

Architettura dell'Informazione

6. - Protocolli a livello applicazione nel WWW

Paolo Milazzo

Dipartimento di Informatica, Università di Pisa

<http://pages.di.unipi.it/milazzo>

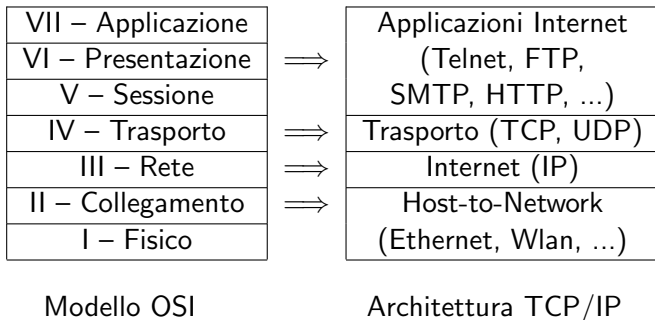
milazzo@di.unipi.it

Master in Turismo e ICT

A.A. 2015/2016

I protocolli di Internet

Il funzionamento di Internet si basa sulla struttura a livelli detta **architettura TCP/IP** (qui confrontata con il modello OSI):



I protocolli a livello applicazione

Applicazioni Internet (Telnet, FTP, SMTP, HTTP, ...)
Trasporto (TCP, UDP)
Internet (IP)
Link (Ethernet, Wlan, ...)

I protocolli di livello applicazione si dividono in due categorie:

- Protocolli di applicazione vera e propria (e.g. SMTP, HTTP, telnet, FTP, etc...) che forniscono il servizio agli utenti finali
- Protocolli di servizio (e.g. DNS) forniscono servizi alle applicazioni usate dagli utenti

Questi protocolli utilizzano TCP/IP come strumento di comunicazione affidabile tra nodi della rete.

I protocolli a livello applicazione

Vediamo un esempio di protocolli a livello applicazione:

- I protocolli per la posta elettronica (SMTP, POP3 e IMAP)

Tutti questi protocolli si scambiano messaggi testuali e definiscono:

- Formato e significato dei messaggi di testo
- Regole per l'invio e la ricezione di tali messaggi

Tutti gli aspetti più tecnici (conversione in binario, instradamento sulla rete, controllo degli errori di trasmissione, ecc...) sono risolti a livello più basso da TCP/IP

Il protocollo SMTP (1)

- SMTP è un protocollo text-based per lo scambio di messaggi di posta elettronica.
- Una connessione SMTP è composta da una apertura, una sequenza di comandi, e una chiusura.
- Ad ogni comando corrisponde una risposta composta da un codice numerico (per le applicazioni) e una stringa leggibile (per gli umani).
 - ▶ MAIL FROM:<Smith@alpha.com>
 - ▶ 250 OK
 - ▶ RCPT TO:<Green@alpha.com>
 - ▶ 550 No such user here

- Esempio di scambio di messaggi tra un client (programma di posta elettronica, in blu) e il server SMTP del mittente (in rosso)
 - ▶ La spedizione di un messaggio SMTP avviene attraverso l'identificazione del mittente (MAIL FROM), del/dei destinatari (RCPT TO), e del messaggio da trasmettere (DATA):
 - ▶ 220 alpha.com Simple Mail Transfer Service Ready
 - ▶ HELO beta.com
 - ▶ 250 alpha.com says: Nice to meet you beta.com
 - ▶ MAIL FROM:<Smith@alpha.com>
 - ▶ 250 OK
 - ▶ RCPT TO:<Green@alpha.com>
 - ▶ 550 No such user here
 - ▶ RCPT TO:<Brown@alpha.com>
 - ▶ 250 OK
 - ▶ DATA
 - ▶ 354 Start mail input; end with <CRLF>.<CRLF>
 - ▶ Blah blah blah....
 - ▶ etc. etc. etc.
 - ▶ .
 - ▶ 250 OK
 - ▶ QUIT
 - ▶ 250 alpha.com Service closing transmission channel

I protocolli POP3 e IMAP4

- SMTP si occupa di recapitare un messaggio di posta elettronica nel disco fisso del server del destinatario
 - ▶ Allo stesso modo in cui un postino recapita la posta nella cassetta della posta
- Per consultare la posta ricevuta (cioè, per andare a ritirare la posta da dentro la cassetta) esistono due protocolli alternativi:
 - ▶ **POP3** permette ad un programma di posta elettronica di accedere alla mailbox posta su un server, consentendo soltanto la possibilità di scaricare e cancellare mail (e poco altro)
 - ▶ **IMAP4** consente a un programma di posta elettronica di effettuare operazioni più “s sofisticate” sulla casella di posta elettronica. Ad esempio: leggere i messaggi senza scaricarli (favorisce l’accesso da diverse postazioni) e creare cartelle sul server in cui suddividere i messaggi ricevuti
- Il protocollo POP3 è sempre meno usato...

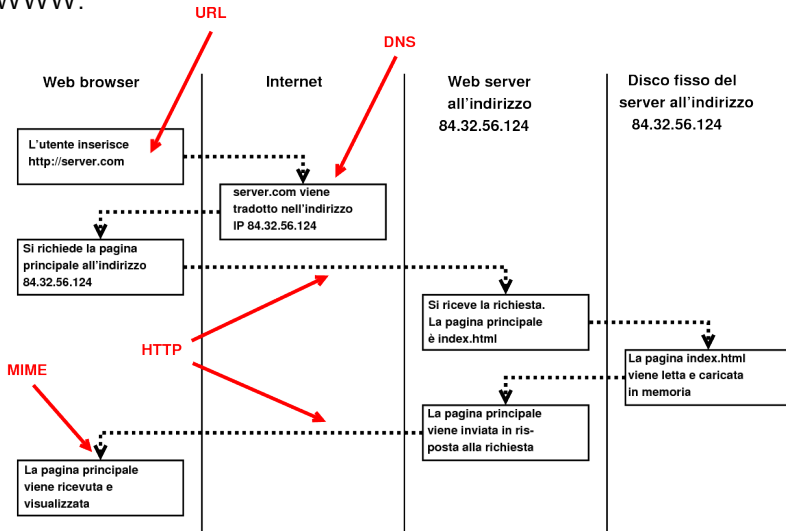
Internet e il WWW (1)

Il WWW usa solo alcuni dei servizi che l'infrastruttura di Internet mette a disposizione. I principali sono:

- L'indirizzamento delle risorse: tramite **URL** (e.g. `http://www.di.unipi.it`)
- La risoluzione degli indirizzi: traduzione di indirizzi testuali tipo `www.di.unipi.it` nei corrispondenti indirizzi IP tramite **DNS**;
- La gestione dei tipi di contenuti: tramite **MIME**;
- Il trasferimento dei contenuti: tramite il protocollo di comunicazione a livello applicazione **HTTP**.

Internet e il WWW (2)

Vediamo dove DNS, MIME e HTTP vengono utilizzati in un tipico utilizzo del WWW:



Indirizzare le risorse: URL

URL (Uniform Resource Locator) è una sintassi standard per la definizione di indirizzi sulla rete;

Un URL consiste di due parti (separate da :)

- una indicazione (detta **scheme**) del protocollo di comunicazione da utilizzare per accedere all'indirizzo descritto
- l'indirizzo vero e proprio (**nome di dominio**)

Tra gli *scheme* abbiamo ad esempio:

- http per gli indirizzi di pagine web
- https per gli indirizzi di pagine web accessibili tramite HTTP sicuro
- mailto per gli indirizzi di posta elettronica
- ftp per gli indirizzi di file da scaricare

Esempi di URL:

- `http://www.di.unipi.it/~milazzo/`
- `mailto:milazzo@di.unipi.it`
- `ftp://fileserv.di.unipi.it:/public/slides.pdf`

Nomi e indirizzi: il DNS (1)

- I computer su Internet sono identificati dal loro indirizzo IP;
- Nel contesto del WWW è molto più comune utilizzare nomi di dominio (e.g. `www.di.unipi.it`) anzichè indirizzi IP;
- Il **Domain Name System (DNS)** è un servizio di naming globale che traduce nomi di dominio in indirizzi IP.

L'utilizzo di nomi di dominio anzichè indirizzi IP offre due vantaggi:

- I nomi di dominio sono più facilmente leggibili e utilizzabili dagli utenti degli indirizzi IP;
- I nomi di dominio offrono un livello di astrazione: se si riorganizza una rete cambiando alcuni indirizzi IP è sufficiente aggiornare il DNS per rendere il cambiamento invisibile agli utenti.

Nomi e indirizzi: il DNS (2)

Il DNS può essere visto come un elenco telefonico globale in cui i nomi anziché essere in ordine alfabetico hanno una struttura gerarchica

- I domini di primo livello sono sigle di nazioni (e.g. `it,uk,fr,...`) o domini generici (e.g. `com,edu,org,...`);
- I domini di livello inferiore sono indicati da destra a sinistra separati da punti: e.g. in `www.di.unipi.it`
 - ▶ `it` è il dominio di primo livello;
 - ▶ `unipi` è il dominio di secondo livello attribuito all'Università di Pisa
 - ▶ `di` è il dominio di terzo livello attribuito al Dipartimento di Informatica;
 - ▶ `www` è il nome convenzionale del computer che ospita il web server in un dominio.

Nomi e indirizzi: il DNS (3)

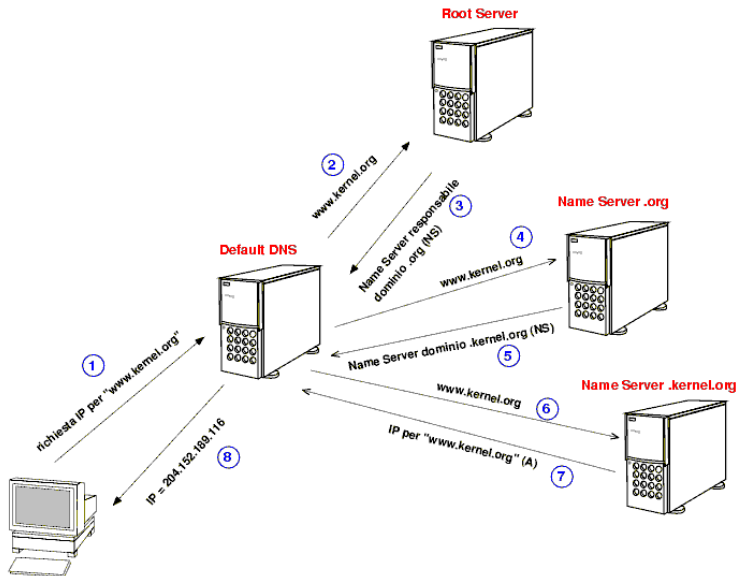
La risoluzione di un nome di dominio tramite DNS funziona come segue:

- Ogni computer collegato a Internet ha un server DNS di riferimento (default DNS, fornito dal fornitore di collegamento a Internet);
- Quando un computer deve risolvere un nome lo inoltra al proprio server DNS di riferimento;
- Se il server DNS non conosce l'indirizzo IP corrispondente al nome richiesto inoltra la richiesta a un *root server* (un server DNS che conosce gli indirizzi dei server DNS responsabili dei domini di primo livello);
- Il server DNS invia ora la richiesta al server responsabile del dominio di primo livello dell'indirizzo richiesto; Tale server fornisce la risposta richiesta, oppure fornisce l'indirizzo del server DNS del dominio di secondo livello di interesse;
- Il procedimento continua finchè la richiesta non riceve risposta.

I vari server DNS possono fare uso di tecniche di caching (memorizzazione) delle informazioni scambiate per evitare successivamente di ripetere tutti questi passaggi

Nomi e indirizzi: il DNS (4)

Un esempio (ricerca dell'indirizzo IP di `www.kernel.org`):



Descrivere il tipo di un contenuto: MIME

Un **tipo MIME** definisce uno specifico tipo per un contenuto associandogli un identificatore di “media type” e un identificatore di “subtype” (rappresentati così: `mediatype/subtype`).

- Il media type dichiara il **tipo** generale del dato trasmesso
- Il subtype dichiara il **formato** specifico del dato trasmesso

Ad esempio:

- il tipo MIME `text/plain` identifica testo semplice;
- il tipo MIME `text/html` identifica testo html;
- il tipo MIME `image/gif` identifica immagini in formato GIF;
- il tipo MIME `video/mpeg` identifica video in formato MPEG.

Il tipo MIME consente al browser di reindirizzare documenti che non sa visualizzare ad applicazioni specifiche.

Questo fa sì, ad esempio, che quando si scarica un file PDF, si apra automaticamente il programma che consente di visualizzare tale documento

Il protocollo HTTP (1)

L'HyperText Transfer Protocol (HTTP 1.1) è un protocollo di comunicazione client-server:

- per la richiesta e la trasmissione di pagine web (e non solo)
- basato su un servizio di trasporto affidabile e orientato alle connessioni (TCP)

Prevede messaggi testuali che consentono:

- A un client (il browser) di richiedere a un server (il sito) di inviare una certa pagina web
- A un server (il sito) di rispondere a una richiesta inviando un documento HTML che descrive la pagina (e/o altro, se serve)

Il protocollo HTTP (2)

Il protocollo HTTP, oltre alle figure del client e del server, può prevedere la presenza di intermediari, tra cui i **proxy**.

Un proxy:

- è un programma che riceve richieste HTTP dai client (opportunamente configurati: l'indirizzo di un eventuale proxy va indicato nel browser) e le inoltra ai server web;
- può applicare tecniche di caching dei contenuti ricevuti dai server per soddisfare autonomamente successive richieste dei client (senza contattare nuovamente i server);
- può applicare filtri ai contenuti ricevuti (e.g. parental control).

Il protocollo HTTP (3)

Il protocollo HTTP:

- E' **stateless** (senza stato): il risultato di una richiesta non può dipendere da richieste precedenti;
- Costringe lo sviluppatore di siti web a trovare metodi alternativi per memorizzare una forma di stato (e.g. cookies).

Un messaggio di richiesta HTTP è costituito da:

- Una riga di richiesta che contiene un “metodo” (descrizione del tipo di richiesta), un URL e l'indicazione della versione di HTTP utilizzata;
- Una sezione di header (contenente informazioni varie)
- Il corpo del messaggio

Il protocollo HTTP (3)

Il metodo può essere uno dei seguenti:

- GET: per richiedere al server l'invio della risorsa indicata nell'URL;
- POST: per inviare informazioni al server;
- HEAD: simile a GET, richiede al server l'invio di informazioni sulla risorsa indicata nell'URL (non la risorsa stessa);
- PUT: per richiedere al server di sostituire una risorsa specificata con un'altra;
- DELETE: per richiedere al server di cancellare una risorsa;
- TRACE: per ricostruire la sequenza di intermediari tra il client e il server;
- OPTIONS: per ottenere informazioni sulle opzioni di comunicazione.

I metodi più comuni (e i soli che approfondiremo) sono GET e POST.

Il protocollo HTTP (4)

Il metodo GET:

- è il metodo più comune, che viene utilizzato ogni volta, ad esempio, che si segue un link in una pagina web o che si inserisce un indirizzo in un browser;
- serve per **richiedere al server** l'invio di una pagina web indicata nell'URL
- nell'header del messaggio:
 - ▶ può porre condizioni sul formato (e.g. lingua della pagina) o le caratteristiche della pagina (dimensioni, data ultimo aggiornamento, ecc..)
 - ▶ può inviare al server informazioni aggiuntive (e.g. browser e sistema operativo utilizzato, pagina di provenienza)
- nel corpo del messaggio:
 - ▶ nulla...

Il protocollo HTTP (5)

Esempio di richiesta che usa il metodo GET:

```
GET www.di.unipi.it/index.html HTTP/1.1
Referer: http://www.google.it/#q=informatica+unipi
User-Agent: Mozilla/4.61 (Macintosh; I; PPC)
Accept-language: en
```

Richiede al server l'invio della pagina `www.di.unipi.it/index.html`. La richiesta è stata generata seguendo un link ottenuto da una ricerca su google (pagina di provenienza), usando il browser Firefox versione 4.61 su un computer Apple. Si richiede inoltre di inviare la versione inglese della pagina.

Il protocollo HTTP (6)

Il metodo POST:

- Viene usato per **trasmettere delle informazioni** dal client al server;
- Esempio tipico di informazione trasmessa: i dati inseriti dall'utente in un modulo (form) in una pagina web;
- I dati trasmessi vengono inseriti nel corpo del messaggio;
- I dati vengono passati alla risorsa (e.g. di solito un programma capace di elaborare i dati ricevuti) specificata dall'URL.

Esempio di richiesta che usa il metodo POST (invio username e password):

```
POST www.di.unipi.it/login.php HTTP/1.1
```

```
User-Agent: Mozilla/4.61 (Macintosh; I; PPC)
```

```
username=milazzo&password=12345678
```

Il protocollo HTTP (7)

Un messaggio di **risposta** a una richiesta di un contenuto è costituito da:

- Una “riga di stato” (status-line) che contiene un codice e una breve descrizione dell’esito della richiesta;
- Una sezione di header (che riporta informazioni varie)
- Il corpo del messaggio di risposta (che contiene di solito il documento HTML richiesto)

Le righe di stato più comuni includono:

200 Ok	metodo eseguito con successo
400 Bad request	errore sintattico nella richiesta
403 Forbidden	richiesta non autorizzabile
404 Not found	URL errato (molto comune)
500 Internal server error	Errore interno (comune con siti web dinamici)

L’header di una risposta a GET include anche il tipo MIME del dato trasmesso

Il protocollo HTTP (8)

Un esempio di richiesta (semplificata):

```
GET www.alpha.com/beta.html HTTP/1.1
```

E relativa risposta:

```
HTTP/1.1 200 OK
```

```
Date: Mon, 28 Jun 2004 10:47:31 GMT
```

```
Server: Apache/1.3.29 (Unix) PHP/4.3.4
```

```
Last-Modified: Mon, 12 Jun 2004 11:32:12 GMT
```

```
Content-Language: it
```

```
Content-Type: text/html; charset=utf-8
```

```
<HTML> ....
```

Il protocollo HTTP (9)

- HTTP trasmette tutti i dati in chiaro;
- Esiste una versione di HTTP per comunicazioni che richiedono crittografia: HTTPS
- HTTPS:
 - ▶ utilizza TCP e SSL (crittografia) per trasmettere i soliti messaggi HTTP;
 - ▶ usa lo schema `https:` al posto di `http:` negli indirizzi.