

What is Character Sets & Character Encodings & Character Reference ?

a. Character Sets

"Unicode (Universal Character Set)" karakterlerin benzersiz decimal sayılarla ifade edildiği listeye denir. Mesela A = 65, B = 66, C = 67 ,... olarak karakter setinde yer alır. Bu şekilde örneğin "hello" string'i Unicode karakter setinde

```
104 101 108 108 111
h e l l o
```

ile karşılık görür.

ASCII, internet tarihindeki ilk karakter seti standardıdır. Bu karakter seti tam olarak 128 farklı alphanumeric karakter tanımlar: Numbers (0-9), English Letters (A-Z) ve special Characters (örn; ! \$+-()@<>).

ANSI (Windows-1252) Windows tarihindeki ilk karakter setidir. Tam olarak 256 farklı karakter tanımlar.

ISO-8859-1 (Latin-1) HTML 4 için ilk default karakter setidir. İngiliz alfabesi dışında "ı", "ü", "ç",... gibi karakterleri de içerir. Tam olarak 256 farklı karakter tanımlar.

b. Character Encodings

Character Encoding karakterlerin benzersiz binary sayılarla ifade edildiği listeye denir. Örneğin UTF-8, UTF-16, UTF-32, ... birer character encoding'tirler.

UTF-8 (Unicode) dünyadaki hemen hemen tüm karakterleri ve sembolleri içerir. Bu yüzden HTML5'te default karakter seti olarak (karakter kodlayıcısı olarak) UTF-8 kullanılmaktadır.

c. "Character Sets" vs. "Character Encodings"

Character set'leri karakterlerin bilgisayarda hangi decimal değerlerle depolanacağını belirler. Character Encoding'ler ise karakterlerin bilgisayarda hangi binary değerlerle depolanacağını belirler. Dolayısıyla ASCII bir karakter setidir ve listesindeki karakterlerin hangi decimal değerlerle bilgisayarda depolanacağını tanımlar. UTF-8 bir karakter setidir ve listesindeki karakterlerin hangi binary değerlerle bilgisayarda depolanacağını tanımlar.

Character Set ve Encodings Üzerine Ekstra

Internet'te kullanılan ANSI ve ISO-8859-1 (Latin-1) karakter setleri çok kısıtlı oldukları için HTML 4 ve sonrası UTF-8 karakter setini desteklemektedir.

```
// HTML 4'de charset belirleme
<meta http-equiv="Content-Type" content="text/html; charset=utf-8">
```

```
// HTML5'de charset belirleme
```

```
<meta charset="UTF-8">
```

```
// Include karabuk 'deki <meta etiketi
```

```
<meta http-equiv="Content-Type" content="text/html; charset=utf-8">
```

Unicode karakter setleri HTML, XML, Java, Javascript, E-mail, ASP, PHP, ... gibi birçok teknolojide kullanılmaktadır. Ayrıca Unicode karakter setlerini birçok işletim sistemi ve tüm modern web tarayıcıları da desteklemektedir.

d. Character Reference

Karakter referansı Unicode karakter setlerindeki karakterleri referans yoluyla çağırmanızı sağlayan kodlamalardır. HTML'deki numerik "karakter referans"ları Unicode (Universal Character Set) 'daki karşılık gelen bir karakteri gösterirler. Numerik karakter referanslarının formatı şu şekildedir:

&#nnnn;

ya da

&#xhhhh;

nnnn olan format decimal form'dur. hhhh olan ise hexadecimal form'dur. nnnn ve hhhh herhangi bir sayı olabilir. Karakter referansı mevcut karakter setinde tanımlı karakterleri örneğin html dökümanına referans yoluyla dahil edebilmemizi sağlar. Mesela Türkçe klavye kullanan bir kimse klavyesinden matematiksel sembolleri normal şartlarda çıkaramaz. Çünkü karakter setinde tanımlı o matematiksel sembolleri klavyeden çıkarmak uzun ve karmaşık tuş kombinasyonları gerektirir. Bu komplike yöntem yerine ilgili karakterin karakter referansı kullanılabilir ve istenilen semboller böylece ekrana verilebilir. Aşağıda bir html dökümanına klavyeden girilebilmesi mümkün olan karakterlerin "karakter referansları" verilmiştir:

Unicode character	Character Reference (decimal)	Character Reference (hexadecimal)	Effect
U+0020	 	 	(space)
U+0021	!	!	!
U+0022	"	"	"
U+0023	#	#	#
U+0024	$	$	\$
U+0025	%	%	%
U+0026	&	&	&
U+0027	'	'	'
U+0028	(((
U+0029)))
U+002A	*	*	*
U+002B	+	+	+
U+002C	,	,	,
U+002D	-	-	-
U+002E	.	.	.
U+002F	/	/	/
U+0030	0	0	0
U+0031	1	1	1
U+0032	2	2	2
U+0033	3	3	3
U+0034	4	4	4
U+0035	5	5	5
U+0036	6	6	6
U+0037	7	7	7
U+0038	8	8	8
U+0039	9	9	9
U+003A	:	:	:
U+003B	;	;	;
U+003C	<	<	<
U+003D	=	=	=
U+003E	>	>	>
U+003F	?	?	?
U+0040	@	@	@
U+0041	A	A	A
U+0042	B	B	B
U+0043	C	C	C
U+0044	D	D	D
U+0045	E	E	E
U+0046	F	F	F
U+0047	G	G	G
U+0048	H	H	H
U+0049	I	I	I
U+004A	J	J	J
U+004B	K	K	K
U+004C	L	L	L

U+004D	M	M	M
U+004E	N	N	N
U+004F	O	O	O
U+0050	P	P	P
U+0051	Q	Q	Q
U+0052	R	R	R
U+0053	S	S	S
U+0054	T	T	T
U+0055	U	U	U
U+0056	V	V	V
U+0057	W	W	W
U+0058	X	X	X
U+0059	Y	Y	Y
U+005A	Z	Z	Z
U+005B	[[[
U+005C	\	\	\
U+005D]]]
U+005E	^	^	^
U+005F	_	_	_
U+0060	`	`	`
U+0061	a	a	a
U+0062	b	b	b
U+0063	c	c	c
U+0064	d	d	d
U+0065	e	e	e
U+0066	f	f	f
U+0067	g	g	g
U+0068	h	h	h
U+0069	i	i	i
U+006A	j	j	j
U+006B	k	k	k
U+006C	l	l	l
U+006D	m	m	m
U+006E	n	n	n
U+006F	o	o	o
U+0070	p	p	p
U+0071	q	q	q
U+0072	r	r	r
U+0073	s	s	s
U+0074	t	t	t
U+0075	u	u	u
U+0076	v	v	v
U+0077	w	w	w
U+0078	x	x	x
U+0079	y	y	y
U+007A	z	z	z
U+007B	{	{	{
U+007C	|	|	
U+007D	}	}	}
U+007E	~	~	~

Örneğin A karakteri ASCII karakter setinde 65 decimal sayısıyla ifade edilmekteydi. Bu karakter setindeki A karakterini referans yoluyla çağırarak için

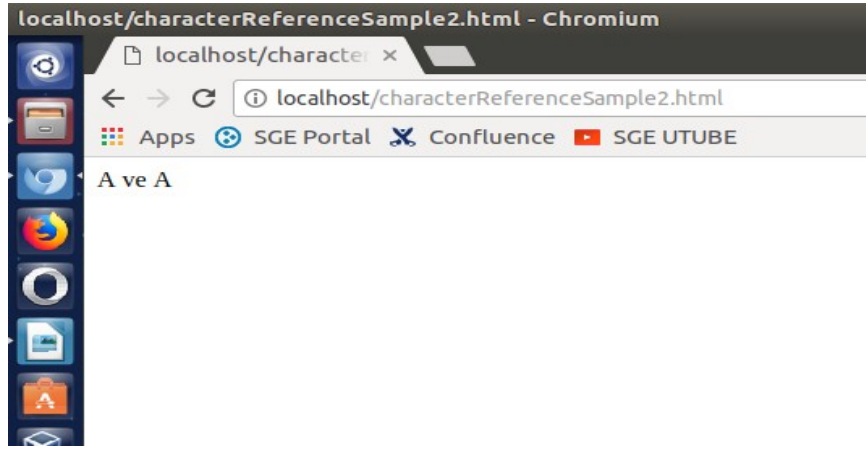
U+0041	A	A	A
--------	-------	--------	---

A decimal referansı ya da A hexadecimal referansı kullanılabilir.

/var/www/characterReferenceSample2.html

A ve A

Çıktı:



Böylece karakter referansları ile biz karakter setinde tanımlı karakterleri html dökümanına ekleyebiliriz.

Aşağıda ise karakter setinde yer alsa bile normal şartlarda klavyeden çıkarılmayacak bir matematiksel sembolün karakter referansı gösterilmiştir:

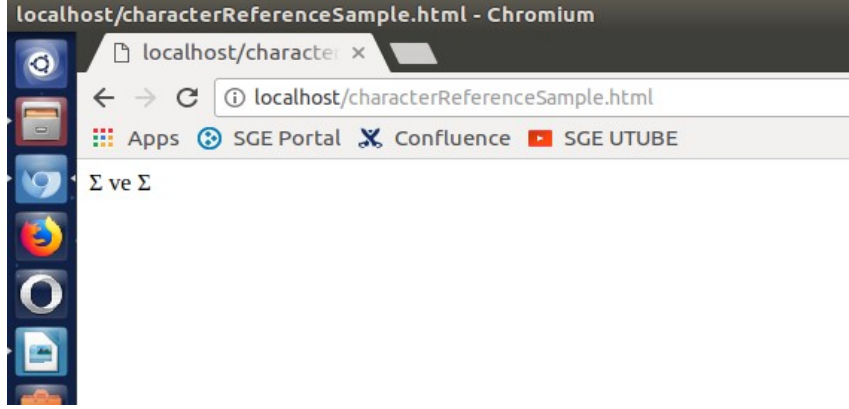
Unicode character	Character Reference (decimal)	Character Reference (hexadecimal)	Effect
U+03A3	Σ	Σ	Σ

Karakter setindeki Σ sembolünü referans yoluyla çağırarak için Σ decimal referansı ya da Σ hexadecimal referansı kullanılabilir.

/var/www/characterReferenceSample.html

Σ ve Σ

Çıktı:



Character Reference'ı yanında bir de Character Entity Reference'ı vardır. Character Reference'ları decimal ya da hexadecimal sayıları kullanarak karakter setindeki bir karakteri göstermeye yararken "Character Entity Reference"ları ise isim kullanarak karakter setindeki bir karakteri göstermeye yarar. Örneğin HTML'de öntanımlı entity'ler (> , " , & , ... v.b.) ilgili karakterleri gösterirken DTD'de öntanımlı entity'ler ve ayrıca explicitly olarak kendi tanımladığımız entity'ler ilgili karakterleri gösterir. Character entity referanslarının formatı şu şekildedir.

&name;

Kaynaklar

https://www.w3schools.com/html/html_charset.asp

https://www.w3schools.com/charsets/ref_html_utf8.asp

<https://stackoverflow.com/questions/2281646/whats-the-difference-between-encoding-and-charset>

<https://tr.wikipedia.org/wiki/Unicode>