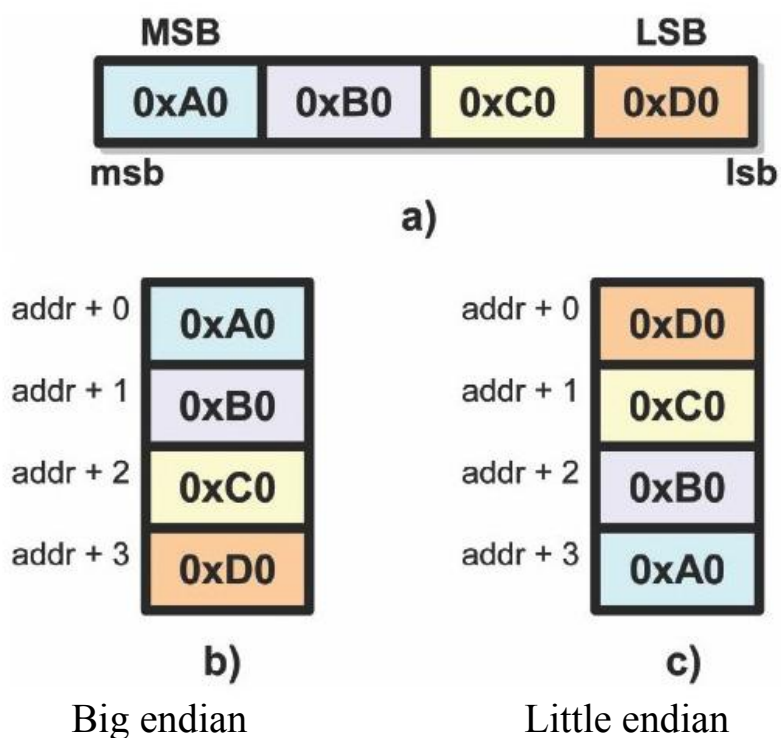


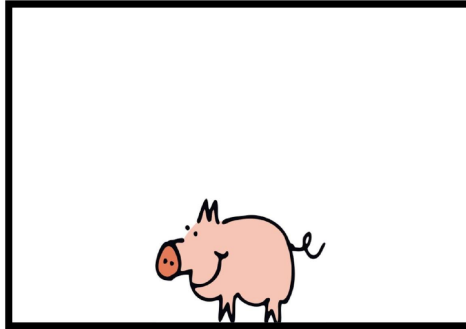
# Il Problema dell'eterogeneità

- Nelle reti di calcolatori vi è un'estrema eterogeneità di sistemi (hardware e software)
- Client e Server possono eseguire su architetture diverse che usano differenti rappresentazioni dei dati:
  - caratteri (ASCII, ISO 8859, Unicode, ...)
  - interi (dimensione 4 o 8 byte, rappresentazione in complemento a 1 o a 2, ...)
  - reali (lunghezza exp e mantissa, formato, ...)
  - ordine byte all'interno di una parola (little endian o big endian)
- Necessità di definire una rappresentazione comune dei dati e di implementare meccanismi per gestirla

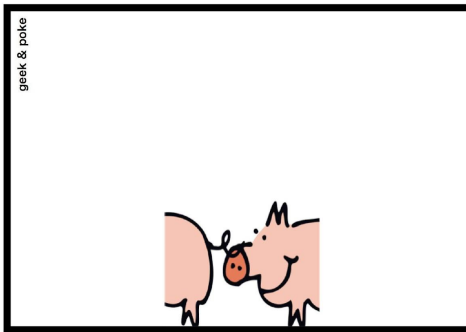
## Es. little endian o big endian



## SIMPLY EXPLAINED



BIG-ENDIAN



LITTLE-ENDIAN

## Caveat

- Attenzione! Il problema della rappresentazione eterogenea dei dati tra diverse piattaforme HW/SW presenta complessità che vanno ben oltre le comunicazioni di rete
- Ad esempio, anche la ricompilazione dello stesso codice sorgente su piattaforme diverse può presentare spiacevoli sorprese. Si veda l'articolo:

“Twice the Bits, Twice the Trouble:  
Vulnerabilities Induced by Migrating to 64-Bit Platforms”  
<https://www.tu-braunschweig.de/Medien-DB/sec/pubs/2016-ccs.pdf>

# Rappresentazione dati

- Per comunicare tra nodi eterogenei sono possibili due tipi di soluzioni:
  - Ogni nodo converte i dati nel formato specifico del destinatario (prestazioni)
  - Si concorda un formato comune di rappresentazione dei dati che i nodi useranno per comunicare tra loro (flessibilità)
- Supponendo di avere  $N$  diversi tipi di nodi, nel primo caso avrò bisogno di  $N*(N-1)$  procedure di conversione dati, nel secondo caso solo di  $2*N$  procedure

## Ma quale ISO livello 6?!?!

- La Socket API è un'interfaccia di basso livello (tra ISO L4 e L5) che purtroppo non fornisce alcuno strumento di questo tipo
- Non esiste un unico standard per il livello di presentazione
  - Molte soluzioni, con caratteristiche molto diverse, sono state sviluppate per ambiti specifici
- La soluzione giusta da adottare andrà valutata caso per caso

# eXternal Data Representation (XDR)

- Sun XDR è una soluzione realizzata all'interno dello stack Sun/ONC RPC
- XDR fornisce un insieme di procedure di conversione per trasformare la rappresentazione nativa dei dati in una *rappresentazione esterna* (XDR) e viceversa
- XDR fa uso di uno stream (contenuto in un buffer) che permette di creare un messaggio con i dati in forma XDR
- I dati vengono inseriti/estratti nello/dallo stream XDR uno alla volta, tramite operazioni di serializzazione e/o deserializzazione

## Esempio di serializzazione XDR

```
int i = 260;
char str[80] = "pippo";
XDR *xdrs;
char buf[BUFSIZE]; /* buffer vuoto, da preparare */
...
/* Creazione stream XDR in memoria */
xdrmem_create(xdrs, buf, sizeof(buf), XDR_ENCODE);

/* Inserimento nello stream di un intero, convertito
in formato XDR */
xdr_int(xdrs, &i);

/* Inserimento nello stream di una stringa,
convertita in formato XDR */
xdr_string(xdrs, &str, strlen(str));

/* Scrittura su socket */
write(sd, buf, xdr_getpos(xdrs));
```

# Esempio di deserializzazione XDR

```
int i;
char str[80];
XDR *xdrs;
char buf[BUFSIZE]; /* buffer vuoto, da preparare */
...
/* Lettura da socket */
read(sd, buf, sizeof(buf));

/* Creazione stream XDR in memoria */
xdrmem_create(xdrs, buf, sizeof(buf), XDR_DECODE);

/* Lettura di un intero dallo stream, convertito dal
formato XDR */
xdr_int(xdrs, &i);

/* Lettura di una stringa dallo stream, convertita
dal formato XDR */
memset(str, 0, sizeof(str));
xdr_string(xdrs, &str, sizeof(str)-1);
```

## IDL

- Oltre ai tipi di dati primitivi, per cui fornisce già routine di serializzazione e di deserializzazione, XDR permette di gestire tipi di dati complessi
- In XDR, il formato delle strutture dati è definito attraverso un apposito linguaggio *IDL (Interface Definition Language)* simile al C
- Uso di *IDL compiler* per generare automaticamente le procedure di codifica e decodifica dei dati complessi

# Esempio di IDL (Sun/ONC RPC)

```
/* definisci massima dimensione stringhe */  
const MAXNAMELEN = 255;  
  
/* parametro: nome directory */  
typedef string nametype<MAXNAMELEN>;  
  
/* valore di ritorno: lista di file */  
typedef struct namenode *namelist;  
  
struct namenode {  
    nametype name; /* nome del file */  
    namelist pNext; /* prossimo file */  
};
```

## Altre soluzioni

- CORBA Common Data Representation (CDR)
- ASN.1/X.680 (ITU-T/OSI)
- Soluzioni Web-oriented:
  - Google Protocol Buffers
  - Apache Avro
  - MessagePack
  - Apache Thrift

# Protocolli Testuali

- Nella realizzazione di applicazioni distribuite, l'adozione di protocolli testuali si è spesso rivelata vincente
  - Facilità di testing e debugging
  - Estendibilità
  - Resilienza alla complessità (sistemi di successo di solito evolvono in sistemi più complessi)

Per approfondire si leggano i seguenti saggi di Eric S. Raymond:

How Not To Design a Wire Protocol, <http://esr.ibiblio.org/?p=8254>

The Art of Unix Programming, <http://www.catb.org/~esr/writings/taoup/>

*“When you feel the urge to design a complex binary file format, or a complex binary application protocol, it is generally wise to lie down until the feeling passes.”*

## US-ASCII

- Per decenni, lo standard per la rappresentazione del testo è stato US-ASCII.
- US-ASCII definisce un set di caratteri e una loro rappresentazione in formato binario a 8-bit.
- Solo i 7 bit meno significativi sono effettivamente utilizzati nello standard US-ASCII (127 caratteri). Il bit più significativo di ciascun byte è sempre settato a 0.

# Tabella US-ASCII

<div><div><div>b<sub>7</sub></div><div>b<sub>6</sub></div><div>b<sub>5</sub></div></div><div><div>b<sub>4</sub></div><div>b<sub>3</sub></div><div>b<sub>2</sub></div><div>b<sub>1</sub></div></div></div>					<div><div>0</div><div>0</div><div>0</div><div>0</div><div>1</div><div>0</div><div>1</div><div>1</div><div>0</div><div>0</div><div>1</div><div>0</div><div>1</div><div>1</div></div>								
<div><div>Bits</div><div>↓</div></div>					<div><div>Column</div><div>→</div></div>								
<div><div>Row ↓</div></div>					0	1	2	3	4	5	6	7	
0	0	0	0	0	0	NUL	DLE	SP	0	@	P	`	p
0	0	0	0	1	1	SOH	DC1	!	1	A	Q	a	q
0	0	1	0	0	2	STX	DC2	"	2	B	R	b	r
0	0	1	1	1	3	ETX	DC3	#	3	C	S	c	s
0	1	0	0	0	4	EOT	DC4	\$	4	D	T	d	t
0	1	0	1	1	5	ENQ	NAK	%	5	E	U	e	u
0	1	1	0	0	6	ACK	SYN	&	6	F	V	f	v
0	1	1	1	1	7	BEL	ETB	'	7	G	W	g	w
1	0	0	0	0	8	BS	CAN	(	8	H	X	h	x
1	0	0	1	1	9	HT	EM	)	9	I	Y	i	y
1	0	1	0	0	10	LF	SUB	*	:	J	Z	j	z
1	0	1	1	1	11	VT	ESC	+	;	K	[	k	{
1	1	0	0	0	12	FF	FC	,	<	L	\	l	
1	1	0	1	1	13	CR	GS	-	=	M	]	m	}
1	1	1	0	0	14	SO	RS	.	>	N	^	n	~
1	1	1	1	1	15	SI	US	/	?	O	_	o	DEL

## ISO 8859

- I caratteri della tabella US-ASCII sono sufficienti per la rappresentazione della lingua inglese, ma non per quella di molte altre lingue europee (non supportano accenti, umlaut, ecc.)
- Lo standard ISO 8859 estende US-ASCII utilizzando anche l'ottavo bit, portando così il set di caratteri supportati a 255
- Diverse mappe di caratteri 8859-*n* per coprire le varie lingue (es. per l'italiano si hanno 8859-1 “Latin 1” e 8859-15 “Latin 9”)



# Tabella ISO 8859-15

	-0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-A	-B	-C	-D	-E	-F
0-																
1-																
2-		!	"	#	\$	%	&	'	(	)	*	+	,	-	.	/
3-	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4-	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5-	P	Q	R	S	T	U	V	W	X	Y	Z	[	\	]	^	_
6-	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7-	p	q	r	s	t	u	v	w	x	y	z	{		}	~	
8-																
9-																
A-	ı	ç	£	€	¥	Š	š	©	ª	«	»	¬	-	®	™	
B-	°	±	²	³	Ž	µ	¶	·	¸	¹	º	»	œ	æ	ÿ	
C-	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
D-	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
E-	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
F-	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

## Un mondo post US-ASCII

- Nel ventunesimo secolo, non è più possibile limitare le nostre applicazioni ai set di caratteri US-ASCII o ISO 8859
- Necessità di nuovi standard, che supportino anche:
  - Nuovi alfabeti (cinese, arabo, ebraico, cirillico, ecc.)
  - Caratteri composti
  - Modi di scrittura right-to-left

# ISO 10646

- Lo *standard normativo ISO 10646* definisce lo *Universal Character Set (UCS)*, un set di caratteri che contiene tutti i caratteri universalmente noti
- Ciascun carattere ha un codice a esso associato, e viene rappresentato con una notazione esadecimale come U+12345678
- Separazione tra Basic Multilingual Plane (BMP) (caratteri da U+0000 a U+FFFF) e plane estesi (astrali)
- Supporto a caratteri composti (ad esempio il carattere Ä è rappresentabile tramite la coppia di caratteri U+0041 U+0308)

## Unicode

- *Unicode* è uno *standard implementativo* sviluppato a partire dagli anni '80 da un consorzio di industrie che realizzavano software multilingua
- Inizialmente sviluppato parallelamente e indipendentemente da ISO 10646, si è allineato con quest'ultimo nel 1991
- *Unicode definisce degli standard (UTF-\*/UCS-\*) per l'encoding* dei caratteri dello UCS
  - Assunzione di lavoro (2003): i caratteri UCS hanno codici rappresentabili con al massimo 21 bit (UTF-16 non è in grado di rappresentare caratteri oltre U+10FFFF)

## UTF-32 / UCS-4

- UTF-32 (anche noto come UCS-4) è l'encoding più semplice: ogni carattere viene rappresentato con 4 byte
- Molto spesso, questo tipo di codifica è troppo “onerosa”
  - I caratteri al di fuori del Basic Multilingual Plane sono talmente rari da poter essere ignorati per molte applicazioni
  - I caratteri all'interno del Basic Multilingual Plane sarebbero rappresentabili con 16 bit
- Codifica *non compatibile con US-ASCII*
- Encoding raramente adottato per le comunicazioni e praticamente utilizzato solo all'interno delle librerie di gestione del testo

## UTF-16 / UCS-2

- UTF-16 (che estende il precedente UCS-2) è una **codifica a lunghezza variabile**:
  - 16 bit per i caratteri del Basic Multilingual Plane (U+0000-U+FFFF)
  - 32 bit per gli altri caratteri (U+10000-U+10FFFF)
- Rappresentazione piuttosto compatta
- Encoding standard di Java e Windows
- Codifica *non compatibile con US-ASCII*
- Due diverse versioni: UTF-16LE (little endian) e UTF-16BE (big endian)
  - Ove non specificato, uso di U+FEFF come Byte Order Mark

# UTF-8

- UTF-8 è una **codifica a lunghezza variabile, che supporta tutto il set di caratteri UCS**
  - Ricordiamo che al momento si assume che i caratteri UCS abbiano codici rappresentabili con al massimo 21 bit
- La codifica UTF-8 associa a ciascun carattere una sequenza di byte di lunghezza variabile (da 1 a 4)
- Encoding standard di XML, di JSON, e della maggior parte dei sistemi Unix moderni

## UTF-8 e retrocompatibilità

- UTF-8 è lo strumento messo a disposizione da UCS e Unicode per fornire un certo livello di compatibilità con il passato
  - I caratteri da U+0000 a U+007F sono identici ai caratteri della tabella US-ASCII
  - UTF-8 usa un solo byte per rappresentare i caratteri da U+0000 a U+007F
  - UTF-8 permette stringhe null-terminated
- Quindi, *l'encoding UTF-8 è del tutto compatibile con l'encoding US-ASCII* per il subset di caratteri da quest'ultimo supportati

# Encoding UTF-8

Caratteri UCS	Rappresentazione binaria in UTF-8
Da U+000000 a U+00007F	0XXXXXXXX
Da U+000080 a U+0007FF	110XXXXX 10XXXXXX
Da U+000800 a U+00FFFF	1110XXXX 10XXXXXX 10XXXXXX
Da U+010000 a U+10FFFF	11110XXX 10XXXXXX 10XXXXXX 10XXXXXX

- Velocità di decodifica: il primo byte mi dice quanti altri byte devo leggere per ottenere un carattere completo
- Auto-sincronizzazione: riesco a capire facilmente e velocemente dove inizia un carattere guardando in un intorno di 3 byte dalla posizione corrente

## UTF-8 - Validazione

- **Attenzione!** Poiché in UTF-8 i caratteri hanno dimensione variabile, si deve fare particolare attenzione a verificare che un buffer di memoria non contenga dei caratteri incompleti prima di utilizzarlo!
  - Con una `write` da un file o da una IPC potrei leggere in un buffer solo una parte di una stringa UTF-8, che non contiene tutti i byte che codificano l'ultimo carattere ricevuto!
- **Attenzione!** Oltre a verificare che un buffer non contenga caratteri incompleti bisogna verificare che i dati rappresentati siano effettivamente validi!
  - I caratteri tra U+D800 e U+DFFF sono riservati per UTF-16 e non possono essere usati in UTF-8
  - Rappresentazioni UTF-8 di caratteri con un numero di bit significativi maggiore di 21 vanno scartate

# UTF-8 - Validazione

UTF-8																
	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	0x00	0x01	0x02	0x03	0x04	0x05	0x06	0x07	0x08	0x09	0x0A	0x0B	0x0C	0x0D	0x0E	0x0F
	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
1	0x10	0x11	0x12	0x13	0x14	0x15	0x16	0x17	0x18	0x19	0x1A	0x1B	0x1C	0x1D	0x1E	0x1F
	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
2	0x20	0x21	0x22	0x23	0x24	0x25	0x26	0x27	0x28	0x29	0x2A	0x2B	0x2C	0x2D	0x2E	0x2F
	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47
3	0x30	0x31	0x32	0x33	0x34	0x35	0x36	0x37	0x38	0x39	0x3A	0x3B	0x3C	0x3D	0x3E	0x3F
	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63
4	0x40	0x41	0x42	0x43	0x44	0x45	0x46	0x47	0x48	0x49	0x4A	0x4B	0x4C	0x4D	0x4E	0x4F
	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79
5	0x50	0x51	0x52	0x53	0x54	0x55	0x56	0x57	0x58	0x59	0x5A	0x5B	0x5C	0x5D	0x5E	0x5F
	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95
6	0x60	0x61	0x62	0x63	0x64	0x65	0x66	0x67	0x68	0x69	0x6A	0x6B	0x6C	0x6D	0x6E	0x6F
	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111
7	0x70	0x71	0x72	0x73	0x74	0x75	0x76	0x77	0x78	0x79	0x7A	0x7B	0x7C	0x7D	0x7E	0x7F
	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127
8	+00	+01	+02	+03	+04	+05	+06	+07	+08	+09	+0A	+0B	+0C	+0D	+0E	+0F
	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143
9	+10	+11	+12	+13	+14	+15	+16	+17	+18	+19	+1A	+1B	+1C	+1D	+1E	+1F
	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
A	+20	+21	+22	+23	+24	+25	+26	+27	+28	+29	+2A	+2B	+2C	+2D	+2E	+2F
	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
B	+30	+31	+32	+33	+34	+35	+36	+37	+38	+39	+3A	+3B	+3C	+3D	+3E	+3F
	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
2-byte C	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000	Latin 0000
2-byte D	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400	Cyril 0400
3-byte E	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000	Indic 1000
4-byte F	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400	SSP, SPMA 2400

La tabella (presa da Wikipedia) mostra che alcuni byte in una sequenza di dati codificata in UTF-8 sono sicuramente (rosso) o possibilmente (rosa) rappresentazioni di caratteri non validi.

## wchar\_t, char16\_t e char32\_t

- Storicamente, C si basava su caratteri wide (tipo di dato `wchar_t`) per l'encoding di caratteri multibyte
- Sfortunatamente, la dimensione di `wchar_t` non è standardizzata (Windows usa 16 bit, Unix 32)
  - Solo implementazioni che usano `wchar_t` di dimensioni maggiori di 20 bit sono well behaved e possono dichiarare la macro `__STDC_ISO_10646__` con un valore maggiore o uguale a 200103L
- C11 introduce i tipi `char16_t` e `char32_t` (e le macro `__STDC_UTF_16__` e `__STDC_UTF_32__`) **ma non le funzioni per gestirli**

# Manipolazione stringhe Unicode

- Chiaramente, non possiamo più usare `strlen` per contare il numero di caratteri in una stringa Unicode – nemmeno in quelle UTF-8 NULL-terminated
  - Le funzioni di libreria `str*` sono state progettate assumendo una codifica ASCII che usa un byte per carattere
  - È necessario convertire la stringa UTF-8 a un array di `wchar_t` e usare funzioni `wcs*` (e per piattaforme non compliant come Windows?)
- Inoltre, poiché UCS ammette caratteri composti, non è detto che la lunghezza di una stringa equivalga al numero effettivo di caratteri che essa stamperà a video
  - Uso di `wcswidth` per determinare il numero di colonne richieste per la stampa di un array di caratteri `wchar_t`
- In generale, **meglio usare librerie specifiche per la gestione del testo UTF-8**

Table 7-4: Narrow and Wide String Functions

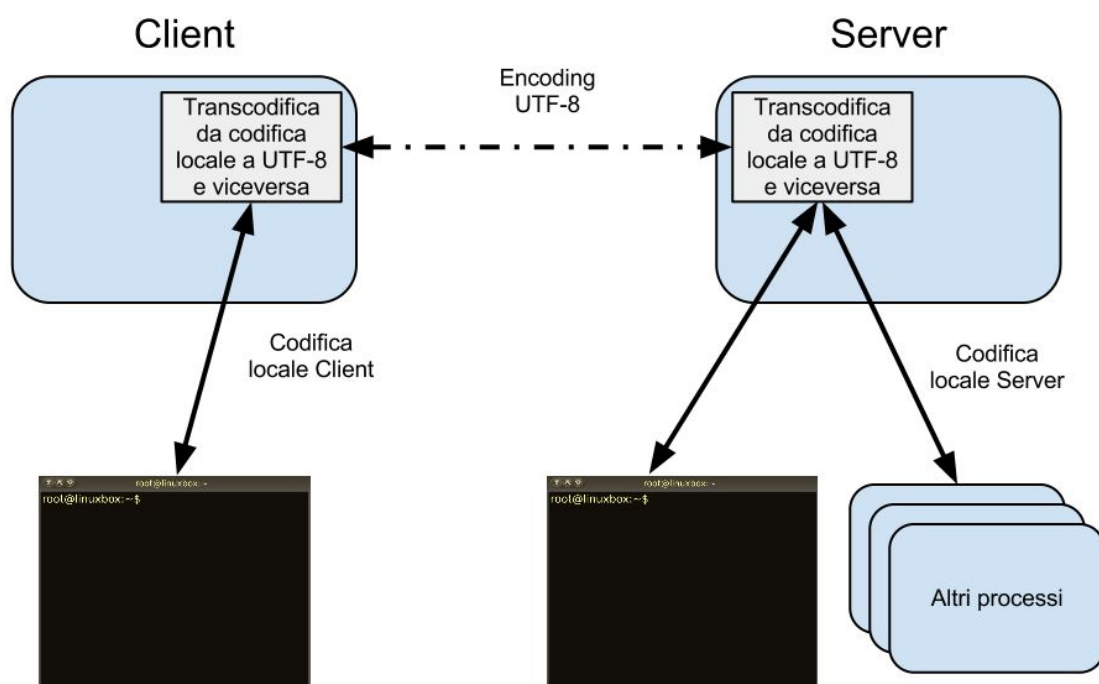
Narrow (char)	Wide (wchar_t)	Description
<code>strcpy</code>	<code>wscpy</code>	String copy
<code>strncpy</code>	<code>wcsncpy</code>	Truncated, zero-filled copy
<code>memcpy</code>	<code>wmemcpy</code>	Copies a specified number of code units
<code>memmove</code>	<code>wmemmove</code>	Copies a specified number of (possibly overlapping) code units
<code>strcat</code>	<code>wscat</code>	Concatenates strings
<code>strncat</code>	<code>wcsncat</code>	Concatenates strings with truncation
<code>strcmp</code>	<code>wscmp</code>	Compares strings
<code>strncmp</code>	<code>wcsncmp</code>	Compares strings to a point
<code>strchr</code>	<code>wcschr</code>	Locates a character in a string
<code>strcspn</code>	<code>wscspn</code>	Computes the length of a complementary string segment
<code>strpbrk</code>	<code>wcspbrk</code>	Finds the first occurrence of a set of characters in a string
<code>strrchr</code>	<code>wcsrchr</code>	Finds the first occurrence of a character in a string
<code>strspn</code>	<code>wcsspn</code>	Computes the length of a string segment
<code>strstr</code>	<code>wcsstr</code>	Finds a substring
<code>strtok</code>	<code>wcstok</code>	String tokenizer (modifies the string being tokenized)
<code>memchr</code>	<code>wmemchr</code>	Finds a code unit in memory
<code>strlen</code>	<code>wcslen</code>	Computes string length
<code>memset</code>	<code>wmemset</code>	Fills memory with a specified code unit

Corrispondenze tra funzioni per manipolazione stringhe con caratteri narrow (ASCII) e wide (`wchar_t`) nel linguaggio C. Courtesy: R. Seacord, “Effective C”, No Starch Press, 2020.

# Unicode e Applicazioni Distribuite

- È bene quindi realizzare applicazioni distribuite che comunichino con l'esterno utilizzando UTF-8
- Nel caso il formato di rappresentazione delle stringhe nella nostra piattaforma di sviluppo non sia UTF-8 (ad esempio su Windows), dobbiamo prevedere una transcodifica da UTF-8 alla codifica locale e viceversa

## Esempio





# Transcodifica

- Per la transcodifica tra UTF-8 e codifica locale:
  - In Java si usano le funzioni di transcodifica fornite da `InputStreamReader` e `OutputStreamWriter` (il secondo parametro del costruttore specifica la codifica “esterna”).
  - In Unix/C lo standard è rappresentato dalla funzione *iconv* fornita dalle librerie di sistema, che però è piuttosto difficile da utilizzare. È pertanto consigliabile valutare l'uso di librerie di più alto livello, come *utf8proc* (<https://julialang.org/utf8proc/>), *libunistring* (<https://www.gnu.org/software/libunistring/>), *glib* (<http://www.gtk.org>), o addirittura *ICU* (<http://site.icu-project.org/>).

## Esempi

- Validazione di una stringa UTF-8 a partire da buffer NULL-terminated con *libunistring*:

```
#include <unistring.h>
#include <inttypes.h>

/* mi assicuro che il buffer sia NULL-terminated */
char buffer[4096];
memset(buffer, 0, sizeof(buffer));
bytes_received = read(socket_fd, buffer, sizeof(buffer)-1);

/* verifico che il messaggio sia UTF-8 valido */
if (u8_check(buffer, bytes_received) != NULL) {
    /* stringa non (interamente) valida, posso provare
       a sanitizzarla o rigettarla: meglio la seconda opzione! */
    fprintf(stderr, "Error 3\n"); fflush(stderr);
    close(socket_fd);
    exit(EXIT_SUCCESS);
}
```

# Esempi

- Validazione di una stringa UTF-8 a partire da buffer NULL-terminated con glib:

```
#include <glib.h>

/* mi assicuro che il buffer sia NULL-terminated */
char buffer[4096];
memset(buffer, 0, sizeof(buffer));
bytes_received = read(socket_fd, buffer,
sizeof(buffer)-1);

if (g_utf8_validate(buffer, NULL, NULL)) {
    /* stringa valida */
} else {
    /* stringa non (interamente) valida, posso provare
    a sanitizzarla o rigettarla: meglio la seconda
    opzione! */
    fprintf(stderr, "Error 3\n"); fflush(stderr);
    close(socket_fd);
    exit(EXIT_SUCCESS);
}
```

# Esempi

- Estrazione sub-stringa UTF-8 valida da un buffer (non necessariamente NULL-terminated) con glib:

```
#include <glib.h>
#include <string.h>

char buff[4096]; int used_buf; char *valid_upto;

used_buf = read(fd, buffer, sizeof(buffer));
if (g_utf8_validate(buff, used_buf, &valid_upto)) {
    /* tutta la stringa è valida */
} else if (valid_upto > buff) {
    /* estraggo porzione di stringa valida e lascio quella non
    valida nel buffer */
    size_t valid_bytes = valid_upto - buff;
    /* devo ricordarmi di deallocare la stringa con free */
    char *valid_str = strdup(buff, valid_bytes);
    memmove(buff, valid_upto, used_buf - valid_bytes);
    used_buf -= valid_bytes;
}
```

# Esempi

- Sanitizzazione di testo UTF-8 in un buffer (non necessariamente NULL-terminated) con glib:

```
#include <glib.h>
```

```
char buff[4096];
```

```
int used_buf;
```

```
char *sanitized_str;
```

```
used_buf = read(socket_fd, buffer, sizeof(buffer));
```

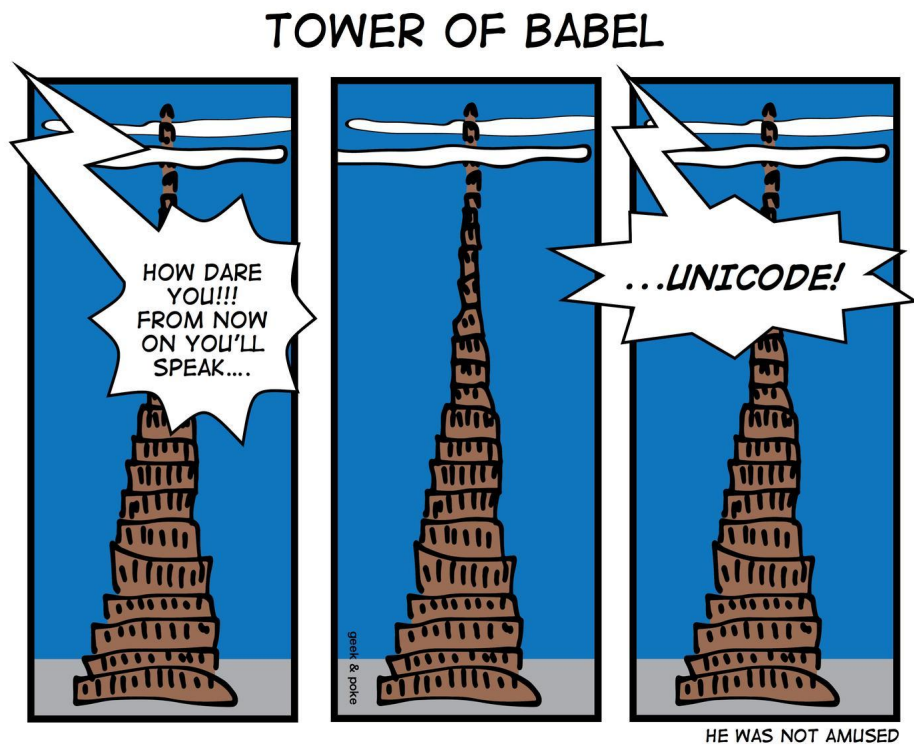
```
/* ricordarsi di deallocare la stringa con free */
```

```
sanitized_str = g_utf8_make_valid(buff, used_buf);
```

**Attenzione!** La sanitizzazione dell'input è un'operazione molto delicata dal punto di vista della sicurezza! Lasciarla a una funzione di libreria di cui non si conosce perfettamente l'implementazione è una pessima idea.

## Per approfondire

- Per ulteriori informazioni su Unicode e UTF-8 si vedano:
  - <http://www.joelonsoftware.com/articles/Unicode.html>
  - <http://utf8everywhere.org/>
  - <http://www.cl.cam.ac.uk/~mgk25/unicode.html>
  - <http://hackaday.com/2013/09/27/utf-8-the-most-elegant-hack/>
  - <https://eev.ee/blog/2015/09/12/dark-corners-of-unicode/>
  - <http://nullprogram.com/blog/2017/10/06/>
  - <http://bjoern.hoehrmann.de/utf-8/decoder/dfa/>



## Dati strutturati

- È possibile scambiare dati di tipo strutturato anche al di sopra di protocolli di tipo testuale, utilizzando standard come XML e JSON
- Di solito, la perdita di performance legata alla trasformazione da rappresentazione binaria a rappresentazione testuale (e viceversa) è più che ripagata dalla flessibilità, dalla robustezza e dalla facilità di debug dei protocolli testuali

# XML

- XML è un linguaggio di descrizione specializzabile per settori specifici (ovverosia un metalinguaggio)
- Standardizzato da W3C, è molto utilizzato anche al di fuori del Web
- Tecnologia sviluppata in ottica machine-oriented, per facilitare la generazione automatica di codice che valida e/o manipola tipi di dati strutturati rigorosamente definiti
- XML è usato sia per la rappresentazione di dati e messaggi scambiati che per la definizione del loro formato

## Esempio XML

```
<? xml version="1.0" encoding="utf-8" ?>
<person>
  <firstName>John</firstName>
  <lastName>Smith</lastName>
  <age>25</age>
  <address>
    <streetAddress>21 2nd Street</streetAddress>
    <city>New York</city>
    <state>NY</state>
    <postalCode>10021</postalCode>
  </address>
  <phoneNumbers>
    <phoneNumber type="home">212 555-1234</phoneNumber>
    <phoneNumber type="fax">646 555-4567</phoneNumber>
  </phoneNumbers>
</person>
```

# XML Schema

- XML Schema è un linguaggio derivato da XML che consente di definire tipi di documento XML contenenti tipi di dati con strutture complesse o non regolari
- XML Schema permette di:
  - definire tipi di dati complessi, basandosi su tipi di dati predefiniti
  - specificare l'ordine in cui gli elementi di un dato devono comparire
  - definire delle regole che specificano il numero di volte che ciascun elemento può o deve comparire
- Uso di XML Namespace per evitare ambiguità

## Esempio XML Schema

```
<?xml version="1.0" encoding="utf-8"?>
<xsd:schema xmlns:xsd="http://www.w3.org/2001/XMLSchema">
  <xsd:element name="Person" type="PersonType"/>
  <xsd:complexType name="PersonType">
    <xsd:sequence>
      <xsd:element name="firstName" type="xsd:string"
        minOccurs="0" MaxOccurs="1"/>
      <xsd:element name="lastName" type="xsd:string"
        minOccurs="1" MaxOccurs="1"/>
      <xsd:element name="age" type="AgeType" />
      ...
    </xsd:sequence>
  </xsd:complexType>
  <xsd:simpleType name="AgeType">
    <xsd:restriction base="xsd:integer">
      <xsd:minInclusive value="0"/>
      <xsd:maxInclusive value="100"/>
    </xsd:restriction>
  </xsd:simpleType>
</xsd:schema>
```

# JSON

- JSON è un formato di rappresentazione dati particolarmente leggero e molto utilizzato sul Web
  - Molto più compatto e significativamente più performante e facile da processare rispetto a XML
- Pensato per applicazioni Web 2.0 e per la manipolazione dei dati in JavaScript
- Compatibile con hash in linguaggio JavaScript

## Esempio JSON

```
{
  "firstName": "John",
  "lastName": "Smith",
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021"
  },
  "phoneNumber": [
    {
      "type": "home",
      "number": "212 555-1234"
    },
    {
      "type": "fax",
      "number": "646 555-4567"
    }
  ]
}
```