podrá entonces alojar documentos en bibliotecas digitales de tipo “no vela”, “obra de teatro”, “libro de texto”, etcétera. A la larga, el proyecto podría continuar para obtener un conjunto de todos los formatos “típicos” de documentos: libros, revistas, mapas, discos, películas, etcétera, como alguna vez MARC lo hizo, pero con divisiones específicas para cada variedad de documento en cada formato, y con alcances para cada parte del contenido y no sólo para los datos de referencia. Esto aseguraría un buen nivel de marcado, la estandarización en trediversas bibliotecas digitales mexicanas que se vayan creando y, lo más importante, una espléndida y poderosa recuperación y explotación por parte de los usuarios, etcétera. Con ello podrían sentarse las bases para un modelo de desarrollo de esquemas que pudiese ser llevado a nivel de recomendaciones para el medio mexicano y tal vez ser imitado en otros países latinoamericanos.

Si no inicia mos un proyecto de este corte, en un futuro cercano las bibliotecas empezarán a hacer sus marcados XML conforme al diseño de cada una de ellas, y como no existeningún estándar, habrá muy diferentes especificaciones y calidades de marcado, incluso algunas que nada tengan que ver con este sector. Esto no sería bueno para esas colecciones digitales, pues cuando se deseara adoptar algún estándar preestablecido

quizá sería tarde y podrían existir muchos documentos ya introducidos en alguna colección dada, lo cual dificultaría la conversión a nuevos estándares y obstaculizaría la estandarización de marcados y calidades homogéneas para las colecciones electrónicas. Peor aún sería tener que esperar e imitar estándares internacionales, que si bien pueden ser bien hechos nos siguen de jando a la sombra de esos desarrollos hechos en otros países, no del todo acordes con nuestra realidad nacional.

Finalmente, las conclusiones pueden resumirse de acuerdo con lo siguiente:

- ❖ Los lenguajes de marcado realmente no son nuevos; tienen más de cuatro décadas y lo que observamos hoy son sus versiones más sofisticadas.
- ❖ Los metalenguajes usados hoy en día en el mundo electrónico-documental tienen diversas características, ventajas y desventajas dependiendo del contexto de su creación y su momento histórico.
- ❖ La tendencia hoy en día es que los lenguajes de marcado ayuden a consignar no tan solo los aspectos catálogos de una obra, sino también los relacionados con su contenido.
- ❖ XML se constituye como la versión más completa y funcional de los lenguajes de marcado hoy en día.
- ❖ Los esquemas permiten la definición de estructuras y relaciones mucho más completas que las de un DTD.
- ❖ El establecimiento de lenguajes de creación de esquemas facilitará la parte operativa y mecánica del desarrollo de esquemas.
- ❖ Es importante ir creando definiciones, DTD o esquemas, de las principales clases de documentos en el medio bibliotecario mexicano, para garantizar su calidad y profundidad.
- ❖ Es importante entender la filosofía de creación de esquemas como un mecanismo de ejercicio mental de organización y sistematización de la información.

REFERENCIAS BIBLIOGRÁFICAS¹

A general Introduction to the MARC Format. Disponible en:
<http://www.tlcdelivers.com/tlc/crs/Gen0009.htm>

Campbell, Ch. 2003. "XML Schema". En: *Sigmond Record*, Vol. 32, No 2, Junio 2003. pp. 96-101.

Eito, R. 2001. *Programación con XML*. Madrid: Anaya. ISBN: 84-415-1186-1. pp. 487-540.

* Extensible Markup Language (XML) 1.0 (2nd. ed.), W3C Recommendations. Octubre 6, 2002. Disponible en:
<http://www.w3c.org/TR/1998/NOTE-XML-data-0105/>

Guérrez, A y Martínez, R. 2001. *XML a través de ejemplos*. Madrid: RA-MA. ISBN: 84-7897-455-5.

¹ Los URLs marcados con un * corresponden a un sitio oficial o con nivel de estándar internacional.

- Free Dictionary of computing. 2000. En tra da por: “Backus-Naur Form”.
 Disponible en:
<http://burks.bton.ac.uk/burks/foldoc/>
- * Iflanet. *What is MARC*. Disponible en:
<http://www.ifla.org/VI/3/p1996-1/unimarc.htm>
- * Library of Congress. MARC Standards. 2003. Disponible en:
<http://www.loc.gov/marc>
- * Library of Congress. MARC 21 XML Schema. 2003. Disponible en:
<http://www.loc.gov/standards/marcxml>
- Morrison, M. 2002. Teach Yourself XML in 24 hours. In dia na: Sams. ISBN:
 0-672-32213-7. Parte II, hora 5, “Using XML Schema”, pp. 69-95.
- MSDN on line. 2000. XML Developing Center Disponible en:
<http://msdn.microsoft.com/library/default.asp?url=/nhp/Default.asp?contentid=28000438>
- * SGML ISO 8879. 1986. Disponible en:
<http://www.iso.org/iso/en/CatalogueDetailPage.CatalogueDetail?CSNUMBER=16387>
- Should it be an element or a type ? Disponible en:
<http://www.xfront.com/ElementVersusType.html>
- Teasdale, G. 1995. L’Hyper texte: his to ri que et appli ca tions en bi bli o thé -
 conomie. Université de Montréal. Disponible en:
<http://www.infotheque.info/ressource/942.html>
- Timeline of XML Events. Disponible en:
<http://www.ukoln.ac.uk/web-focus/events/workshops/webmaster-2002/materials/savory/slides/img18.html>
- * XHTML1.0:TheeXtensibleHypertextMarkupLanguage.Disponible en:
<http://www.w3.org/TR/xhtml1/#xhtml>
- The XML Cover Pages - Home page. Disponible en:
<http://oasis-open.org/cover/>
- [Malhotra, A. 1999]. *XML Schema Requirements*. Disponible en:
<http://www.w3.org/TR/NOTE-xml-schema-req>
- Van der Vlist, E. 2001. *Using W3C XML Schema*. Disponible en:
<http://www.xml.com/pub/a/2000/11/29/schemas/part1.html>
- Thompson, H. 1997. *Why I demand “Schemata”*. Disponible en:
<http://www.ltg.ed.ac.uk/~ht/sgml97.html>
- * W3C Schools. 2003. XML Tutorial. Disponible en:
<http://www.w3schools.com/xml/default.asp>
- Young, M.J. 2002. XML Step by Step. 2ª ed. Redmond: Microsoft Press.
 ISBN: 0-7356-1465-2. Capítulo 7: “Creating valid XML documents
 using XML Schemas”. pp. 163-192.