



LMSGI

LENGUAJES DE MARCAS

Y

SISTEMAS DE GESTIÓN DE LA INFORMACIÓN

0. Introducción a los lenguajes de marcas

Alejandro Amat Reina

ÍNDICE

1. ORDENADOR E INFORMACIÓN	2
DATOS EN FORMA DE TEXTO Y DATOS BINARIOS	2
DATOS BINARIOS	2
TEXTO	3
EL CÓDIGO ASCII	3
UNICODE	7
ARCHIVOS BINARIOS Y ARCHIVOS DE TEXTO.....	7
VENTAJAS DE LOS ARCHIVOS BINARIOS	7
VENTAJAS DE LOS ARCHIVOS DE TEXTO	8
2. EL INTERCAMBIO DE INFORMACIÓN	8
EL TEXTO COMO EL FORMATO MÁS VERSÁTIL	9
3. LENGUAJES DE MARCAS	9
ORÍGENES DE LOS LENGUAJES DE MARCAS	10
TEX Y LATEX	10
RTF	11
SGML	12
POSTSCRIPT	12
HTML	13
XML.....	13
JSON	14
TIPOS DE LENGUAJES DE MARCAS.....	14
ORGANIZACIONES DESARROLLADORAS DE LENGUAJES DE MARCAS	15

1. Ordenador e información

El ordenador es una máquina digital, por lo tanto sólo es capaz de representar información utilizando el sistema binario de numeración (secuencias de 1 y 0). Esto obliga a que, para poder almacenar información en un ordenador, previamente haya que codificarla en forma de números binarios.

El problema de los números binarios es que están muy alejados del ser humano; es decir, que las personas no estamos capacitadas para manejar información en binario. Nosotros usamos sistema decimal para los números y formas de representación mucho más complejas para otra información (como el texto, las imágenes, la música,...)

Sin embargo, actualmente un ordenador es capaz de manejar información de todo tipo: música, imágenes, texto,.... Esto es posible porque se ha conseguido que casi cualquier tipo de información sea codificable en binario.

Los seres humanos tenemos la capacidad de diferenciar claramente lo que es un texto de una imagen, lo que es un número de una canción,... Pero en un ordenador todo es más complicado, porque todo es binario.

Desde los inicios de la informática, la **codificación** (el paso de información humana a información digital) ha sido problemática debido a la falta de acuerdo en la representación. Pero hoy día ya tenemos numerosos estándares.

Fundamentalmente, la información que un ordenador maneja son **Números** y **Texto**. Pero curiosamente a nivel formal se consideran **datos binarios** a cualquier tipo de información representable en el ordenador, que no es texto (imagen, sonido, vídeo,...), aunque como ya hemos comentado, en realidad toda la información que maneja un ordenador es binaria, incluido el texto.

Datos en forma de texto y datos binarios

Datos binarios

Cualquier dato que no sea texto se considera dato binario. Por ejemplo: música, vídeo, imagen, un archivo Excel, un programa,...

La forma de codificar ese tipo de datos a su forma binaria es muy variable. Por ejemplo, en el caso de las imágenes cada punto (píxel) de la imagen se codifica utilizando su nivel de rojo, verde y azul. De modo que una sola imagen produce millones de dígitos binarios.

En cualquier caso, sea cual sea la información que estamos codificando en binario, para poder acceder a dicha información, el ordenador necesita el software que sepa como decodificar la misma, es decir saber qué significa cada dígito binario para traducirlo a una forma más humana. Eso sólo es posible utilizando el mismo software con el que se codificó o bien otro software, pero que sea capaz de entender la información codificada.

Texto

El texto es quizá la forma más humana de representar información. Antes de la llegada del ordenador, la información se transmitía mediante documentos o libros en papel. Esa forma de transmitir es milenaria y sigue siendo la forma más habitual de transmitir información entre humanos; incluso con la tecnología actual aplicaciones como twitter, WhatsApp,... siguen usando el texto como formato fundamental para transmitir información.

En cuanto apareció la informática como una ciencia digital, apareció también el problema de cómo codificar texto en forma de dígitos binarios para hacerlo representable en el ordenador. La forma habitual ha sido codificar cada carácter en una serie de números binarios, de modo que, por ejemplo el carácter A se codifica como 01000001 y la B como 01000010.

El problema surgió por la falta de estandarización, la letra A se podía codificar de forma distinta en diferentes ordenadores y así nos encontrábamos con un problema al querer pasar datos de un ordenador a otro. Poco a poco aparecieron estándares para intentar que todo el hardware y software codificara los caracteres igual.

El código ASCII

El problema de la codificación de texto que hacía incompatibles los documentos de texto entre diferentes sistemas, se redujo cuando se ideó en 1967 un código estándar por parte de la Agencia de Estándares Norteamericana (ANSI), dicho código es el llamado ASCII (American Standard Code for Information Interchange, código estándar americano para el intercambio de información). El código utiliza el alfabeto inglés (que utiliza caracteres latinos) y para codificar todos los posibles caracteres necesarios para escribir en inglés se ideó un sistema de 7 bits (con 7 bits se pueden representar 128 símbolos, suficientes para todas las letras del

alfabeto inglés, en minúsculas y mayúsculas, caracteres de puntuación, símbolos especiales e incluso símbolos de control).

El código ASCII es el siguiente:

Caracteres ASCII de control			Caracteres ASCII imprimibles					
00	NULL	(carácter nulo)	32	espacio	64	@	96	`
01	SOH	(inicio encabezado)	33	!	65	A	97	a
02	STX	(inicio texto)	34	"	66	B	98	b
03	ETX	(fin de texto)	35	#	67	C	99	c
04	EOT	(fin transmisión)	36	\$	68	D	100	d
05	ENQ	(consulta)	37	%	69	E	101	e
06	ACK	(reconocimiento)	38	&	70	F	102	f
07	BEL	(timbre)	39	'	71	G	103	g
08	BS	(retroceso)	40	(72	H	104	h
09	HT	(tab horizontal)	41)	73	I	105	i
10	LF	(nueva línea)	42	*	74	J	106	j
11	VT	(tab vertical)	43	+	75	K	107	k
12	FF	(nueva página)	44	,	76	L	108	l
13	CR	(retorno de carro)	45	-	77	M	109	m
14	SO	(desplaza afuera)	46	.	78	N	110	n
15	SI	(desplaza adentro)	47	/	79	O	111	o
16	DLE	(esc.vínculo datos)	48	0	80	P	112	p
17	DC1	(control disp. 1)	49	1	81	Q	113	q
18	DC2	(control disp. 2)	50	2	82	R	114	r
19	DC3	(control disp. 3)	51	3	83	S	115	s
20	DC4	(control disp. 4)	52	4	84	T	116	t
21	NAK	(conf. negativa)	53	5	85	U	117	u
22	SYN	(inactividad sinc)	54	6	86	V	118	v
23	ETB	(fin bloque trans)	55	7	87	W	119	w
24	CAN	(cancelar)	56	8	88	X	120	x
25	EM	(fin del medio)	57	9	89	Y	121	y
26	SUB	(sustitución)	58	:	90	Z	122	z
27	ESC	(escape)	59	;	91	[123	{
28	FS	(sep. archivos)	60	<	92	\	124	
29	GS	(sep. grupos)	61	=	93]	125	}
30	RS	(sep. registros)	62	>	94	^	126	~
31	US	(sep. unidades)	63	?	95	_		
127	DEL	(suprimir)						

Pero, en países con lenguas distintas del inglés, surgió el problema de que parte de los símbolos de sus alfabetos quedaban fuera del ASCII (como la letra ñe). Por ello, se diseñaron códigos de 8 bits que añadían 128 símbolos más y así aparecieron los llamados códigos ASCII extendidos. En ellos, los 128 símbolos primeros son los mismos de la tabla ASCII original y los 128 siguientes se corresponden a símbolos extra. Así, por ejemplo, el sistema MS-DOS utilizaba el llamado código 437 que incluía símbolos y caracteres de otras lenguas de Europa Occidental y caracteres que permitían hacer marcos y bordes en pantallas de texto, entre otros símbolos.

ASCII extendido (Página de código 437)					
128	Ç	160	á	192	Ł
129	ü	161	í	193	ł
130	é	162	ó	194	Ť
131	â	163	ú	195	ť
132	ä	164	ñ	196	—
133	à	165	Ñ	197	†
134	â	166	°	198	ã
135	ç	167	º	199	Ä
136	ê	168	¿	200	Ł
137	ë	169	®	201	Œ
138	è	170	¬	202	℥
139	ĩ	171	½	203	Ŧ
140	î	172	¼	204	Ŧ
141	ì	173	¡	205	=
142	Ä	174	«	206	‡
143	Å	175	»	207	□
144	É	176	⋮	208	ø
145	æ	177	⋮	209	Ð
146	Æ	178	⋮	210	È
147	ô	179	⋮	211	Ê
148	ö	180	⋮	212	Ë
149	ò	181	À	213	Ì
150	û	182	Â	214	Í
151	ù	183	Ã	215	Î
152	ÿ	184	©	216	Ï
153	Ö	185	ª	217	Ɔ
154	Ü	186	»	218	Ɔ
155	ø	187	»	219	■
156	£	188	»	220	■
157	Ø	189	¢	221	■
158	×	190	¥	222	■
159	f	191	γ	223	■
				224	Ó
				225	ß
				226	Ô
				227	Ò
				228	ö
				229	Õ
				230	μ
				231	þ
				232	Þ
				233	Ú
				234	Ù
				235	Û
				236	ý
				237	Ý
				238	—
				239	·
				240	≡
				241	±
				242	≡
				243	¾
				244	¶
				245	§
				246	÷
				247	°
				248	°
				249	·
				250	·
				251	¹
				252	³
				253	²
				254	■
				255	nbsp

Sin embargo, 8 bits siguen siendo insuficientes para codificar todos los alfabetos del planeta. Por lo que cada zona usaba su propia tabla ASCII extendida (siempre los 128 primeros son el ASCII original). Ante el caos consiguiente, la ISO decidió normalizar dichas tablas de códigos para conseguir versiones estándares de los mismos. Lo hizo mediante las siguientes normas

- 🚩 **8859-1.** ASCII extendido para Europa Occidental (incluye símbolos como ñ o ß)
- 🚩 **8859-2.** ASCII extendido para Europa Central y del Este (incluye símbolos como Ž o ě)
- 🚩 **8859-3.** ASCII extendido para Europa del Sur (incluye símbolos como Ġ o ï)
- 🚩 **8859-4.** ASCII extendido para Europa del Norte (incluye símbolos como ø o å)

- ✚ **8859-5.** ASCII extendido para alfabeto cirílico (incluye símbolos como д o Ж)
- ✚ **8859-6.** ASCII extendido para alfabeto árabe (incluye símbolos como ù o ُ)
- ✚ **8859-7.** ASCII extendido para alfabeto griego moderno (incluye símbolos como φ o α)
- ✚ **8859-8.** ASCII extendido para alfabeto hebreo (incluye símbolos como ך o ם)
- ✚ **8859-9.** ASCII extendido, versión de 8859-1 que incluye símbolos turcos en lugar de otros poco utilizados
- ✚ **8859-10.** ASCII extendido, versión de 8859-4 que incluye símbolos más utilizados en las lenguas nórdicas actuales
- ✚ **8859-11.** ASCII extendido para alfabeto tailandés (incluye símbolos como ณ o ฅ)
- ✚ **8859-12.** ASCII extendido para alfabeto devanagari de India y Nepal que ya no se usa
- ✚ **8859-13.** ASCII extendido para alfabetos bálticos con símbolos que no estaban en 8859-4
- ✚ **8859-14.** ASCII extendido para alfabeto celta (incluye símbolos como ŵ o W)
- ✚ **8859-15.** ASCII extendido, versión de 8859-1 que incluye el símbolo del euro y símbolos de lenguas bálticas. Es el recomendado actualmente para Europa Occidental.
- ✚ **8859-16.** ASCII extendido, versión de 8859-1 pensada para los países del sureste de Europa
- ✚ **2022-JP.** Símbolos japoneses (parte 1)
- ✚ **2022-JP-2.** Símbolos japoneses (parte 2)
- ✚ **2022-KR.** Símbolos coreanos

Este problema sigue existiendo en la actualidad, de modo que en los documentos de texto hay que indicar el sistema de codificación utilizado (el caso más evidente son las páginas web), para saber cómo interpretar los códigos del archivo. Así en 8859_1 el código 245 es el carácter ñ y en 8859_2 es el carácter ó

Unicode

La complicación de las tablas de código se intenta resolver gracias al sistema Unicode que ha conseguido incluir los caracteres de todas las lenguas del planeta a cambio de que cada carácter ocupe más de un byte (ocho bits). En Unicode a cada símbolo se le asigna un número (evidentemente los 128 primeros son los originales de ASCII para mantener la compatibilidad con los textos ya codificados y de hecho los 256 primeros son la tabla ISO-8859-1).

Para ello el organismo también llamado Unicode, participado por numerosas e influyentes empresas informáticas y coordinado por la propia ISO, ha definido tres formas de codificar los caracteres:

🚦 **UTF-8.** Es la más utilizada (y la más compleja de usar para el ordenador).

Utiliza para cada carácter de uno a cuatro caracteres, de forma que:

- Utilizan uno los que pertenecen al código ASCII original
- Dos los pertenecientes a lenguas latinas, cirílicas, griegas, árabes, hebreas y otras de Europa, Asia Menor y Egipto
- Tres para símbolos fuera de los alfabetos anteriores como el chino o el japonés
- Cuatro para otros símbolos: por ejemplo los matemáticos y símbolos de lenguas muertas como el fenicio o el asirio o símbolos asiáticos de uso poco frecuente.

🚦 **UTF-16.** Utiliza para cada carácter dos (para los dos primeros grupos del punto anterior) o cuatro caracteres (para el resto). Es más sencillo que el anterior

🚦 **UTF-32.** La más sencilla de todas. Cada carácter independientemente del grupo al que pertenezca ocupa 4 caracteres. No se utiliza.

Archivos binarios y archivos de texto

Ventajas de los archivos binarios

- 🚦 Ocupan menos espacio que los archivos de texto, ya que optimizan mejor su codificación a binario (por ejemplo el número 213 ocupa un solo byte y no tres como ocurriría si fuera un texto).
- 🚦 Son más rápidos de manipular por parte del ordenador (se parecen más al lenguaje nativo del ordenador)

- ✚ Permiten el acceso directo a los datos. Los archivos de texto siempre se manejan de forma secuencial, más lenta
- ✚ En cierto modo, permiten ocultar el contenido que de otra forma sería totalmente visible por cualquier aplicación capaz de entender textos (como el bloc de notas). Es decir los datos no son fácilmente entendibles

Ventajas de los archivos de texto

- ✚ Son ideales para almacenar datos para exportar e importar información a cualquier dispositivo electrónico ya que todos son capaces de interpretar texto
- ✚ Son directamente modificables, sin tener que acudir a software específico
- ✚ Su manipulación es más sencilla que la de los archivos binarios
- ✚ Son directamente transportables y entendibles por todo tipo de redes

2. El intercambio de información

Los problemas relacionados con el intercambio de información entre aplicaciones y máquinas informáticas es tan viejo como la propia informática.

El problema parte del hecho de haber realizado un determinado trabajo con un software en un ordenador concreto y después querer pasar dicho trabajo a otro software en ese u otro ordenador.

Para hacer ese proceso con archivos binarios el origen y el destino de los datos deben comprender cómo codificar y decodificar la información. En muchos casos esto ha sido un gran problema que ha obligado a que todos los trabajadores y trabajadoras hayan tenido que adaptarse al software de la empresa, y, por supuesto, en toda la empresa utilizar dicho software.

En la informática actual el problema es aún mayor al tener una necesidad de disponibilidad global del trabajo y además la posibilidad de ver dicho trabajo en dispositivos de todo tipo como mini ordenadores, tablets o, incluso, teléfonos móviles.

Por ello, poco a poco, han aparecido formatos binarios de archivo que han sido estándares de facto (aunque no han sido reconocidos por ningún organismo de estándares), como por ejemplo el formato documental PDF, el formato de imagen JPEG, la música MP3 o el formato MPEG de vídeo.

Pero sigue habiendo empresas que utilizan formato propio por la idea de que sus formatos de archivo están directamente relacionados con la calidad de su software, es decir, razonan que el software que fabrican es muy potente y necesitan un formato binario propio compatible con esa potencia. De ahí que muchas veces la opción para exportar e importar datos sea utilizar conversores, capaces de convertir los datos de un formato a otro (por ejemplo de Word a Open Office; de MP3 a MOV de Apple, etc.).

El texto como el formato más versátil

Hay un formato de archivo que cualquier dispositivo es capaz de entender: el texto. La cuestión es que los archivos de texto sólo son capaces de almacenar texto plano, es decir, sólo texto sin indicar ningún formato o añadir información no textual.

Debido a la facilidad de ser leído con cualquier aparato, se intenta que el propio texto sirva para almacenar otros datos. Para ello, dentro del archivo habrá contenido que no se interpretará como texto, sino que hay texto en el archivo que se marca de manera especial haciendo que signifique otra cosa. Desde hace muchos años hay dos campos en los que esta idea ha funcionado bien: en las bases de datos y en los procesadores de texto. Actualmente, el éxito de Internet ha permitido espolpear esta tecnología a otros campos.

3. Lenguajes de marcas

Como se ha comentado en el punto anterior, el problema de la exportación de datos ha puesto en entredicho a los archivos binarios como fuente para exportar e importar información.

En su lugar parece que los archivos de texto poseen menos problemas. Por ello, se ha intentado que los archivos de texto plano (archivos que sólo contienen texto y no otros datos binarios) pudieran servir para almacenar otros datos, como por ejemplo detalles sobre el formato del propio texto u otras indicaciones.

Los procesadores de texto fueron el primer software en encontrarse con este dilema. Puesto que son programas que sirven para escribir texto, parecía que lo lógico era que sus datos se almacenaran como texto. Pero necesitan guardar datos

referidos al formato del texto, tamaño de la página, márgenes, etc. La solución clásica ha sido guardar la información de formato de forma binaria, lo que provoca los ya comentados problemas de portabilidad.

Algunos procesadores de texto optaron por guardar toda la información como texto, haciendo que las indicaciones de formato no se almacenen de forma binaria sino textual. Dichas indicaciones son caracteres marcados de manera especial para que así un programa adecuado pueda traducir dichos caracteres, no como texto, sino como operaciones que finalmente permitirán mostrar el texto del documento de forma adecuada.

La idea del marcado procede del inglés *marking up* término con el que se referían a la técnica de marcar manuscritos con lápiz de color para hacer anotaciones, como por ejemplo la tipografía a emplear en las imprentas. Este mismo término se ha utilizado para los documentos de texto que contienen comandos u anotaciones.

Las posibles anotaciones o indicaciones incluidos en los documentos de texto han dado lugar a lenguajes (entendiendo que en realidad son formatos de documento y no lenguajes en el sentido de los lenguajes de programación de aplicaciones) llamados lenguajes de marcas, lenguajes de marcado o lenguajes de etiquetas.

Orígenes de los lenguajes de marcas

Se considera a Charles Goldfarb como al padre de los lenguajes de marcas. Se trata de un investigador de IBM que propuso ideas para que los documentos de texto tuvieran la posibilidad de indicar el formato del mismo. Al final ayudó a realizar el lenguaje GML de IBM, el cual puso los cimientos del futuro SGML ideado por el propio Goldfarb.

TeX y LaTeX

En la década de los 70 Donald Knuth (uno de los ingenieros informáticos más importantes de la historia, padre del análisis de algoritmos) creó *TeX* para producir documentos científicos utilizando una tipografía y capacidades que fueran iguales en cualquier computadora, asegurando además una gran calidad en los resultados. Para ello apoyó a *TeX* con tipografía especial (fuentes Modern Computer) y un lenguaje de definición de tipos (METAFONT). *TeX* ha tenido cierto éxito en la comunidad científica gracias a sus 300 comandos que permiten crear documentos

con tipos de gran calidad, para ello se necesita un programa capaz de convertir el archivo *TeX* a un formato de impresión.

El éxito de *TeX* produjo numerosos derivados de los cuales el más popular es *LaTeX*. Se trata de un lenguaje que intenta simplificar a *TeX*, fue definido en 1984 por Leslie Lamport, aunque después ha sido numerosas veces revisado. Al utilizar comandos de *TeX* y toda su estructura tipográfica, adquirió rápidamente notoriedad y sigue siendo utilizado para producir documentos con expresiones científicas, de gran calidad. La idea es que los científicos se centren en el contenido y no en la presentación.

Ejemplo de código LaTeX:

```
\documentclass[12pt]{article}
\usepackage{amsmath}
\title{\Ejemplo}
\begin{document}
Este es el texto ejemplo de \LaTeX{}
Con datos en \emph{cursiva} o \textbf{negrita}.
Ejemplo de f'ormula
\begin{align}
E &= mc^2
\end{align}
\end{document}
```

Que con un traductor daría lugar al resultado:

Este es el texto ejemplo de \LaTeX
Con datos en *cursiva* o **negrita**. Ejemplo de fórmula

$$E = mc^2 \quad (1)$$

RTF

RTF es el acrónimo de Rich Text Format (Formato de Texto Enriquecido) un lenguaje ideado por Microsoft en 1987 para producir documentos de texto que incluyan anotaciones de formato.

Actualmente se trata de un formato aceptado como texto con formato y en ambiente Windows es muy utilizado como formato de intercambio entre distintos procesadores por su potencia.

El procesador de texto Word Pad incorporado por Windows lo utiliza como formato nativo.

Ejemplo:

```
{\rtf1\ansi\ansicpg1252\deff0\deflang3082{\fonttbl{\f0\fnil\fcharset0Calibri;}}\viewkind4\uc1\pard\sa200\sl276\slmult1\lang10\f0\fs22 soy \icursiva\i0\par }
```

Produce el resultado:

soy cursiva

SGML

Se trata de la versión de GML que estandarizaba el lenguaje de marcado y que fue definida finalmente por ISO como estándar mundial en documentos de texto con etiquetas de marcado. La estandarización la hace el subcomité SC24 que forma parte del comité JTC1 del organismo IEC de ISO que se encarga de los estándares electrónicos e informáticos (en definitiva se trata de una norma ISO/IEC JTC1/SC24, concretamente la 8879).

Su importancia radica en que es el padre del lenguaje XML y la base sobre la que se sostiene el lenguaje HTML.

En SGML las etiquetas que contienen indicaciones para el texto se colocan entre símbolos < y >. Las etiquetas se cierran con el signo /. Es decir las reglas fundamentales de los lenguajes de etiquetas actuales ya las había definido SGML.

En realidad (como XML) no es un lenguaje con unas etiquetas concretas, sino que se trata de un lenguaje que sirve para definir lenguajes de etiquetas; o más exactamente, es un lenguaje de marcado que sirve para definir formatos de documentos de texto con marcas. Entre los formatos definidos mediante SGML, sin duda HTML es el más popular.

PostScript

Se trata de un lenguaje de descripción de páginas, de hecho es el más popular. Permite crear documentos en los que se dan indicaciones potentísimas sobre como mostrar información en el dispositivo final. Se inició su desarrollo en 1976 por John Warnock y dos años más tarde se continuó con la empresa Xerox, hasta que en 1985 el propio Warnock funda Adobe Systems y desde esa empresa se continúa su desarrollo.

Es en realidad todo un lenguaje de programación que indica la forma en que se debe mostrar la información que puede incluir texto y el tipo de letra del mismo, píxeles individuales y formas vectoriales (líneas, curvas). Sus posibilidades son muy amplias.

Ejemplo:

%!PS	
/Courier	<i>% Elige el tipo de letra</i>
20 selectfont	<i>% Establece el tamaño de la letra y</i>
	<i>% la toma como el tipo de letra en uso</i>
72 500 moveto	<i>% Coloca el cursor en las coordenadas</i>
	<i>% 72, 500 (contando los píxeles desde</i>
	<i>% la esquina izquierda de la página)</i>
(Hola mundo!) show	<i>% Escribe el texto entre paréntesis,</i>
showpage	<i>% Imprime el resultado</i>

HTML

Tim Bernes Lee utilizó SGML para definir un nuevo lenguaje de etiquetas que llamó Hypertext Markup Language (lenguaje de marcado de hipertexto) para crear documentos transportables a través de Internet en los que fuera posible el hipertexto; es decir, que determinadas palabras marcadas de forma especial permitieran abrir un documento relacionado con ellas.

A pesar de tardar en ser aceptado, HTML fue un éxito rotundo y la causa indudable del éxito de Internet. Hoy en día casi todo en Internet se ve a través de documentos HTML, que popularmente se denominan páginas web.

Inicialmente estos documentos se veían con ayuda de intérpretes de texto (como por ejemplo el Lynx de Unix), que simplemente coloreaban el texto y remarcaban el hipertexto. Después el software se mejoró y aparecieron navegadores con capacidad más gráfica para mostrar formatos más avanzados y visuales.

XML

Se trata de un subconjunto de SGML ideado para mejorar el propio SGML y con él definir lenguajes de marcado con sintaxis más estricta, pero más entendibles. Su popularidad le ha convertido en el lenguaje de marcado más importante de la actualidad y en el formato de documentos para exportación e importación más exitoso. Aunque el formato JSON prácticamente lo ha desbancado en lo que se refiere al intercambio de datos entre aplicaciones.

JSON (Abreviatura de JavaScript Object Notation)

Se trata de una notación de datos procedente del lenguaje JavaScript estándar (concretamente ECMA Script de 1999). En el año 2002 se le daba soporte desde muchos de los navegadores y su fama ha sido tal que ahora se ha convertido en una notación independiente de JavaScript que compite claramente con XML.

Se trata de una notación que realmente no se considera lenguaje de marcas, ya que no hay diferencia en el texto a través de etiquetas, sino que se basa en que el texto se divide en dato y metadato. De modo que el símbolo de los dos puntos separa el metadato del dato. Por otro lado, los símbolos de llave y corchete permiten agrupar de manera correcta los datos.

Ejemplo de JSON:

```
{
  "nombre": "Jorge",
  "apellido1": "Sánchez",
  "dirección": {
    "calle": "C/ Falsa nº 0",
    "localidad": "Palencia",
    "código Postal": "34001",
    "país": "España"
  },
  "teléfonos": [
    {
      "tipo": "fijo",
      "número": "999 999 999"
    },
    {
      "tipo": "móvil",
      "number": "666 666 666"
    }
  ]
}
```

Tipos de lenguajes de marcas

- ✚ **Orientados a la presentación.** En ellos al texto común se le añaden palabras encerradas en símbolos especiales que contienen indicaciones de formato que permiten a los traductores de este tipo de documentos generar un documento final en el que el texto aparece con el formato indicado. Es el

caso de los archivos generados por los procesadores de texto tradicionales, en los que al texto del documento se le acompaña de indicaciones de formato (como negrita, cursiva,...)

✚ ***Orientados a la descripción.*** En ellos las marcas especiales permiten dar significado al texto pero no indican cómo se debe presentar en pantalla el mismo. Sería el caso de XML (o de SGML) y JSON en el que la presentación nunca se indica en el documento; simplemente se indica una semántica de contenido que lo hace ideal para almacenar datos (por ejemplo, si el texto es un nombre de persona o un número de identificación fiscal).

✚ ***Orientados a procedimientos.*** Se trata de documentos en los que hay texto marcado especialmente que en realidad se interpreta como órdenes a seguir y así el archivo en realidad contiene instrucciones a realizar con el texto. Es el caso de LaTeX o PostScript donde, por ejemplo, se puede indicar una fórmula matemática.

En su origen el lenguaje HTML era un lenguaje orientado a la descripción, pero a lo largo de su evolución se fueron introduciendo etiquetas que indicaban como debía visualizarse cada elemento, con lo que paso a convertirse en un lenguaje híbrido: orientado a la descripción con elementos de visualización. En la actualidad se han desechado las etiquetas orientadas a la presentación, siendo un lenguaje utilizado para introducir elementos semánticos y no visuales.

Organizaciones Desarrolladoras de Lenguajes de Marcas

✚ ***Organización Internacional para la Estandarización (ISO).*** Organización con sede en Ginebra (Suiza). Se crea en 1947. Engloba a 163 países. Sus normas son voluntarias y no obligatorias pues no es gubernamental.

Después del éxito del GML y después de un largo proceso, publicó en 1986 el Standard Generalized Markup Language (SGML) con rango de Estándar Internacional con el código ISO 8879.

✚ ***World Wide Web Consortium (W3C).*** Se crea en 1994 por Tim Berners-Lee. Su función principal es tutelar el crecimiento y organización de la web. Su primer trabajo fue normalizar el lenguaje HTML, al crecer la web aumentaron las presiones para ampliar HTML, en vez de ello, creo unas reglas para que cualquiera pudiera crear lenguajes de marcas adecuados a

sus necesidades, pero manteniendo unas estructuras y sintaxis comunes que permitieran compatibilizarlos y tratarlos con las mismas herramientas. Ese conjunto de reglas es el XML, cuya primera versión se publicó en 1998.