

# Python - Web

# ¿Por qué web?

Disponemos de infinidad de páginas de donde obtener datos, ya sea porque almacenan archivos estructurados, o información en su página.

Dependiendo de cómo esté la información almacenada, existen diferentes técnicas para obtener datos de la web.

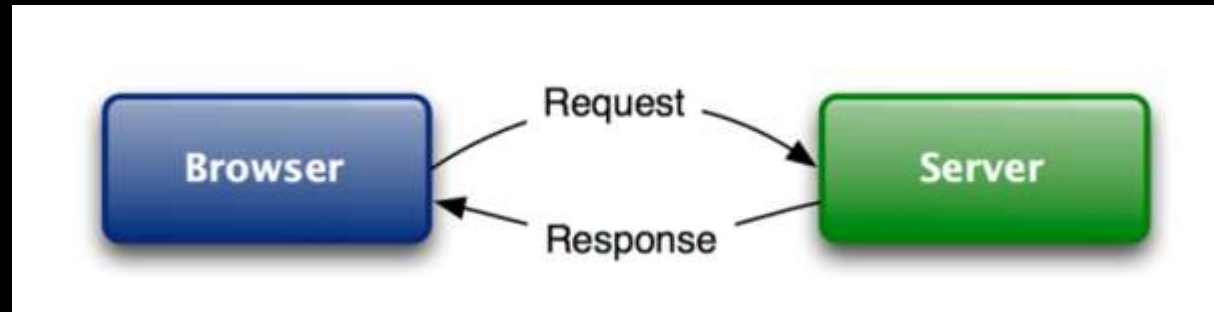


# Técnicas para obtener datos web

- Manual
  - Datos
  - Archivos
- Archivos web con Pandas
- Web Scrapping
- APIs
  - Abiertas
  - Restringidas
  - De pago

# Protocollo HTTP

# Petición respuesta



```
1 GET /home?pageId=c5789534 HTTP/1.1
2 Host: www.buildvsbreak.com
3 User-Agent: Mozilla/5.0
4 Accept: text/html,application/xhtml+xml,application/xml;
5 Accept-Language: en;q=0.5
6 Accept-Encoding: gzip, deflate
7 DNT: 1
8 Connection: keep-alive
```

```
1 HTTP/1.1 200 OK
2 Content-Type: text/html
3 Date: Mon 15 Jul 2013 20:48:49 GMT
4 Server: Apache/2.2.22 (Ubuntu)
5 X-Powered-By: PHP/5.3.10-1 ubuntu3.2
6 Content-Length: 2838
7
8 <!DOCTYPE html>
9 <html xmlns="http://www.w3.org/1999/xhtml">
10 <head>
11 <META http-equiv="Content-Type" content="text/html;
12 <title>Build vs Break Technical Training</title>
13 <meta name="keywords" content="Programming Security
14 <meta name="description" content="We provide trainin
```



# Petición respuesta

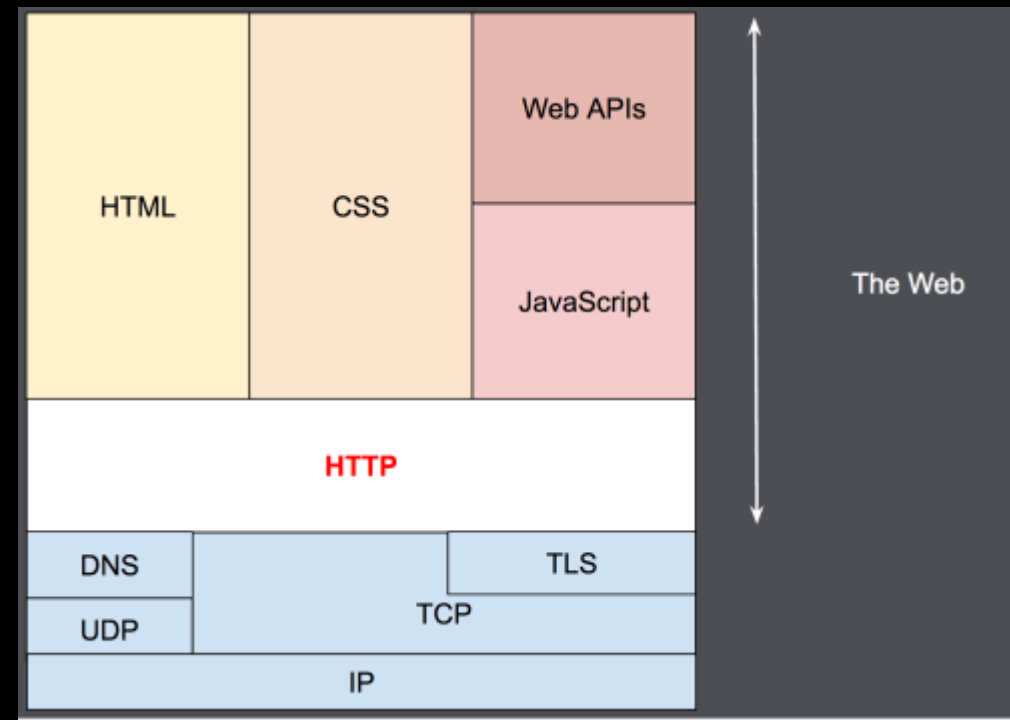
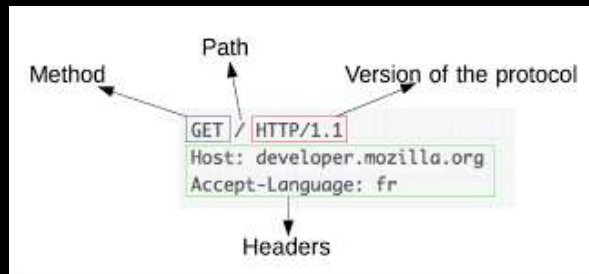


Python server

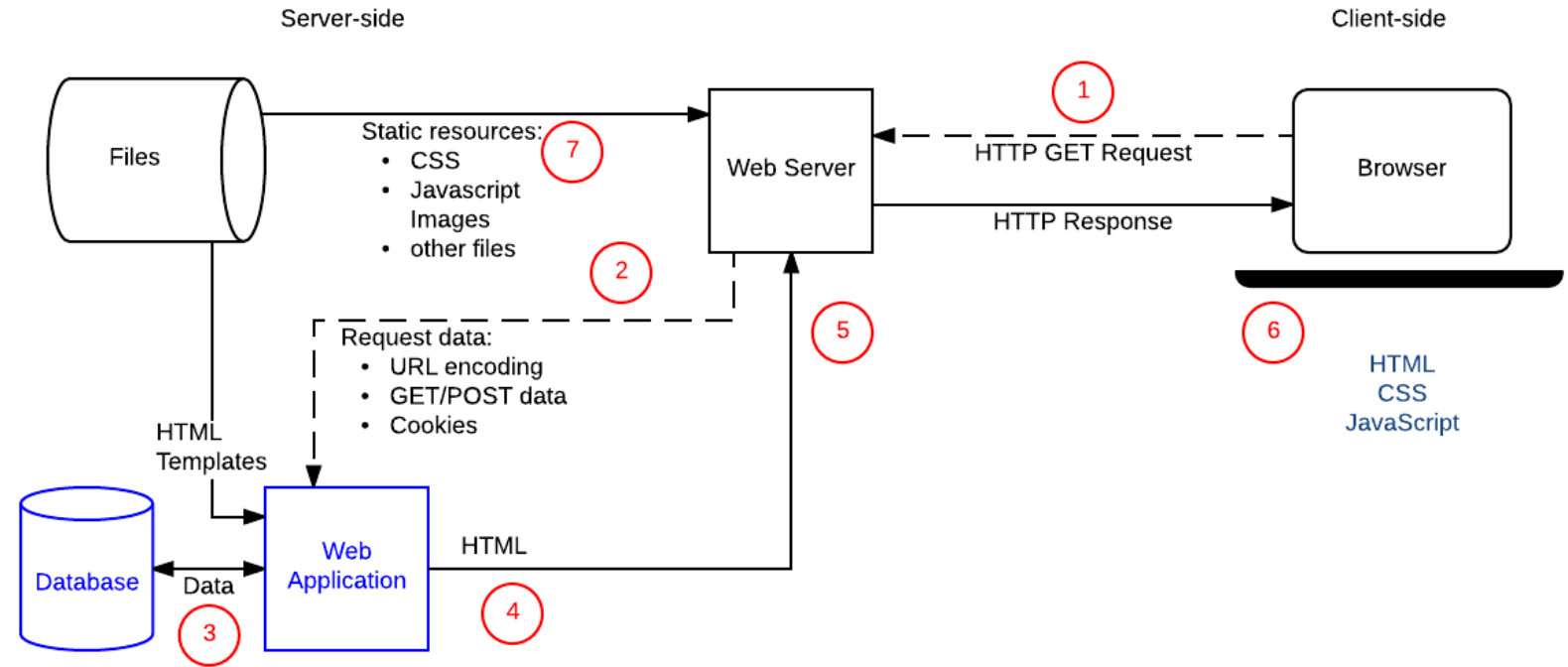


# Protocolo HTTP

Protocolo de aplicación diseñado en los 90s. Se usa para transmisión de datos, documentos, imágenes o vídeos



# Arquitectura cliente servidor



Los mensajes que manda un navegador son peticiones (**request**), y los del servidor respuestas (**responses**)



# HTTP Códigos respuesta

## HTTP Status Codes



## HTTP Status Codes

### Level 200 (Success)

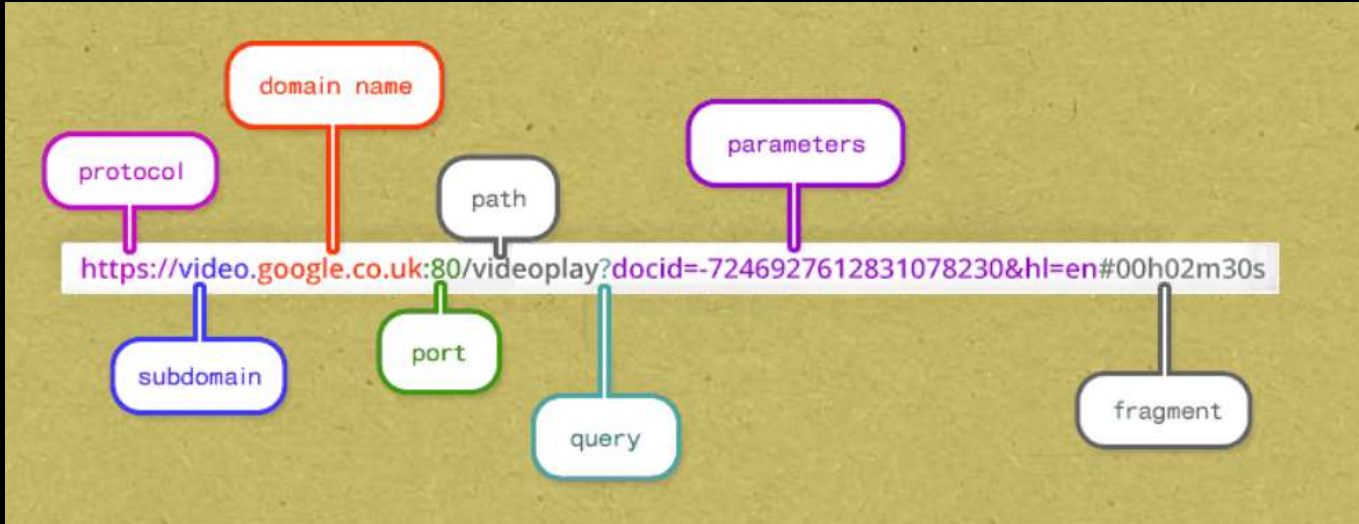
200 : OK  
201 : Created  
203 : Non-Authoritative  
Information  
204 : No Content

### Level 400

400 : Bad Request  
401 : Unauthorized  
403 : Forbidden  
404 : Not Found  
409 : Conflict

### Level 500

500 : Internal Server Error  
503 : Service Unavailable  
501 : Not Implemented  
504 : Gateway Timeout  
599 : Network timeout  
502 : Bad Gateway



La URL no solo sirve para identificar el protocolo de comunicación con el servidor, y la dirección del servidor, sino que también permite establecer ciertos parámetros que usa el servidor para hacer consultas

# URL

## /books

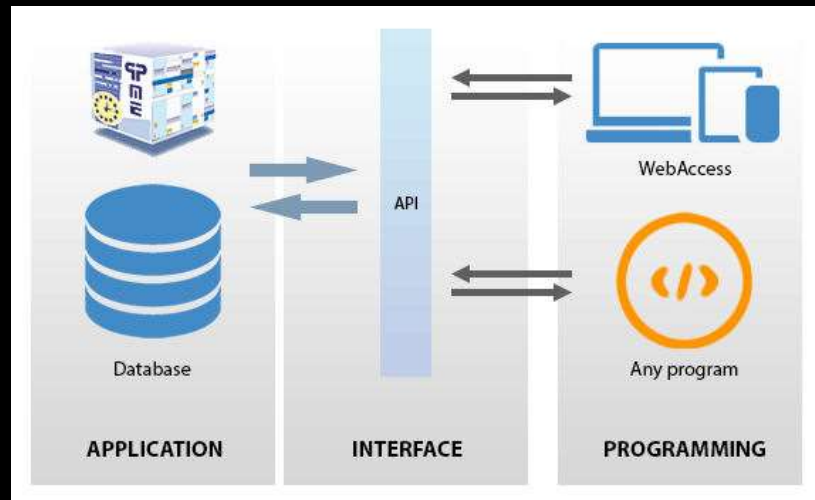
GET	/books	Lists all the books in the database
DELETE	/books/{bookId}	Deletes a book based on their id
POST	/books	Creates a Book
PUT	/books/{bookId}	Method to update a book
GET	/books/{bookId}	Retrieves a book based on their id

# Peticiones HTTP

Dependiendo de la acción que se quiera realizar sobre el servidor, habrá un tipo de petición diferente

# API (Application Programming Interface)

Se trata de una pieza de software intermedia que permite que dos aplicaciones se hablen entre ellas. Las APIs corren en el servidor y tienen un conjunto de operaciones y subrutinas bien definidas. El lado del cliente necesita saber cómo está definida esa API, para poder acceder a los datos



```
    },  
    "formatted_address" : "3100 E Fletcher Ave, Tampa, FL 33613,  
USA",  
    "geometry" : {  
      "location" : {  
        "lat" : 28.0711061,  
        "lng" : -82.4234235  
      },  
      "location_type" : "ROOFTOP",  
      "viewport" : {  
        "northeast" : {  
          "lat" : 28.0724558882915,  
          "lng" : -82.42207451978849  
        },  
        "southwest" : {  
          "lat" : 28.06975711970849,  
          "lng" : -82.42477248029151  
        }  
      }  
    },  
    "place_id" : "ChIJ42XOo5nHwogR3hdzAmex1JU",  
  },  
}
```

[Explicación en 3 minutos](#)



Técnica con la que podemos extraer información de una web de manera automatizada

Entre las aplicaciones prácticas estarían:

- Monitorización de precios de la competencia
- Localización de items o stock en eCommerce
- Detección de cambios en sitios web
- Registrar lanzamientos y novedades
- Analizar los enlaces de un sitio para buscar links rotos

# Web Scrapping

# Herramientas en Python

1. Archivos Web: nymphy pandas
2. APIs: librería request o librería con funciones propias de la API
3. Web Scraping: Selenium o BeautifulSoup



# APIs interesantes

- Idealista: <https://developers.idealista.com/access-request>
- Tripadvisor: <https://developer-tripadvisor.com/content-api/>
- Twitter: <https://developer.twitter.com/es?lang=browser>
- Facebook developers: <https://developers.facebook.com/>
- Youtube: <https://developers.google.com/youtube/v3>
- IA Google Cloud: <https://cloud.google.com/products/ai?hl=es>
- Servicios cognitivos Azure: <https://azure.microsoft.com/es-es/services/cognitive-services/>
- IBM Watson: <https://www.ibm.com/watson/products-services>