

# Web Technologies

INFO310

# URLs

<http://www.google.com/search?q=facebook#result>



The diagram illustrates the components of the URL `http://www.google.com/search?q=facebook#result`. Red curly brackets are placed below the URL to group its parts. Below each bracket is a label: 'protocol' for 'http:', 'domain' for 'www.google.com', 'path' for '/search', 'parameters' for '?q=facebook', and 'fragment' for '#result'.

protocol      domain      path      parameters      fragment



# SEO Cheat Sheet: Anatomy of A URL

## 1 SEO-FRIENDLY URL

 <sup>1</sup> <http://store.example.com/topics/subtopic/descriptive-product-name#top>

- <sup>1</sup> Protocol
- <sup>2</sup> Subdomain
- <sup>3</sup> Domain
- <sup>4</sup> Top-Level Domain
- <sup>5</sup> Folders / Paths
- <sup>6</sup> Page
- <sup>7</sup> Named Anchor

### Keyword Priority<sup>1</sup>

Observed Google priority of keyword placement:

- (1) Domain
- (2) Subdomain
- (3) Folder
- (4) Path/Page

### SEO Tips for URLs

- Use **subdomains** carefully. They may be treated as separate entities, splitting domain authority.
- Separate **path** & **page** keywords with hyphens ("-").
- **Anchors** may help engines understand page structure.
- Keyword effectiveness in URLs decreases as URL length and keyword position increases.<sup>1</sup>

<sup>1</sup> SEOMoz correlational data (2009)

## 2 OLD DYNAMIC URL

 <sup>1</sup> <http://www.example.com/index.php?product=1234&sort=price&print=1>

- <sup>1</sup> Protocol
- <sup>2</sup> Subdomain
- <sup>3</sup> Domain
- <sup>4</sup> Top-Level Domain
- <sup>5</sup> Page / File Name
- <sup>6</sup> File Extension
- <sup>7</sup> CGI Parameters

### Popular TLDs<sup>2</sup>

**.com** - commercial  
**.net** - infrastructure  
**.org** - non-profit  
**.edu** - schools  
**.info** - informational  
**.biz** - small business  
**.name** - personal sites

<sup>2</sup> Verisign domain report (2009)

### Popular ccTLDs\*

**.cn** - China  
**.de** - Germany  
**.uk** - United Kingdom  
**.nl** - Netherlands  
**.eu** - European Union  
**.ru** - Russian Federation  
**.ar** - Argentina

\* ccTLD = Country Code TLD

### Popular Extensions

**.htm** - Static HTML  
**.html** - Static HTML  
**.php** - PHP code  
**.asp** - ASP code  
**.aspx** - ASP.NET  
**.cfm** - ColdFusion  
**.jsp** - Java Code

①	②	③	④	⑤	⑥	⑦	⑧
scheme:	//	login.password@	address:	port	/path/to/resource	?query_string	#fragment

- ① Scheme/protocol name
  - ② Indicator of a hierarchical URL (constant)
  - ③ Credentials to access the resource (optional)
  - ④ Server to retrieve the data from
  - ⑤ Port number to connect to (optional)
  - ⑥ Hierarchical Unix path to a resource
  - ⑦ "Query string" parameters (optional)
  - ⑧ "Fragment identifier" (optional)
- } "Authority"

# URL Characters

- Unreserved
  - The alphanumerical upper and lower case character may optionally be encoded:

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z  
a b c d e f g h i j k l m n o p q r s t u v w x y z  
0 1 2 3 4 5 6 7 8 9 - \_ . ~

- Reserved
  - Special symbols must sometimes be percent-encoded:

! \* ' ( ) ; : @ & = + \$ , / ? % # [ ]

- Further details can for example be found in
  - RFC 3986
  - <http://www.w3.org/Addressing/URL/uri-spec.html>
- Source: [https://en.wikipedia.org/wiki/Uniform\\_resource\\_locator](https://en.wikipedia.org/wiki/Uniform_resource_locator)

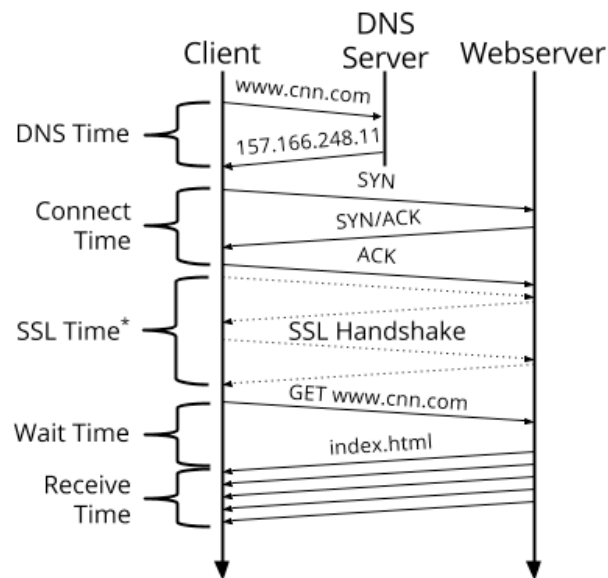
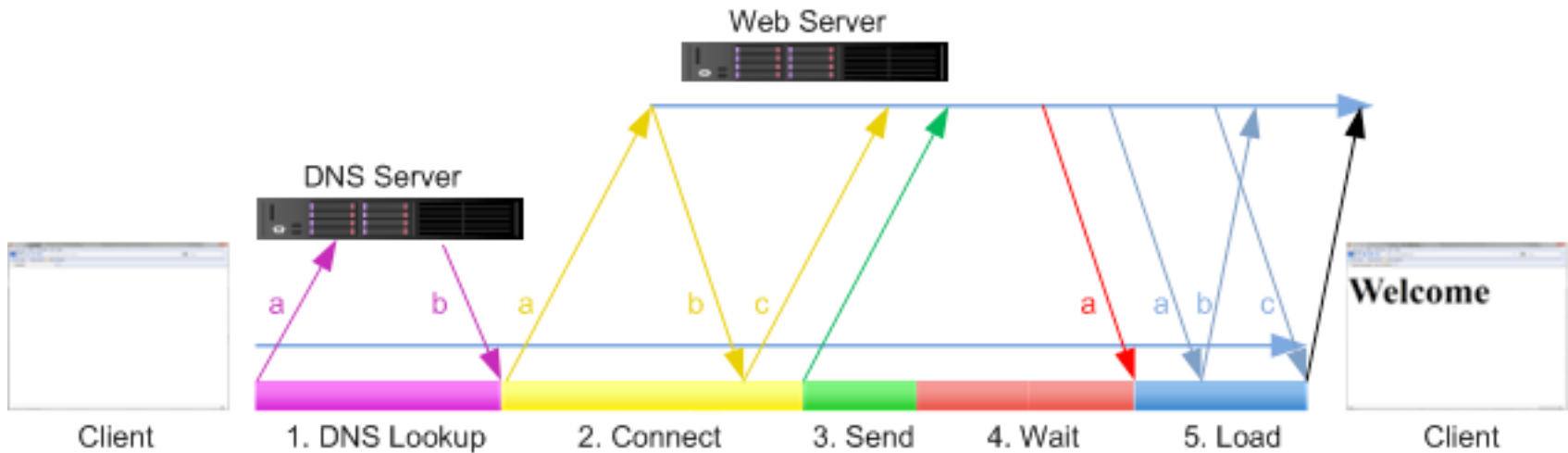
# URL Schemes

- Tons of supported schemes
  - <https://www.iana.org/assignments/uri-schemes/uri-schemes.xhtml>
- Supporting these can lead so some weirdness
- Common ones you may see:
  - file://
  - ftp://
  - http://
  - https://
  - mailto://
  - sms://

# Things can get weird

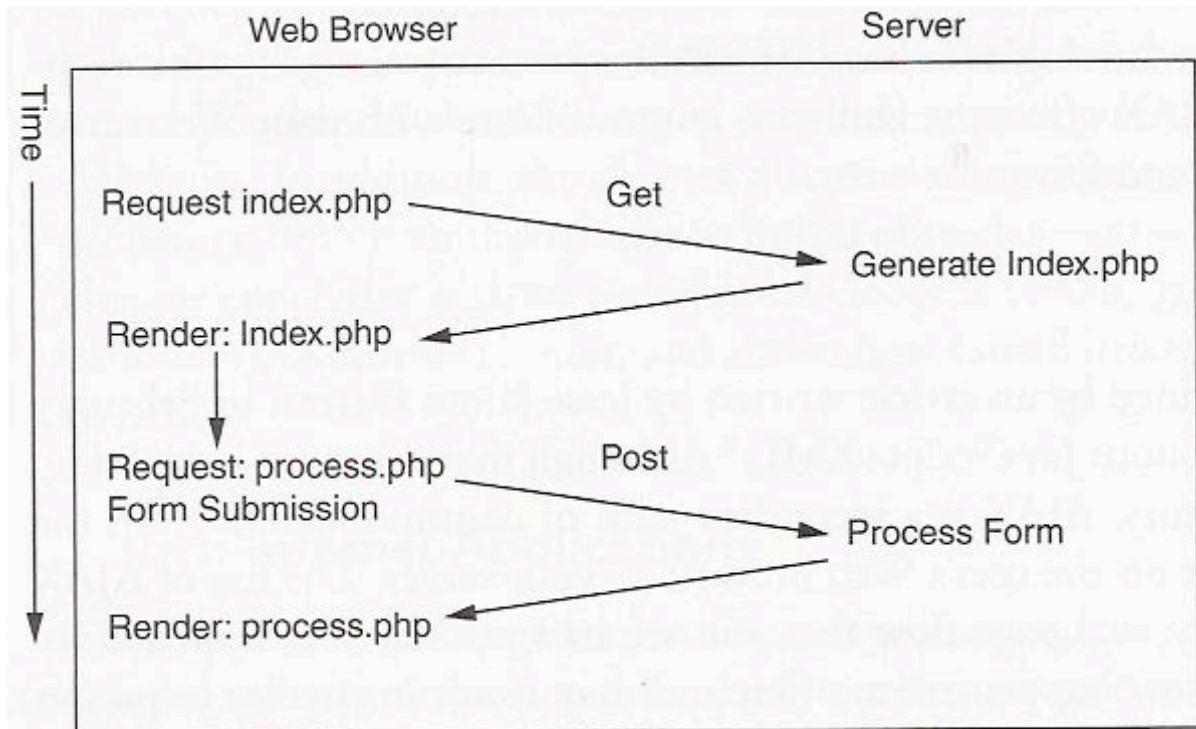
- `http://127.0.0.1/`
  - This is a canonical representation of an IPv4 address.
- `http://0x7f.1/`
  - This is a representation of the same address that uses a hexadecimal number to represent the first octet and concatenates all the remaining octets into a single decimal value.
- `http://017700000001/`
  - The same address is denoted using a 0-prefixed octal value, with all octets concatenated into a single 32-bit integer.
- `http://example.com&gibberish=1234@167772161/`
  - Where do you think this goes?
- `http://example.com\@coredump.cx/`
  - How about this one?
- `http://example.com;.coredump.cx/`
  - And this?
- Source: Tangled Web by Michal Zalewski (pages 26 and 30)

# Browser Requests





# HTTP Requests



**FIGURE 1-1**

Web application request flow

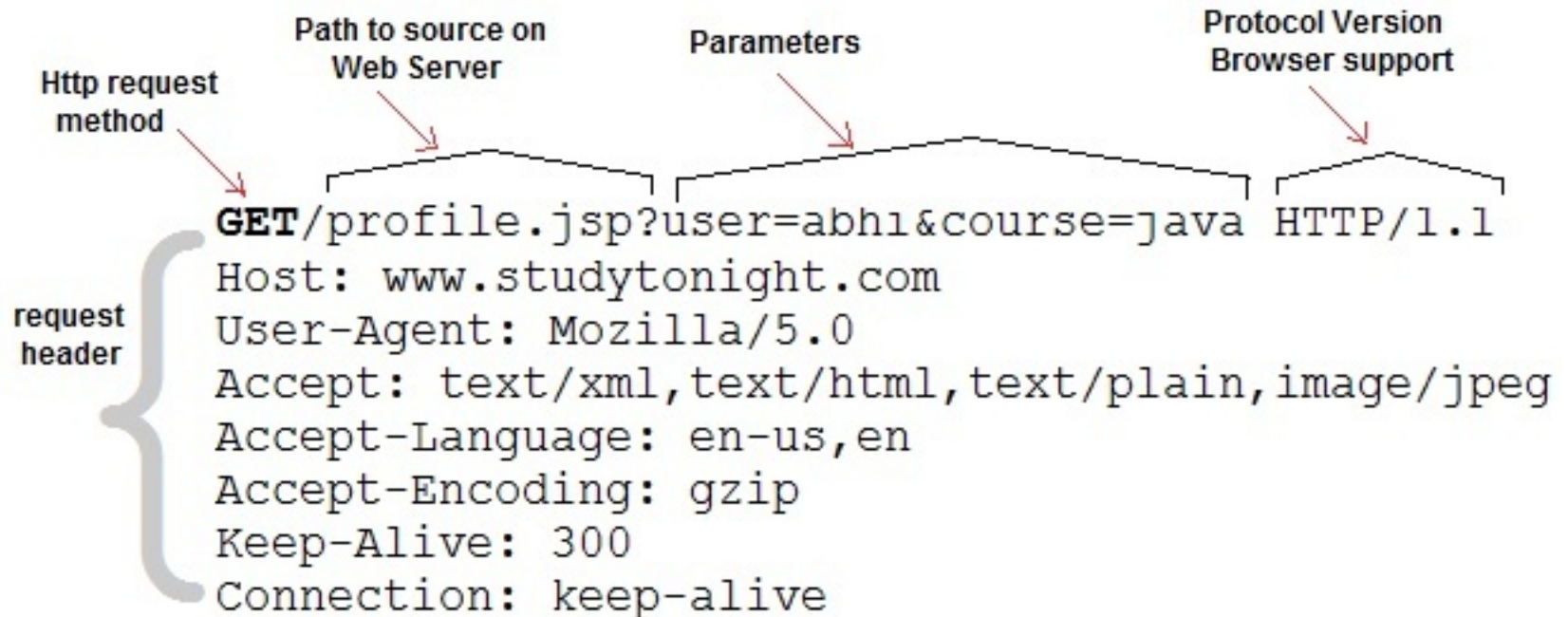
# HTTP Request/Response

```
POST /fuzzy_bunnies/bunny_dispenser.php
HTTP/1.1
Host: www.fuzzybunnies.com
User-Agent: Bunny-Browser/1.7
Content-Type: text/plain
Content-Length: 17
Referer:
http://www.fuzzybunnies.com/main.html
I REQUEST A BUNNY
```

---

```
HTTP/1.1 200 OK
Server: Bunny-Server/0.9.2
Content-Type: text/plain
Connection: close
BUNNY WISH HAS BEEN GRANTED
```

# GET Request



# POST Request

The diagram illustrates the structure of an HTTP POST request. It shows the request line, headers, and the message body. Annotations with arrows point to specific parts of the request:

- Http request method**: Points to the `POST` method in the request line.
- Path to source on Web Server**: Points to the `/profile.jsp` path in the request line.
- Protocol Version Browser support**: Points to the `HTTP/1.1` version in the request line.
- request header**: A bracket on the left side groups the header lines: `Host`, `User-Agent`, `Accept`, `Accept-Language`, `Accept-Encoding`, `Keep-Alive`, and `Connection`.
- parameter inside message body**: Points to the `user=abhi&course=java` string in the message body.

```
POST/profile.jsp HTTP/1.1
Host: www.studytonight.com
User-Agent: Mozilla/5.0
Accept: text/xml,text/html,text/plain,image/jpeg
Accept-Language: en-us,en
Accept-Encoding: gzip
Keep-Alive: 300
Connection: keep-alive
user=abhi&course=java
```

# HTTP Methods

Method	Description
GET	Request to read a Web page
HEAD	Request to read a Web page's header
PUT	Request to store a Web page
POST	Append to a named resource (e.g., a Web page)
DELETE	Remove the Web page
TRACE	Echo the incoming request
CONNECT	Reserved for future use
OPTIONS	Query certain options



# HTTP Headers

- Define the operating parameters of the HTTP transaction
- There are tons “official” ones:
  - [https://en.wikipedia.org/wiki/List\\_of\\_HTTP\\_header\\_fields](https://en.wikipedia.org/wiki/List_of_HTTP_header_fields)
- Colon separated
- Ultimately they can be whatever you want
- No limit on size of name or value

# Cookies

- A small bit of data sent by a web server to a browser that is stored by the browser and sent back with subsequent requests
- Designed to provide a storage mechanism for stateful information and record a user's browsing activity
- Structure
  - Name
  - Value
  - 0+ attributes



# Cookie Attributes

- Domain and Path
  - Defines scope of cookie
- Expires and Max-age
  - Defines when the browser should delete the cookie
- Secure
  - Directs the browser on whether or not to send the cookie over encrypted connection only or not
- HttpOnly
  - Directs the browser on JavaScripts access to the cookie

# Cookies

GET /index.html HTTP/1.1

Host: www.example.org

...

HTTP/1.0 200 OK

Content-type: text/html

Set-Cookie: theme=light

Set-Cookie: sessionToken=abc123;

Expires=Wed, 09 Jun 2021 10:18:14 GMT

...

GET /spec.html HTTP/1.1

Host: www.example.org

Cookie: theme=light; sessionToken=abc123

...