

Alejandro Ramírez Aldariz

Consultor tecnológico de desarrollo de proyectos informáticos.

Puedes encontrarme en Autentia: Ofrecemos servicios de soporte a desarrollo, factoría y formación

Somos expertos en Java/J2EE

Ver todos los tutoriales del autor



Fecha de publicación del tutorial: 2015-02-24

Tutorial visitado 32 veces [Descargar en PDF](#)

## Índice

- Introducción
  - El porqué del sistema binario
  - Bits y bytes
  - ¿Qué es un caracter?
  - Tipos de caracteres
- ASCII
  - ¿Qué es ASCII?
  - Variantes del ASCII
- Unicode
  - ¿Qué es Unicode?
  - Unidades y puntos de código
    - ¿Como escribir puntos de código en Mac?
  - Estructura de Unicode
    - Planos
    - Bloques
      - Jeroglifos egipcios
      - Emoji
  - ¿Qué quieren decir cuando dicen...?
  - Carácter compuesto
  - Grafema
  - Glifos
  - Carácter percibido por el usuario
- Formatos de transformación Unicode
  - UTF-32
  - UTF-16
  - UTF-8
  - Punto de código a UTF-16
  - Punto de código a UTF-8
  - ¿Qué UTF usa este fichero?
- Editores hexadecimales
- Lenguajes de Programación
  - Objective-C
  - Otros Lenguajes de Programacion

## Introducción

Este artículo explica como representa texto el ordenador. Desde porque los ordenadores usan el sistema binario, qué es ASCII, a conceptos generales de Unicode. Unicode es un estandar complejo, pero aquí tienes la cultura general necesaria para hablar del tema, y orientarte por ti mismo.

### El sistema binario

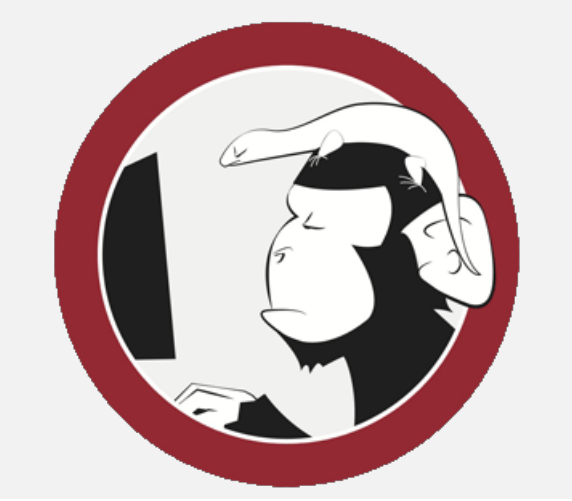
Un **componente electrónico** es un dispositivo capaz de manipular electrones o sus campos asociados. Por ejemplo, transistores, diodos, condensadores, etc.

Los componentes electrónicos se combinan y conectan mediante cables para formar **circuitos electrónicos** capaces de realizar operaciones complejas. Por ejemplo, una radio es un circuito electrónico que transforma una señal electromagnética en una señal sonora.

A un circuito electrónico lo llamamos **analógico** si opera con rangos continuos de señales electricas. Por ejemplo, un

Catálogo de servicios Autentia





Síguenos a través de:



Últimas Noticias

» 2015: ¡Volvemos a la oficina!

» Curso JBoss de Red Hat

» Si eres el responsable o líder técnico, considérate desafortunado. No puedes culpar a nadie por ser gris

» Portales, gestores de contenidos documentales y desarrollos a medida

» Comentando el libro Start-up Nation, La historia del milagro económico de Israel, de Dan Senor & Salu Singer

Histórico de noticias

Últimos Tutoriales

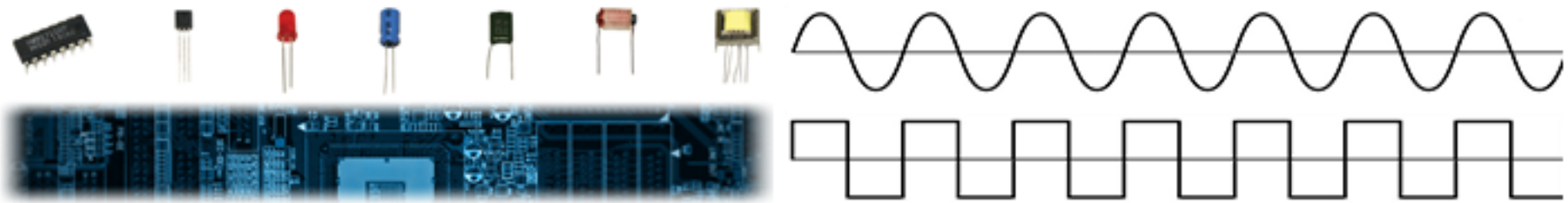
» Crea interfaces web amigables con Twitter Bootstrap

» Experimenta con tu código en Eclipse utilizando Scrapbooks

» Curso de WatchKit ¡ahora

altavoz amplifica la señal a un número indefinido de valores entre un volumen mínimo y máximo.

En cambio, si opera con señales discretas, decimos que el circuito electrónico es **digital**. Por ejemplo, un ordenador opera con dos voltajes de valores exactos: cero voltios, y el número de voltios que proporcione la fuente de voltaje, por ejemplo 1.4 voltios. Esto es así porque están contruidos con **transistores**, que son componentes electrónicos cuyo funcionamiento es más fiable cuando representa solo dos estados. A grandes rangos, cuantos más estados representas con un rango finito de voltaje, más posibilidades existen de confundir un estado con el otro.



En resumen, los ordenadores expresan información usando solo dos estados porque así sus componentes son más fiables.

El sistema de numeración más fácil de implementar con dos estados es el **sistema numérico binario**. El binario es similar al decimal, pero solo utiliza los dígitos uno y cero. Por ejemplo:

decimal	0	1	2	3	4	...
binario	0	1	10	11	100	...

### Bits y bytes

Cuando transferimos texto usando un ordenador, existen componentes electrónicos y cables que leen y escriben estas cadenas de ceros y unos. Cada uno de los dígitos 0 o 1 que circulan por el ordenador, se considera la unidad mínima de información en informática, y se conoce como **bit**.

Para acelerar la transferencia, los bits se transfieren en grupos de 8, 16, 32, o 64. Cuantos más bits lees a la vez, más cables y componentes necesitas para transferirlos. El equilibrio entre miniaturización, coste de construcción, y otras consideraciones prácticas, hace que los ordenadores actuales operen con bloques de solo 64 bits. Sin embargo hace 50 años, el límite práctico estaba en 8. Un grupo de 8 bits se conoce como **byte**.

El ordenador usa el sistema de numeración binario tanto para expresar datos, como para referenciar esos datos por su posición en memoria. Este es el motivo por el que la información se transfiere en multiplos de dos. ¿Sabes porqué? Es porque la capacidad expresiva de una base de numeración es  $\text{base}^{\text{número de dígitos}}$ . Por ejemplo con un dígito decimal puedes expresar 10 valores (0, 1, 2, 3, 4, 5, 6, 7, 8, 9), por tanto lo más eficiente es manejar información en multiplos de diez, de otro modo desperdicias la capacidad expresiva de un dígito decimal.

El motivo por el que además son multiplos de 8 es historico. Han existido computadores en los que un byte tenía 1, 6, 7, 8, 9, 12, 18, y 20 bits, pero la popularidad del chip 8080 (el cual usaba 8-bit) hizo que byte fuera sinónimo de 8 bits.

### ¿Qué es un caracter?

Un **caracter** es una unidad de información abstracta. Algunos ejemplos de caracteres son los números y letras del alfabeto, los signos de puntuación, los emojis, y los caracteres de control como un salto de línea.

Un **juego de caracteres** es una colección de caracteres usada para escribir en un lenguaje particular. Algunos idiomas usan uno, otros usan varios (por ejemplo, kanji y kana del japones), y otros usan variantes de un mismo juego de caracteres (los lenguajes europeos usan variantes del alfabeto romano).

Un **glifo** es una representación visual de un caracter. En un ordenador, los glifos asociados a caracteres suelen estar almacenados en **fuentes**, también llamados tipos de letra. Cada fuente contiene un estilo homogeneo de glifos para un alfabeto. Este es un ejemplo de varios glifos que representan un mismo caracter:



### Tipos de caracteres

A grandes rasgos hay tres sistemas de escritura

- **Alfabético**. Un alfabeto es el conjunto de letras usadas para escribir un lenguaje. Cada una representa un fonema o se usa junto a otras para representar un fonema. La mayoría de alfabetos en uso se basan en el alfabeto romano.
- **Silábico**. Contiene caracteres que representan cada una de las sílabas de un lenguaje. Una sílaba es una combinación de vocal o vocal y consonante. Es apropiado para lenguajes con pocas silabas, como el japones, que contiene un centenar. El inglés en cambio, contiene miles de combinaciones de conjuntos de vocales y consonantes.
- **Ideográfico**. Contiene caracteres que representan ideas. Por ejemplo, el chino, o el egipcio antiguo.

Estas categorías no son rígidas. Por ejemplo, el español es alfabético, pero también usa un ideograma como el símbolo de euro, que representa la estabilidad de Europa (es una e de Europa con dos líneas indican estabilidad).

## ASCII

### ¿Qué es ASCII?

El único tipo de dato que un ordenador es capaz de almacenar y manipular es el número binario. El truco para representar texto es asignar un número a cada letra del alfabeto, e implementar un sistema de renderizado que visualiza el caracter correspondiente a cada número.

Llamamos **código** al sistema de normas para convertir información simple en información codificada. Una **codificación de caracteres** por ejemplo, es un código que convierte caracteres a números.

La codificación de caracteres más popular se llama **ASCII** (American Standard Code for Information Interchange). ASCII asigna números decimales del 0 al 127 a las letras del alfabeto inglés, los números, y algunos signos de puntuación y de control. Los signos de control son indicadores para procesar información, por ejemplo, el salto de línea. La letra “A” en ASCII tiene asignado el número 65 (0x41 en hexadecimal), y las letras consecutivas del alfabeto tienen asignados números consecutivos:

sólo 9 dólares!

- » [Cómo implementar una nube de etiquetas con D3.js](#)
- » [Test de servicios REST con Spring MVC y Spring Test](#)

### Últimos Tutoriales del Autor

- » [¿Qué es Go?](#)



letra	binario	decimal
A	1000001	65
B	1000010	66
C	1000011	67
D	1000100	68
etc.		

Low Ascii

000:	013:J	026:→	039:’	052:4	065:A	078:M	091:í	104:h	117:u
001:Q	014:P	027:↔	040:(	053:5	066:B	079:O	092:\	105:i	118:v
002:Q	015:*	028:⌞	041:)	054:6	067:C	080:P	093:]	106:j	119:w
003:▼	016:►	029:↗	042:✱	055:7	068:D	081:Q	094:^	107:k	120:x
004:✚	017:◄	030:▲	043:+	056:8	069:E	082:R	095:	108:l	121:y
005:✚	018:‡	031:▼	044:,	057:9	070:F	083:S	096:˘	109:n	122:z
006:✚	019:!!	032:	045:-	058::	071:G	084:T	097:a	110:n	123:(
007:•	020:¶	033:!	046:.	059::	072:H	085:U	098:b	111:o	124:í
008:¶	021:§	034:”	047:↙	060:<	073:I	086:V	099:c	112:p	125:}
009:o	022:≡	035:¶	048:Q	061:=	074:J	087:W	100:d	113:q	126:˘
010:¶	023:‡	036:\$	049:1	062:>	075:K	088:X	101:e	114:r	127:Δ
011:Q	024:†	037:×	050:2	063:?	076:L	089:Y	102:f	115:s	
012:Q	025:↓	038:Δ	051:3	064:Q	077:M	090:Z	103:g	116:t	

Según este código, la representación de las letras “abc” en binario es la siguiente:

1000001 1000010 1000011

El documento que estas leyendo en este momento tiene una representación en binario en tu ordenador en forma de ceros y unos. A nivel físico se representa en la memoria mediante la presencia o ausencia de voltaje, y en el disco mediante magnetismo.

Cuando ASCII se publicó en 1964, las máquinas de entonces transferían 8 bits a la vez. Sin embargo ASCII usó solo 7, y reservó el octavo bit para detectar errores de transmisión. Esos 7 bits permiten expresar 2<sup>7</sup> = 128 valores, que ASCII usó para codificar 33 señales de control (usadas entre máquinas), y 95 números, símbolos de puntuación, y letras del alfabeto inglés.

### Variantes del ASCII

En los ’80 aparecieron variantes de ASCII que usaban el octavo bit para representar caracteres adicionales, lo que permitió representar caracteres adicionales propios de alfabetos no ingleses, como por ejemplo el español.

Debido a que 8 bits no son suficientes para representar todos los alfabetos del mundo, continuaron apareciendo variantes ASCII de 8 bits incompatibles entre sí. Estas variantes se llaman a veces **ASCII extendido**, pero no forman parte del estandar ANSI.

Hay varios conjuntos de “ASCII extendido”, cada uno de los cuales contiene codificaciones para muchos lenguajes:

- Windows code pages**, usado en aplicaciones gráficas Windows.
- OEM code pages**, usando en aplicaciones de consola Windows.
- ISO-8859** es un estandar **ISO** para codificación en 8 bits. Tiene 16 partes. La primera se llama ISO-8859-1, también conocida como Latin-1, que cubre la mayoría de lenguajes de Europa occidental.

En ASCII solo es posible trabajar con un alfabeto a la vez. ASCII tampoco es válido para representar alfabetos asiáticos, porque contienen miles de caracteres. Estas deficiencias, junto a la popularidad de la informática, y la potencia creciente de los ordenadores, hicieron deseable y posible la creación de un estandar más exhaustivo.

## Unicode

### ¿Qué es Unicode?

**Unicode** es una codificación de caracteres que asigna un número a cada uno de los caracteres de prácticamente todos los alfabetos existentes, incluyendo las lenguas muertas como el egipcio antiguo y otros. El estandar Unicode lo publica y mantiene el **consorcio Unicode**, una entidad sin ánimo de lucro formada por empresas e individuos.

El estandar Unicode también especifica algoritmos sobre como manejar texto, que probablemente solo necesites si realizas tareas especializadas. Por ejemplo:

- Dividir palabras y líneas.
- Ordenar texto.
- Formatear números, fechas, y horas.
- Mostrar texto que fluye de derecha a izquierda.
- Mostrar texto cuya forma escrita se combina y reordena.
- Tratar problemas de seguridad relativos a caracteres parecidos.

### Unidades y puntos de código

Unicode comenzó como un sistema de codificación de 16 bits, aunque desde Unicode 2.0 (publicado en 1996), es una codificación de 21 bits.

Un **punto de código** es el número con el que se identifica un caracter en el estandar Unicode. El punto de código se escribe con el formato U+xxxx donde las xxxx son de cuatro a seis dígitos en sistema de numeración hexadecimal.

En Unicode 7.0 el rango válido de puntos de código va de 0 a 10FFFF<sub>16</sub>. El hexadecimal se usa por conveniencia en lugar del binario, porque es más fácil recordar. Por ejemplo, Unicode asigna el número 65 a la letra a latina mayúscula. El punto de código correspondiente es U+0041 porque 65<sub>decimal</sub> = 0x41<sub>hexadecimal</sub>.

Unicode es compatible con la codificación ISO-8859-1, porque los 256 primeros caracteres de Unicode coinciden con los caracteres de ISO-8859-1. Esto hace que la mayoría de texto en uso requiera solo un byte por caracter.

El punto de código se representa con grupos de 8, 16, o 32 bits dependiendo respectivamente, de si el tipo de codificación es UTF-8, UTF-16, o UTF-32. Cada uno de estos grupos se llama unidad de código. Una **unidad de código** es el mínimo grupo de bits necesario para representar una unidad de texto codificado. En UTF-8 es 8 bit, en UTF-16 es 16 bit, y en UTF-32 es 32 bit. Por ejemplo, U+0041 requiere 2 bytes, es decir, dos unidades de código en UTF-8, y una unidad de código en UTF-16.

#### ¿Como escribir puntos de código en Mac?

*Para teclear caracteres en Mac usando puntos de código* sigue estos pasos:

- Ve a System Preferences > Keyboard > Input Sources
  - Añade Unicode Hex Input.
  - Activa Show Input menu in menu bar.
- Activa el teclado Unicode en la barra de menu.
- Para probarlo, deja pulsado alt y pulsa el código numérico del punto de código que desees (por ejemplo el 0041 para





Un **carácter compuesto** es una entidad Unicode que puede definirse como una secuencia de otros caracteres.

- Por ejemplo, U+00E9 (letra minúscula latina e con acento agudo) tiene el mismo significado y apariencia que U+0065 U+0301 (letra minúscula latina e, y acento agudo de combinación). Se consideran **canónicamente equivalentes**, pero no iguales porque están hechos de diferentes puntos de código.
- Los caracteres compuestos existen para permitir implementaciones de Unicode en sistemas que solo permiten representar cada caracter por separado.

Los **caracteres compatibles equivalentes** son aquellos que representan el mismo caracter abstracto, pero tienen diferentes apariencias visuales. Por ejemplo, el carácter **ff** (ligadura latina minúscula **f f**, U+FB00) es compatible con (pero no canónicamente equivalente a) la secuencia **f f** (dos letras latinas minúsculas **f**, U+0066 U+0066).

Si queremos comparar texto, por ejemplo durante una búsqueda, quizás queramos tratar caracteres canónicamente equivalentes, y caracteres compatibles equivalentes como un mismo caracter. En ese caso necesitaremos aplicar un **algoritmo de normalización** que convierten una cadena a una representación única que puede ser comparada binariamente. Hay cuatro formas de normalización llamadas C, D, KD, y KC. Estos algoritmos permiten comparar caracteres, pero pueden eliminar distinciones de formato que evitan conversiones en ambos sentidos.

## Grafema

Un **grafema** es una unidad atómica en un lenguaje.

Un **cluster de grafemas** es una secuencia de puntos de código que representa un grafema.

## Glifos

Un **glifo** es un carácter gráfico que representa un caracter abstracto.

- Una **forma variante** es cada uno de los glifos de un conjunto de glifos que representan el mismo caracter abstracto.
- Una **secuencia de variaciones** es un mecanismo para escoger una forma variante. Consiste de un caracter base seguido de uno de los 256 **selectores de variación**, que son VS1 a VS256, U+FE00 a U+FE0F, y U+E0100 a U+E01EF).

Por ejemplo, los glifos 🌂 y 🌂 representan el caracter abstracto “Paraguas con las gotas de agua”, y son por tanto formas variantes de dicho caracter. La secuencia de variaciones para representarlos son:

- U+2614 U+FE0F para 🌂
- U+2614 U+FE0E para 🌂

## Caracter percibido por el usuario

Lo que el usuario percibe como caracter. Por ejemplo, en inglés “ch” son dos letras, pero en checo y eslovaco se considera una.

## Formatos de transformación Unicode

Los formatos de transformación Unicode (UTF) especifican como serializar un punto de código abstracto a bytes en memoria. Se les conoce popularmente como “codificaciones”, aunque este nombre es incorrecto.

Debido a que Unicode es una codificación que asigna números expresables con 21 bits, y que los ordenadores transfieren datos en multiples de 8 bits (8, 16, 32, ...), hay tres posibles modos de expresar Unicode:

- Usando una unidad de código de 32 bits (UTF-32).
- Usando una o dos unidades de código de 16 bits (UTF-16).
- Usando de una a cuatro unidades de código de 8 bits (UTF-8).

## UTF-32

- Formato con unidades de código de 32 bits.
- Representa cada punto de código con exactamente una unidad de código.
- Se usa poco porque los puntos de código más usados necesitan mucho menos de 32 bits.

## UTF-16

- Formato con unidades de código de 16 bits.
- Usa una unidad de código para representar puntos de código del BMP y dos para el resto. Al par de unidades de código usadas para representar un punto de código se les conoce como **par sustituto** (surrogate pair).
- El orden de los bytes es normalmente el de la CPU en la que se ejecuta una implementación dada.

## UTF-8

Características de UTF-8:

- Formato con unidades de código de 8 bits.
- Usa una unidad de código para los primeros 256 caracteres, dos para los 1920 siguientes, y tres para el resto de puntos de código del BMP. Usa cuatro unidades de código para codificar los puntos de código de otros planos.
- Cuando el texto utiliza unicamente caracteres ISO-8859-1, no hay diferencia entre un texto codificado en ISO-8859-1 y UTF-8.
- UTF-8 es con diferencia el formato más popular debido a su compatibilidad con ASCII e ISO-8859-1, y a su soporte en librerías.

## Punto de código a UTF-16

El BMP es un plano de 16 bits con potencial para codificar  $2^{16} = 65536$  caracteres en puntos de código U+0000 a U+FFFF. Dentro de este rango, los 2048 puntos de código que van de U+D800 a U+DFFF no tienen caracter asignado, y se consideran invalidos.

El algoritmo para convertir un punto de código a UTF-16 es este:

- Si es un punto de código valido del BMP se traduce directamente a UTF-16 usando una unidad de código.
- Si es un punto de código fuera del BMP se procede como sigue:
  - Resta 0x10000 del punto de código (0x10000 es el rango del BMP).

- Expresa el resultado como binario de 20 bits (rellena con ceros si hace falta).
- Usa los patrones 110110xxxxxxxx y 110111xxxxxxxx para los 10 bits más y menos significativos del resultado del punto anterior.

Por ejemplo, dado el caracter “grinning face” (😄) con Unicode U+1F600, su expresión en UTF-16 es U+D83D U+DE00 porque:

```
0x1F600 - 0x10000 = 0b 0000111101 1000000000
0b1101100000111101 to hex = 0xD83D
0b1101111000000000 to hex = 0xDE00
```

Punto de código a UTF-8

Dependiendo del rango del punto de código a codificar, usaremos uno, dos, tres, o cuatro unidades de código UTF-8.

0000—007F	0xxxxxxx	7 bits en 1 byte
0080—07FF	110xxxxx, 10xxxxxx	5+6 bits = 11 bits en 2 bytes
0800—FFFF	1110xxxx, 10xxxxxx, 10xxxxxx	4+6+6 bits = 16 bits en 3 bytes
10000—1FFFF	11110xxx, 10xxxxxx, 10xxxxxx, 10xxxxxx	3+6+6+6 bits = 21 bits en 4 bytes

Observa que si el bit más significativo no es cero, el número de unos antes del primer cero indica cuantas unidades de código codifican el punto de código.

Por ejemplo, dado el caracter “grinning face” (😄) con Unicode U+1F600, su expresión en UTF-8 es 0xF0 0x9F 0x98 0x83 porque

```
0x1F600 to binary = 0b111110110000000000
    000    011111    011000    000000
111110000 10011111 10011000 10000000
    0xF0    0x9F    0x98    0x80
```

¿Qué UTF usa este fichero?

El modo algoritmico de identificar los puntos de código de un texto es teniendo en cuenta las secciones anteriores, pero el modo rápido que yo uso es este otro.

Cargalo en un editor cualquiera (*Textmate*, *BEdit*, *TextWrangler*, *SublimeText*, *Pages*, *Xcode*, ...) y que te lo diga el editor. —este era mi gran secreto, ahora tú también lo conoces.

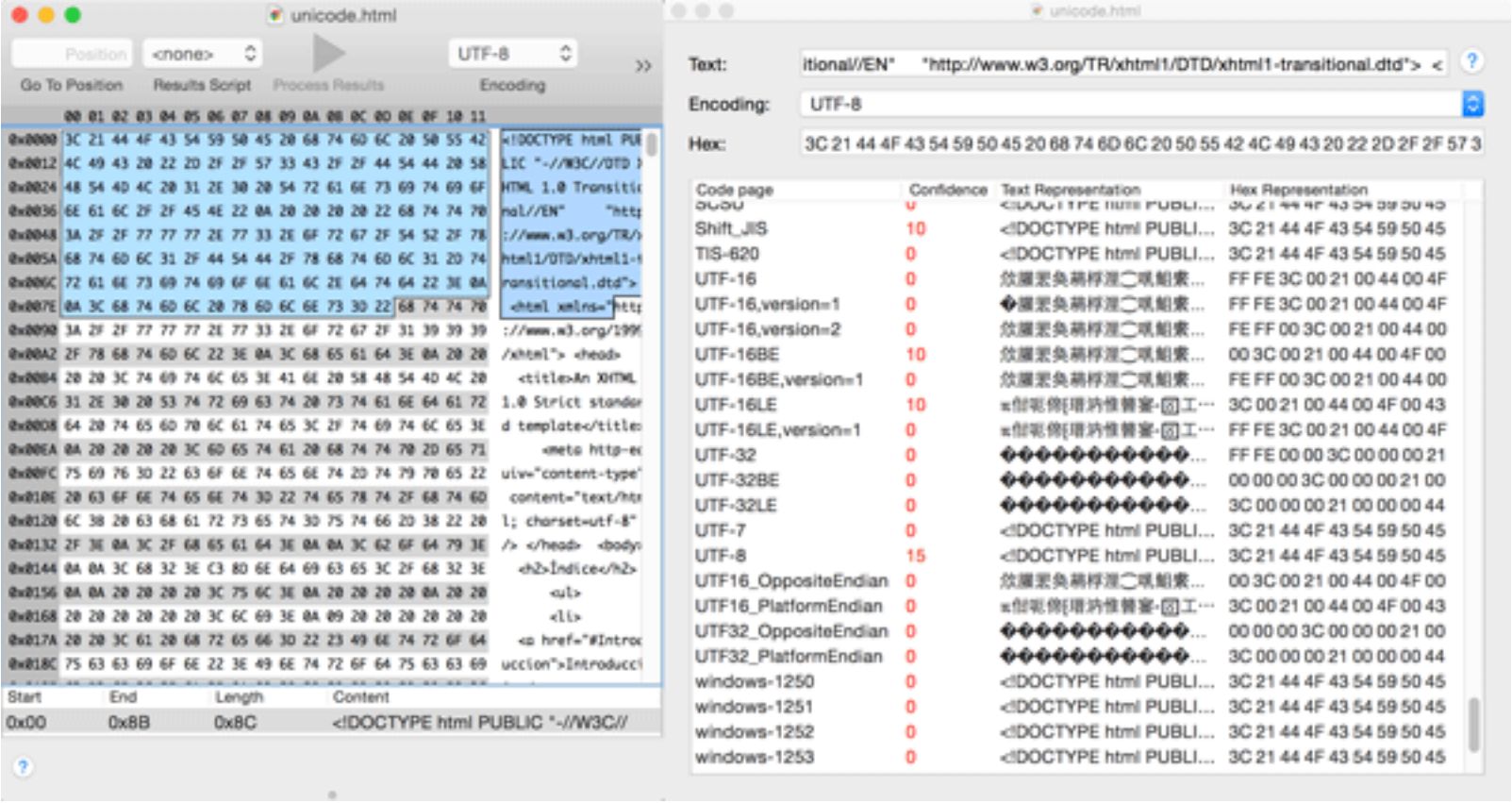
Deduce el formato a partir del punto de código de un caracter conocido. Por ejemplo, dado que un documento XML comienza con un caracter <, los puntos de código serán:

Encoding	Puntos de código
UTF-8	U+3C
UTF-16 Low Endian	U+3C00
UTF-16 Big Endian	U+003C
UTF-32 (BE)	U+0000 U+003C
UTF-32 (LE)	U+3C00 U+0000

Comprueba si la cadena o fichero comienza con un BOM. El BOM o **Marca de Orden de Byte** es un punto de código opcional que aparece al comienzo de un fichero Unicode. Indica el formato Unicode y el orden de sus bytes (endianness). En la práctica, la codificación más usada es UTF-8, y debido a que usa unidades de código de un byte, no necesita BOM. El estandar Unicode no requiere ni recomienda el uso de BOMs, pero define los siguientes:

Encoding	BOM (hexadecimal)
UTF-8	EF BB BF
UTF-16 Low Endian	EF BB
UTF-16 Big Endian	FE FF
UTF-32 (BE)	00 00 FE FF
UTF-32 (LE)	FF FE 00 00

Carga el texto en *Synalyze It!* (\$39), abre el comparador de encodings dandole a *Windows > Code Page Compare*, y pega parte del texto del documento en el campo *Text*. Te mostrará el texto interpretado en multitud de encodings, y podrás ver visualmente en cual se corresponde con lo que esperabas.







ser humano, y ya yo lo he alcanzado y superado.

HALA!, CADA UNO A SU CASA!

## Referencias

- [Unicode 7.0 Chapter One: Introduction](#)
- [UnicodeChecker](#) Una herramienta para explorar y convertir Unicode.
- [NSString y Unicode](#) (objc.io)
- [UTF-8](#) (Wikipedia)
- [UTF-8, UTF-16, UTF-32 y BOM](#) (unicode.org)
- [Usted también puede hablar Unicode](#) - Ross Carter
- [Glosario de terminos Unicode](#) (unicode.org)

### A continuación puedes evaluarlo:

[Regístrate para evaluarlo](#)



### Por favor, vota +1 o compártelo si te pareció interesante

Share | 0 0

Anímate y coméntanos lo que pienses sobre este **TUTORIAL**:

» **Regístrate** y accede a esta y otras ventajas «



Esta obra está licenciada bajo licencia [Creative Commons de Reconocimiento-No comercial-Sin obras derivadas 2.5](#)