

Università di Ferrara

# Architettura di Reti

## Gestione dell'Eterogeneità

Carlo Giannelli

carlo.giannelli@unife.it

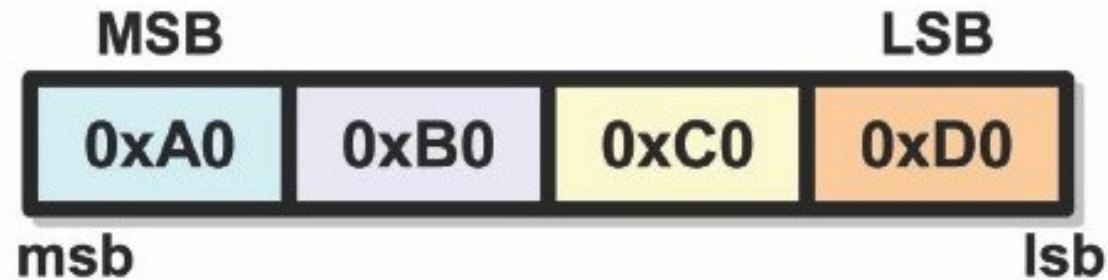
<http://www.unife.it/scienze/informatica/insegnamenti/architettura-reti/>

<http://docente.unife.it/carlo.giannelli>

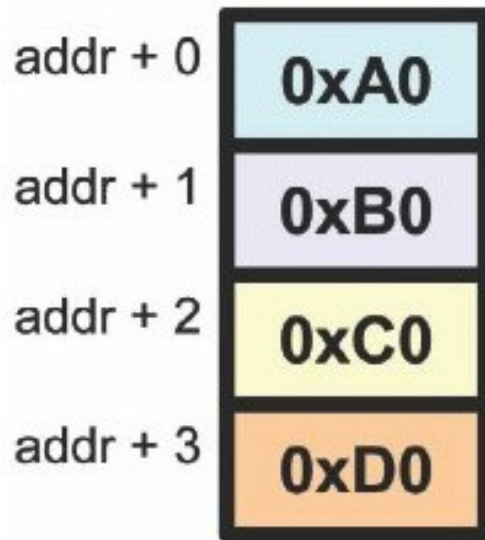
# Il Problema dell'eterogeneità

- Nelle reti di calcolatori vi è un'estrema eterogeneità di sistemi hw/sw
- Client e Server possono eseguire su architetture diverse che usano differenti rappresentazioni dei dati
  - **rappresentazione informazioni**
    - caratteri: US-ASCII, ISO 8859, Unicode...
    - interi: rappresentazione in complemento a 1 o a 2...
  - **dimensione dati**
    - interi: dimensione 4 o 8 byte
    - reali: lunghezza esponente e mantissa
    - caratteri: US-ASCII 1 byte, UTF-8 da 1 a 4 byte, UTF-16 da 2 a 4 byte
  - **ordine byte** all'interno di una parola di più byte
    - little endian vs. big endian
- Necessità di definire una rappresentazione comune dei dati e di implementare meccanismi per gestirla

# Es. little endian o big endian



a)



b)

Big endian



c)

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?!?!

- 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 (marshalling/unmarshalling)

# Esempio di serializzazione XDR

```
int i = 260;
char str[80] = "pippo";
XDR *xdrs; /* stream XDR */
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);
```

# Interface Definition Language - 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 (protobuf)
  - Avro <https://avro.apache.org/>
  - MessagePack <https://msgpack.org/>
  - Apache Thrift <https://thrift.apache.org/>

# Google Protocol Buffers - protobuf

- <https://developers.google.com/protocol-buffers/>
- "Protocol buffers are Google's **language-neutral, platform-neutral, extensible** mechanism for serializing structured data"
- Step principali:
  1. definire i tipi di messaggio "protocol buffer" in un file .proto
  2. compilare il file .proto col compilatore di protobuf per generare codice language-specific
  3. utilizzare il codice generato per serializzare e deserializzare
  4. compilare e linkare tutto assieme

# protobuf – pre-requisiti

- Installazione librerie
  - `sudo apt install protobuf-c-compiler`
  - `sudo apt-get install libprotobuf-dev`
  - `sudo apt-get install libprotobuf-c-dev`

# File .proto per definizione messaggi

- Ad esempio, file di testo “AMessage.proto”
- Ogni campo ha a) un nome, b) un tipo, c) un identificativo univoco numerico
- required → uno ed uno solo di questi campi
- optional → non più di uno di questi campi

```
message AMessage {  
    required int32 a = 1;  
    optional int32 b = 2;  
    required string c = 3;  
}
```

# Tipi di dato

| <b>.proto</b> | <b>C++</b> | <b>Java</b> | <b>Python</b> | <b>C#</b>  | <b>PHP</b>     |
|---------------|------------|-------------|---------------|------------|----------------|
| <b>double</b> | double     | double      | float         | double     | float          |
| <b>float</b>  | float      | float       | float         | float      | float          |
| <b>int32</b>  | int32      | int         | int           | int        | integer        |
| <b>int64</b>  | int64      | long        | int/long      | long       | integer/string |
| <b>bool</b>   | bool       | boolean     | bool          | bool       | boolean        |
| <b>string</b> | string     | String      | str/unicode   | string     | string         |
| <b>bytes</b>  | string     | ByteString  | str           | ByteString | string         |

Dettagli e note in <https://developers.google.com/protocol-buffers/docs/proto>

“[...] A string must always contain UTF-8 encoded or 7-bit ASCII text. [...]”



# Compilazione file .proto (1)

- "protoc-c --c\_out=. amessage.proto" per generare "amessage.pb-c.c" e "amessage.pb-c.h"
- **AMessage**: struct che contiene i campi a, b, has\_b, c

amessage.pb-c.h

```
typedef struct _AMessage AMessage;  
  
struct _AMessage  
{  
    ProtobufCMessage base;  
    int32_t a;  
    protobuf_c_boolean has_b; // b is optional  
    int32_t b;  
    char *c;  
};
```

# Compilazione file .proto (2)

- "protoc-c --c\_out=. amessage.proto" per generare "amessage.pb-c.c" e "amessage.pb-c.h"

**amessage.pb-c.h**

```
void    amessage__init (AMessage *message);
size_t  amessage__get_packed_size (const AMessage *message);
size_t  amessage__pack (const AMessage
                        *message, uint8_t *out);
size_t  amessage__pack_to_buffer (const AMessage *message,
                                   ProtobufCBuffer *buffer);

AMessage * amessage__unpack (ProtobufCAllocator *allocator,
                              size_t len,
                              const uint8_t *data);
void     amessage__free_unpacked (AMessage *message,
                                   ProtobufCAllocator *allocator);
```

# Funzioni/macro amessage.pb-c.h

- Serializzazione

- **AMESSAGE\_\_INIT**: macro per inizializzare correttamente una variabile di tipo AMessage
- **amessage\_\_get\_packed\_size**: funzione per ottenere quale sarà la dimensione di una variabile di tipo AMessage a valle del processo di serializzazione
- **amessage\_\_pack**: funzione che effettua la serializzazione

- Deserializzazione

- **unpack**: funzione che effettua la deserializzazione
- **amessage\_\_free\_unpacked**: funzione che dealloca la memoria allocata da unpack

# Uso del codice generato per serializzare

```
AMessage msg = AMESSAGE__INIT; // AMessage
void *buf;           // Buffer to store serialized data
unsigned len;       // Length of serialized data

msg.a = atoi(argv[1]);
if(argc==3) { msg.has_b = 1; msg.b = atoi(argv[2]); }
msg.c = "ciao";

len = amessage__get_packed_size(&msg);
buf = malloc(len);
// buf will store the serialized version of msg
amessage__pack(&msg, buf);

// Use buf, e.g., sending it via socket */

free(buf); // Free the allocated serialized buffer
```

# Uso del codice generato per deserializzare

```
AMessage *msg;

// Read packed message and size
uint8_t buf[MAX_MSG_SIZE];
size_t msg_len = read_buffer (MAX_MSG_SIZE, buf);

// Unpack the message (implicit memory allocation)
msg = amessage__unpack(NULL, msg_len, buf);
if (msg==NULL){ fprintf(stderr, "error\n");exit(1); }

// Now it is possible to use msg fields
printf("Received: a=%d",msg->a);
if (msg->has_b) printf("  b=%d",msg->b);
printf("  c=%s",msg->c);

// Free the unpacked message
amessage__free_unpacked(msg, NULL);
```

# Funzione di supporto read\_buffer

```
read_buffer (unsigned max_length, uint8_t *out) {  
  
    size_t cur_len = 0;  
    size_t nread;  
  
    while ((nread=fread(out + cur_len, 1, max_length - cur_len,  
                        stdin)) != 0) {  
        cur_len += nread;  
        if (cur_len == max_length) {  
            fprintf(stderr, "max message length exceeded\n");  
            exit(1); }  
    }  
  
    return cur_len;  
}
```

# Compilazione e linking

- Compilazione file generati
  - "gcc -c amessage.pb-c.c"
- Compilazione e linking 1) del codice che utilizza protobuf, 2) dei file generati e 3) delle librerie protobuf
  - gcc -o example.out example.c amessage.pb-c.c -L/usr/lib -lprotobuf-c

# Message boundary

- Protobuf risolve il problema dell'eterogeneità
- Protobuf non risolve il problema del message boundary
  - OK singolo messaggio protobuf in singolo datagram
  - KO più messaggi protobuf in singolo datagram
  - KO più messaggi protobuf consecutivi sullo stesso stream
- Soluzioni per message boundary
  - ZeroMQ <http://zeromq.org/> come soluzione production-ready: **“ZMTP delimits the TCP stream as 'frames'. [...] A frame consists of a flags field, followed by a 'length' field and a frame body of 'length' octets”**. <https://rfc.zeromq.org/spec:15/ZMTP/>
  - Ipotesi semplificative: ack applicativo, prima si invia dimensione del pacchetto e poi pacchetto, buffer di dimensione costante...
  - ...



# Protobuf & ZeroMQ – Client (1)

```
#include <stdio.h>
#include <stdlib.h>
#include <zmq.h>
#include "amessage.pb-c.h"

#define MAX_MSG_SIZE 4096

int main (int argc, const char *argv[])
{
    AMessage msg = AMESSAGE__INIT; // AMessage
    uint8_t *send_buf; // Buffer to store serialized data
    // uint8_t recv_buf[MAX_MSG_SIZE];
    int rc, cc;
    size_t len;
    void *context, *requester;

    // Allow one or two integers
    if (argc != 2 && argc != 3) {
        fprintf(stderr, "Usage: %s a [b]\n", argv[0]);
        exit(EXIT_FAILURE);
    }

    // Create REQREP socket
    context = zmq_ctx_new();
    if (context == NULL) {
        perror("zmq_ctx_new");
        exit(EXIT_FAILURE);
    }
}
```

# Protobuf & ZeroMQ – Client (2)

```
requester = zmq_socket(context, ZMQ_REQ);  
if (requester == NULL) {  
    perror("zmq_socket");  
    exit(EXIT_FAILURE);  
}  
  
rc = zmq_connect(requester, "tcp://localhost:5555");  
if (rc != 0) {  
    perror("zmq_connect");  
    exit(EXIT_FAILURE);  
}  
  
msg.a = atoi(argv[1]);  
if (argc == 3) { msg.has_b = 1; msg.b = atoi(argv[2]); }  
msg.c = "ciao";  
len = amessage__get_packed_size(&msg);  
  
send_buf = malloc(len);  
amessage__pack(&msg, send_buf);
```

# Protobuf & ZeroMQ – Client (3)

```
// invio richiesta contenente messaggio
fprintf(stdout, "Writing %lu serialized bytes\n", len);
cc = zmq_send(requester, send_buf, len, 0);
if (cc < 0) {
    perror("zmq_send");
    exit(EXIT_FAILURE);
}

// Free the allocated serialized buffer
free(send_buf);

zmq_close(requester);
zmq_ctx_destroy(context);

return 0;
}
```

# Protobuf & ZeroMQ – Server (1)

```
#include <stdio.h>
#include <stdlib.h>
#include <zmq.h>
#include "amessage.pb-c.h"

#define MAX_MSG_SIZE 4096

int main(void)
{
    uint8_t recv_buf[MAX_MSG_SIZE];
    AMessage *msg;
    void *context, *responder;
    int cc, rc;
    // char *risposta = "Ok!";

    // Create REQREP socket
    context = zmq_ctx_new();
    if (context == NULL) {
        perror("zmq_ctx_new");
        exit(EXIT_FAILURE);
    }
}
```

# Protobuf & ZeroMQ – Server (2)

```
responder = zmq_socket(context, ZMQ_REP);
if (responder == NULL) {
    perror("zmq_socket");
    exit(EXIT_FAILURE);
}

rc = zmq_bind(responder, "tcp://*:5555");
if (rc != 0) {
    perror("zmq_bind");
    exit(EXIT_FAILURE);
}

// Receive request
cc = zmq_recv(responder, recv_buf, sizeof(recv_buf), 0);
if (cc < 0) {
    perror("zmq_recv");
    exit(EXIT_FAILURE);
}
```

# Protobuf & ZeroMQ – Server (3)

```
// Unpack the message using protobuf-c.
msg = amessage__unpack(NULL, cc, recv_buf);
if (msg == NULL)
{
    fputs("Error unpacking incoming message", stderr);
    exit(EXIT_FAILURE);
}

// Display the message's fields.
printf("Received: a=%d", msg->a); // required field
if (msg->has_b) printf("  b=%d", msg->b); // handle optional field

printf("  c=%s\n", msg->c);

// Free the unpacked message
amessage__free_unpacked(msg, NULL);

zmq_close(responder);
zmq_ctx_destroy(context);

return 0;
}
```

# Protocolli Testuali

- Nella realizzazione di applicazioni distribuite, l'adozione di protocolli testuali si è spesso rivelata vincente
  - facilità di testing e debugging
  - estendibilità
- *“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.”* – Eric Raymond, “The Art of Unix Programming”  
<http://www.catb.org/~esr/writings/taoup/html/>

# 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

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

# ISO 8859

- I caratteri della tabella US-ASCII sono sufficienti per la rappresentazione dell'alfabeto inglese, ma non per quella di molte altre lingue europee. Ad esempio, US-ASCII non supporta accenti (à, è, è...) e umlaut (ä, ö, ü).
- **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- |      | 0001 | 0002 | 0003 | 0004 | 0005 | 0006 | 0007 | 0008 | 0009 | 000A | 000B | 000C | 000D | 000E | 000F |
| 1- | 0010 | 0011 | 0012 | 0013 | 0014 | 0015 | 0016 | 0017 | 0018 | 0019 | 001A | 001B | 001C | 001D | 001E | 001F |
| 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    | ¶    | ·    | ž    | ı    | º    | »    | Œ    | œ    | 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
  - lingue 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 supplementary plane (astral plane)
- 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 ed essenzialmente 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 (2 o 4 byte)**:
  - 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 (da 1 a 4 byte) 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**, quelli della tabella US-ASCII
  - 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 dice quanti altri byte bisogna leggere per ottenere un carattere completo
- **Auto-sincronizzazione:** si riesce 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

|              | _0                       | _1                   | _2                    | _3                        | _4                      | _5                   | _6                    | _7                    | _8                    | _9                    | _A                    | _B                    | _C                     | _D                     | _E                    | _F                   |
|--------------|--------------------------|----------------------|-----------------------|---------------------------|-------------------------|----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|------------------------|------------------------|-----------------------|----------------------|
| 0_           | NUL<br>0000<br>0         | SOH<br>0001<br>1     | STX<br>0002<br>2      | ETX<br>0003<br>3          | EOT<br>0004<br>4        | ENQ<br>0005<br>5     | ACK<br>0006<br>6      | BEL<br>0007<br>7      | BS<br>0008<br>8       | HT<br>0009<br>9       | LF<br>000A<br>10      | VT<br>000B<br>11      | FF<br>000C<br>12       | CR<br>000D<br>13       | SO<br>000E<br>14      | SI<br>000F<br>15     |
| 1_           | DLE<br>0010<br>16        | DC1<br>0011<br>17    | DC2<br>0012<br>18     | DC3<br>0013<br>19         | DC4<br>0014<br>20       | NAK<br>0015<br>21    | SYN<br>0016<br>22     | ETB<br>0017<br>23     | CAN<br>0018<br>24     | EM<br>0019<br>25      | SUB<br>001A<br>26     | ESC<br>001B<br>27     | FS<br>001C<br>28       | GS<br>001D<br>29       | RS<br>001E<br>30      | US<br>001F<br>31     |
| 2_           | SP<br>0020<br>32         | !<br>0021<br>33      | "<br>0022<br>34       | #<br>0023<br>35           | \$<br>0024<br>36        | %<br>0025<br>37      | &<br>0026<br>38       | '<br>0027<br>39       | (<br>0028<br>40       | )<br>0029<br>41       | *<br>002A<br>42       | +<br>002B<br>43       | ,<br>002C<br>44        | -<br>002D<br>45        | .<br>002E<br>46       | /<br>002F<br>47      |
| 3_           | 0<br>0030<br>48          | 1<br>0031<br>49      | 2<br>0032<br>50       | 3<br>0033<br>51           | 4<br>0034<br>52         | 5<br>0035<br>53      | 6<br>0036<br>54       | 7<br>0037<br>55       | 8<br>0038<br>56       | 9<br>0039<br>57       | :                     | ;<br>003A<br>58       | <<br>003B<br>59        | =<br>003C<br>60        | ><br>003D<br>61       | ?<br>003E<br>62      |
| 4_           | @<br>0040<br>64          | A<br>0041<br>65      | B<br>0042<br>66       | C<br>0043<br>67           | D<br>0044<br>68         | E<br>0045<br>69      | F<br>0046<br>70       | G<br>0047<br>71       | H<br>0048<br>72       | I<br>0049<br>73       | J<br>004A<br>74       | K<br>004B<br>75       | L<br>004C<br>76        | M<br>004D<br>77        | N<br>004E<br>78       | O<br>004F<br>79      |
| 5_           | P<br>0050<br>80          | Q<br>0051<br>81      | R<br>0052<br>82       | S<br>0053<br>83           | T<br>0054<br>84         | U<br>0055<br>85      | V<br>0056<br>86       | W<br>0057<br>87       | X<br>0058<br>88       | Y<br>0059<br>89       | Z<br>005A<br>90       | [<br>005B<br>91       | \<br>005C<br>92        | ]<br>005D<br>93        | ^<br>005E<br>94       | _<br>005F<br>95      |
| 6_           | `<br>0060<br>96          | a<br>0061<br>97      | b<br>0062<br>98       | c<br>0063<br>99           | d<br>0064<br>100        | e<br>0065<br>101     | f<br>0066<br>102      | g<br>0067<br>103      | h<br>0068<br>104      | i<br>0069<br>105      | j<br>006A<br>106      | k<br>006B<br>107      | l<br>006C<br>108       | m<br>006D<br>109       | n<br>006E<br>110      | o<br>006F<br>111     |
| 7_           | p<br>0070<br>112         | q<br>0071<br>113     | r<br>0072<br>114      | s<br>0073<br>115          | t<br>0074<br>116        | u<br>0075<br>117     | v<br>0076<br>118      | w<br>0077<br>119      | x<br>0078<br>120      | y<br>0079<br>121      | z<br>007A<br>122      | {<br>007B<br>123      | <br>007C<br>124        | }                      | ~<br>007D<br>125      | DEL<br>007F<br>127   |
| 8_           | •<br>+0<br>128           | •<br>+01<br>129      | •<br>+02<br>130       | •<br>+03<br>131           | •<br>+04<br>132         | •<br>+05<br>133      | •<br>+06<br>134       | •<br>+07<br>135       | •<br>+08<br>136       | •<br>+09<br>137       | •<br>+0A<br>138       | •<br>+0B<br>139       | •<br>+0C<br>140        | •<br>+0D<br>141        | •<br>+0E<br>142       | •<br>+0F<br>143      |
| 9_           | •<br>+10<br>144          | •<br>+11<br>145      | •<br>+12<br>146       | •<br>+13<br>147           | •<br>+14<br>148         | •<br>+15<br>149      | •<br>+16<br>150       | •<br>+17<br>151       | •<br>+18<br>152       | •<br>+19<br>153       | •<br>+1A<br>154       | •<br>+1B<br>155       | •<br>+1C<br>156        | •<br>+1D<br>157        | •<br>+1E<br>158       | •<br>+1F<br>159      |
| A_           | •<br>+20<br>160          | •<br>+21<br>161      | •<br>+22<br>162       | •<br>+23<br>163           | •<br>+24<br>164         | •<br>+25<br>165      | •<br>+26<br>166       | •<br>+27<br>167       | •<br>+28<br>168       | •<br>+29<br>169       | •<br>+2A<br>170       | •<br>+2B<br>171       | •<br>+2C<br>172        | •<br>+2D<br>173        | •<br>+2E<br>174       | •<br>+2F<br>175      |
| B_           | •<br>+30<br>176          | •<br>+31<br>177      | •<br>+32<br>178       | •<br>+33<br>179           | •<br>+34<br>180         | •<br>+35<br>181      | •<br>+36<br>182       | •<br>+37<br>183       | •<br>+38<br>184       | •<br>+39<br>185       | •<br>+3A<br>186       | •<br>+3B<br>187       | •<br>+3C<br>188        | •<br>+3D<br>189        | •<br>+3E<br>190       | •<br>+3F<br>191      |
| 2-byte<br>C_ | •<br>192                 | •<br>193             | Latin<br>0080<br>194  | Latin<br>00C0<br>195      | Latin<br>0100<br>196    | Latin<br>0140<br>197 | Latin<br>0180<br>198  | Latin<br>01C0<br>199  | IPA<br>0200<br>200    | IPA<br>0240<br>201    | IPA<br>0280<br>202    | IPA<br>02C0<br>203    | accents<br>0300<br>204 | accents<br>0340<br>205 | Greek<br>0380<br>206  | Greek<br>03C0<br>207 |
| 2-byte<br>D_ | Cyril<br>0400<br>208     | Cyril<br>0440<br>209 | Cyril<br>0480<br>210  | Cyril<br>04C0<br>211      | Cyril<br>0500<br>212    | Armen<br>0540<br>213 | Hebrew<br>0580<br>214 | Hebrew<br>05C0<br>215 | Arabic<br>0600<br>216 | Arabic<br>0640<br>217 | Arabic<br>0680<br>218 | Arabic<br>06C0<br>219 | Syriac<br>06C0<br>220  | Arabic<br>0740<br>221  | Thaana<br>0780<br>222 | N'Ko<br>07C0<br>223  |
| 3-byte<br>E_ | Indic<br>0800<br>224     | Misc.<br>2000<br>225 | Symbol<br>2000<br>226 | Kana, CJK<br>3000<br>227  | CJK<br>4000<br>228      | CJK<br>5000<br>229   | CJK<br>6000<br>230    | CJK<br>7000<br>231    | CJK<br>8000<br>232    | CJK<br>9000<br>233    | Asian<br>A000<br>234  | Hangul<br>B000<br>235 | Hangul<br>C000<br>236  | Hangul<br>D000<br>237  | PUA<br>E000<br>238    | Forms<br>F000<br>239 |
| 4-byte<br>F_ | SMP, SIP<br>10000<br>240 | 40000<br>241         | 80000<br>242          | SSP, SPUA<br>C0000<br>243 | SPUA-B<br>100000<br>244 | 245                  | 246                   | 247                   | 248                   | 249                   | 250                   | 251                   | 252                    | 253                    | 254                   | 255                  |

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.

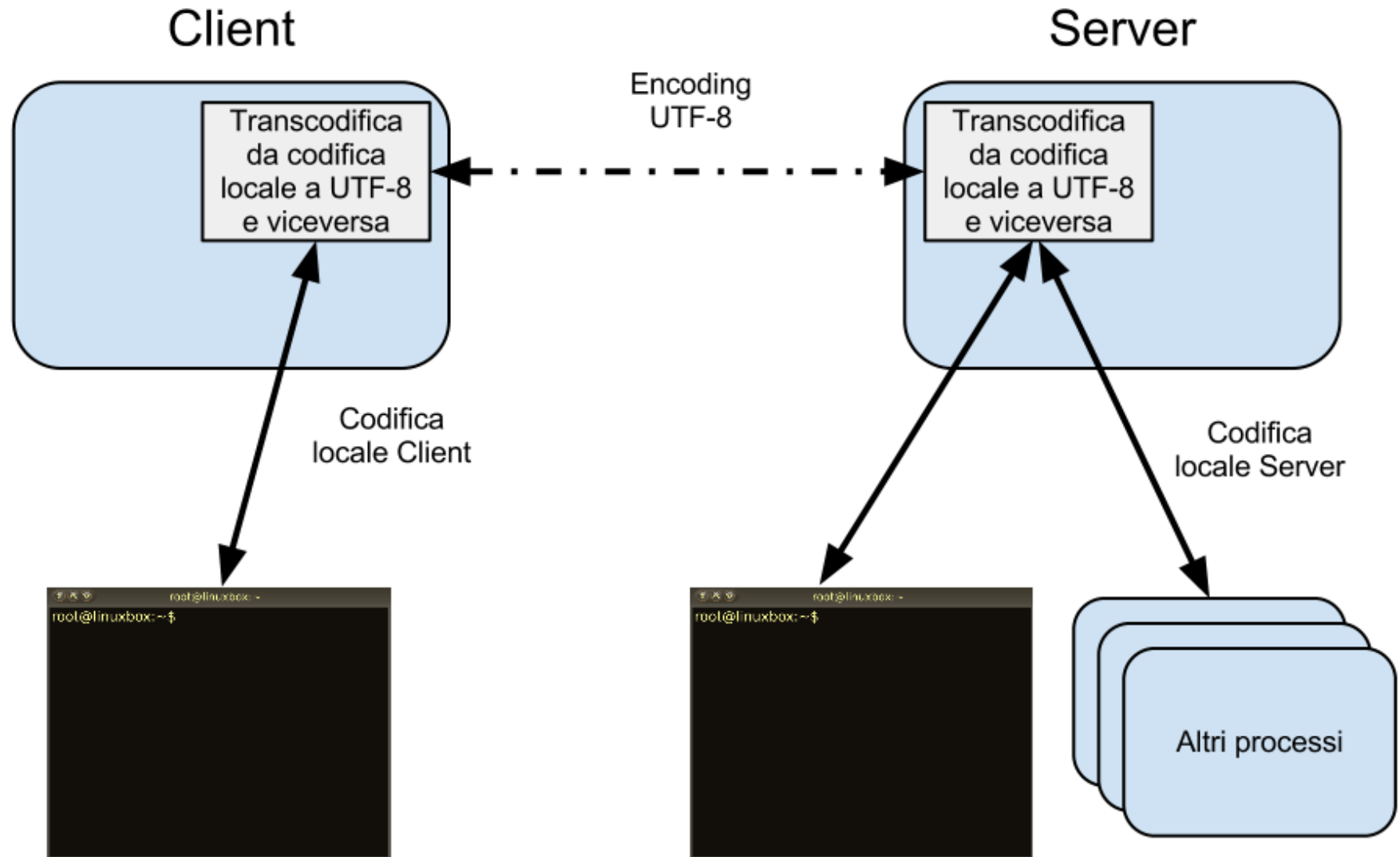
# UTF-8 – Lunghezza stringhe

- **Chiaramente, non possiamo più usare `strlen` per contare il numero di caratteri in una stringa UTF-8 NULL-terminated**
  - le funzioni di libreria `str*` sono state progettate assumendo una codifica US-ASCII che usa un byte per carattere
  - è necessario convertire la stringa UTF-8 a un array di caratteri UTF-32 (`wchar_t` su Unix) e usare funzioni `wcs*`
  - per il caso specifico di `strlen`, al posto di `strlen(str)` si può usare `mbstowcs(NULL, str, 0)`
- 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 UTF-32

# 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>, *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 glib:

```
#include <glib.h>

/* mi assicuro che il buffer sia NULL-terminated */
char buffer[4096];
memset(buffer, 0, sizeof(buffer));
read(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 */
}
```

# 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(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 potrebbe non essere una buona 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" />
      minOccurs="0" MaxOccurs="1"/>
      ...
    </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 nel 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

# 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"
    }
  ]
}
```