

Webscraping w/ Python

Webscraping & API

Sarra Ben Yahia | José Ángel García Sánchez

Introduction

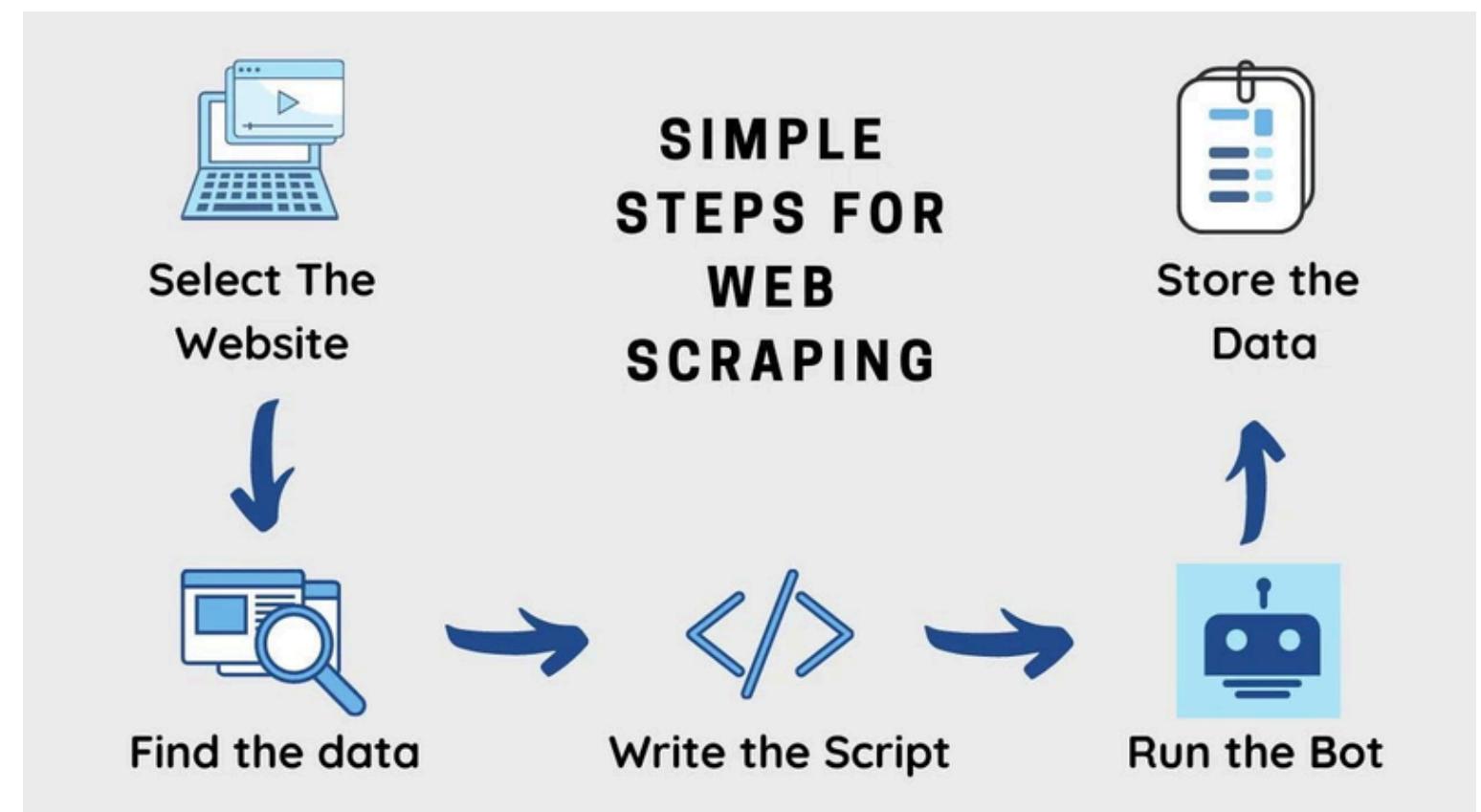
Definition

Web scraping, web harvesting, or web data extraction is data scraping used for extracting data from websites.

Web scraping software may directly access the World Wide Web using the **Hypertext Transfer Protocol (HTTP)** or a **web browser**.

The term typically refers to **automated processes implemented using a bot or web crawler**. It is a form of copying in which specific data is gathered and copied from the web, typically into a central local database or spreadsheet, for later retrieval or analysis.

Example



Webscraping with Python

Several Python frameworks
webscrap



BeautifulSoup Request/BeautifulSoup



Scrapy



Selenium

Request and BeautifulSoup



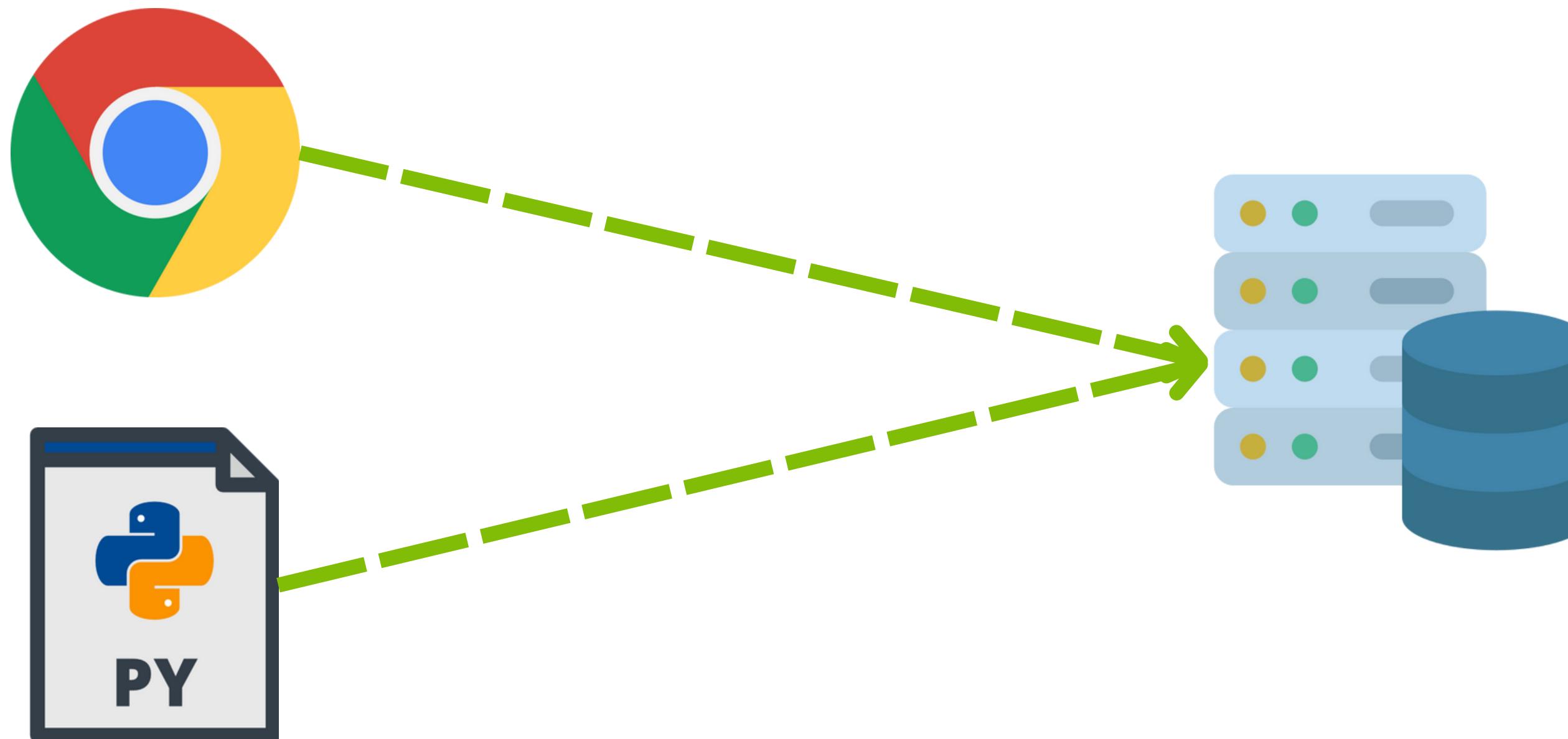
DEMO IN PYTHON : BOOKS



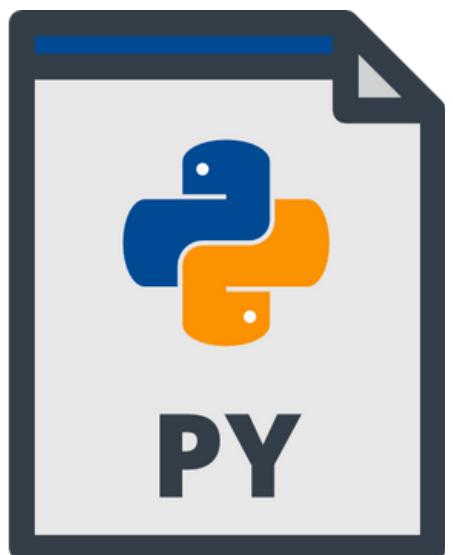
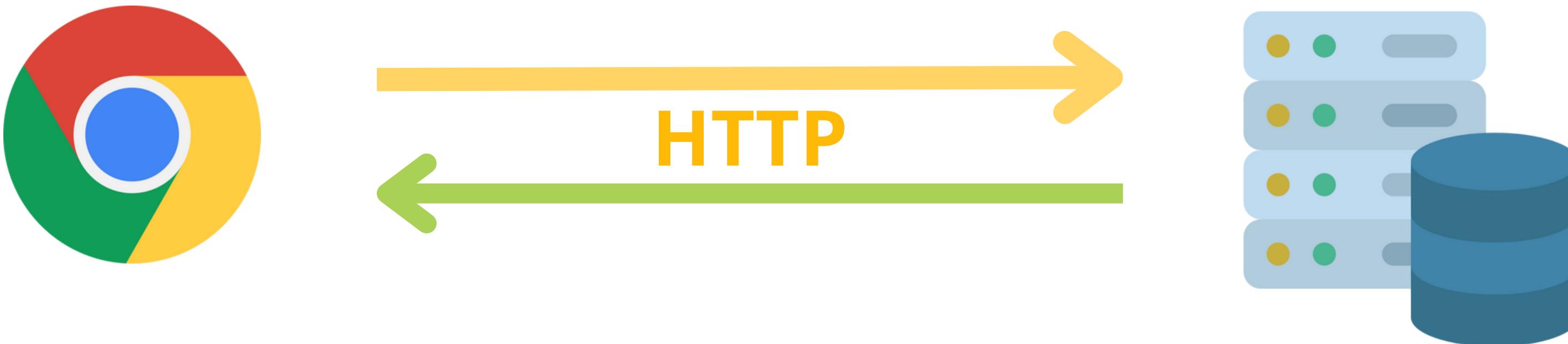
[HTTPS://WWW.SCRAPETHIS SITE.COM/PAGES/SIMPLE/](https://www.scrapethissite.com/pages/simple/)

**SCRAPE THIS WEBSITE AND SAVE RESULT
IN EXCEL**

Scraping Protections

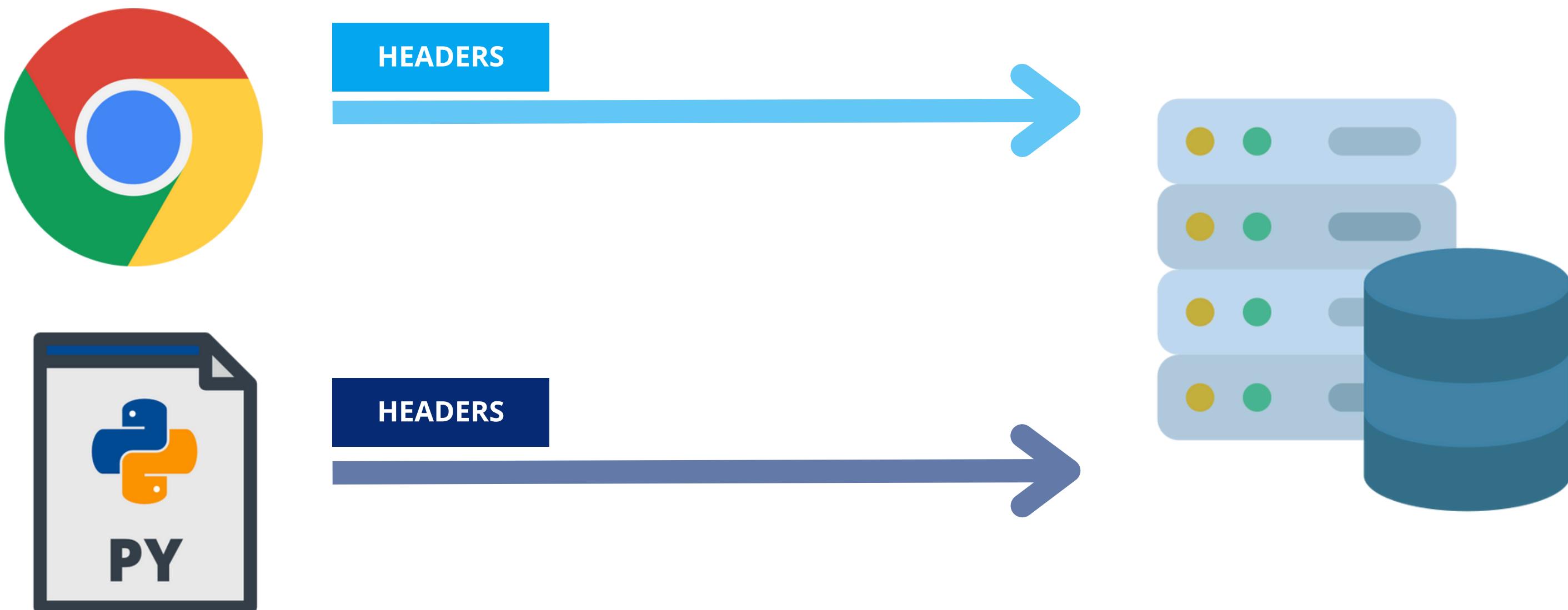


Scraping Protections



- Principal Scraping Protections:**
- **HTTP Headers**
 - **IP/Proxys**
 - **Javascript**

HTTP Headers



HTTP Headers example

General	
Request URL:	https://books.toscrape.com/
Request Method:	GET
Status Code:	304 Not Modified
Remote Address:	35.211.122.109:443
Referrer Policy:	strict-origin-when-cross-origin
Response Headers	
Date:	Sun, 06 Oct 2024 16:42:04 GMT
Etag:	"63e40de8-c85e"
Last-Modified:	Wed, 08 Feb 2023 21:02:32 GMT
Strict-Transport-Security:	max-age=0; includeSubDomains; preload
Request Headers	
:authority:	books.toscrape.com
:method:	GET
:path:	/
:scheme:	https
Accept:	text/html,application/xhtml+xml,application/xml;q=0.9,image/avif,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3;q=0.7
Accept-Encoding:	gzip, deflate, br, zstd
Accept-Language:	fr-FR,fr;q=0.9,en-US;q=0.8,en;q=0.7,ar;q=0.6
Cache-Control:	max-age=0
Dnt:	1
If-Modified-Since:	Wed, 08 Feb 2023 21:02:32 GMT
If-None-Match:	W/"63e40de8-c85e"
Priority:	u=0,i
Sec-Ch-Ua:	"NotA;Brand";v="99", "Google Chrome";v="127", "Chromium";v="127"
Sec-Ch-Ua-Mobile:	?0
Sec-Ch-Ua-Platform:	"macOS"
Sec-Fetch-Dest:	document
Sec-Fetch-Mode:	navigate
Sec-Fetch-Site:	none
Sec-Fetch-User:	?1
Upgrade-Insecure-Requests:	1
User-Agent:	Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/127.0.0.0 Safari/537.36



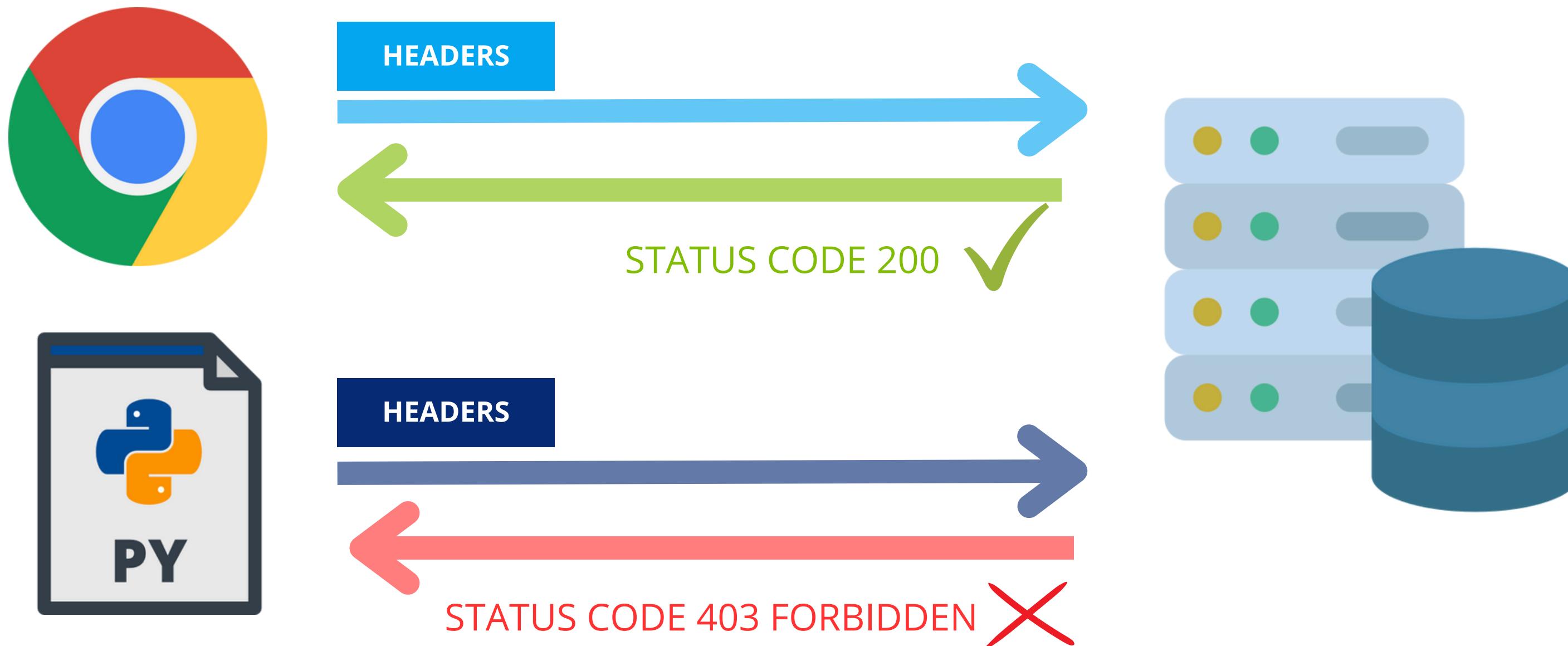
Script

```
✓ _store: OrderedDict([('user-agent', ('User-Agent', 'python-requests/2.28.1')), ('accept-encoding', ('Accept-Encoding', 'gzip, deflate')), ('accept', ('Accept', '*/*')), ('connection', ('Connection', 'keep-alive'))])  
len(): 4
```

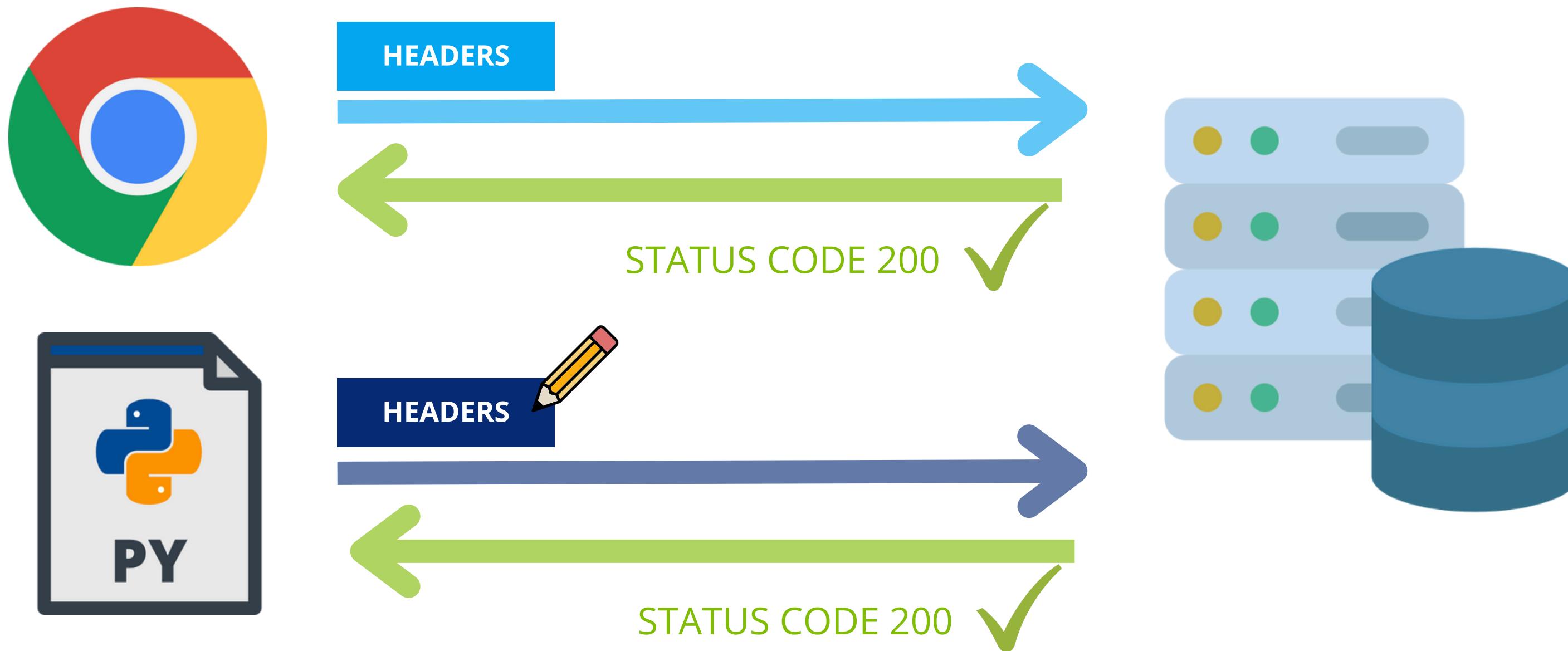


Navigator

HTTP Headers



HTTP Headers



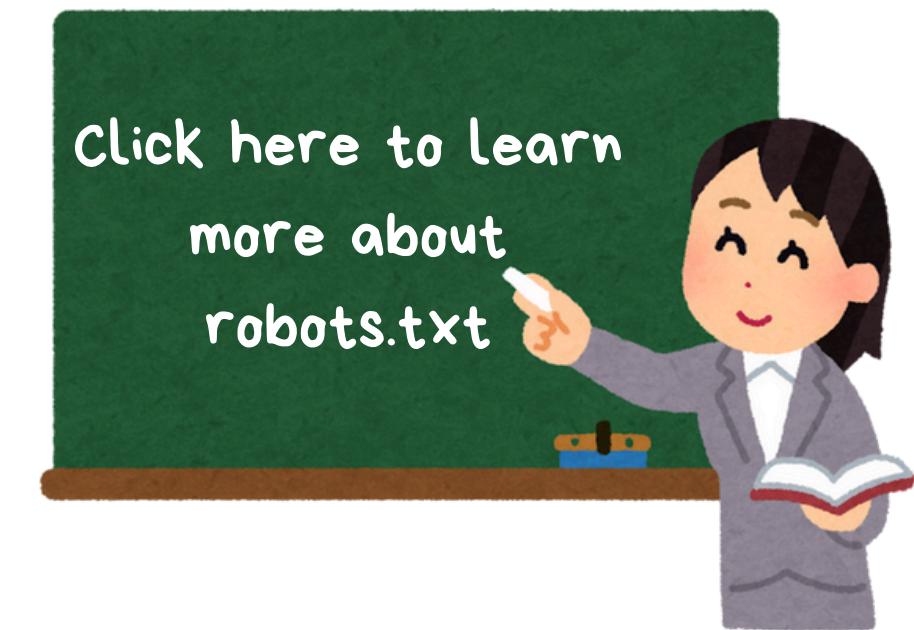
Scraping w/ User-Agents



DEMO IN PYTHON : IMDB

```
User-agent: *
Disallow: /OnThisDay
Disallow: /ads/
Disallow: /ap/
Disallow: /mymovies/
Disallow: /rf/
Disallow: /register/
Disallow: /search/name-text
Disallow: /search/title-text
Disallow: /find
Disallow: /finds
Disallow: /find/
Disallow: /tvshowschedule
Disallow: /updates
Disallow: /watch/_ajax/option
Disallow: /json/video/mon
Disallow: /json/getAdsForMediaViewer/
Disallow: /list/_ls/_ajax
Disallow: /list/_ls/_export
Disallow: /*/_rgs/_mediaviewer/rm*/tr
Disallow: /*/_rgs/_mediaviewer/rm*/tr
Disallow: /*/_mediaviewer/*_tr
Disallow: /title/tte/_mediaviewer/rm*/tr
Disallow: /name/nm/_mediaviewer/rm*/tr
Disallow: /gallery/_rgs/_mediaviewer/rm*/tr
Disallow: /tr/
Disallow: /title/tte/_watchoptions
Disallow: /search/title/?title_type=feature,tv_movie
Disallow: /lists=%21ls538187658,%21ls539867036,%21ls538186
Disallow: /name/nm/_filmtotype/*
Disallow: /user/_ur/_ratings
Disallow: /user/_ur/_checkins
Disallow: /json/*
User-agent: Baiduspider
Disallow: /list/*
Disallow: /user/*
User-agent: bingbot
Disallow: /list/*
Disallow: /user/*
Disallow: /
```

In the robots.txt file of IMDB, we see that there are restrictions for crawlers



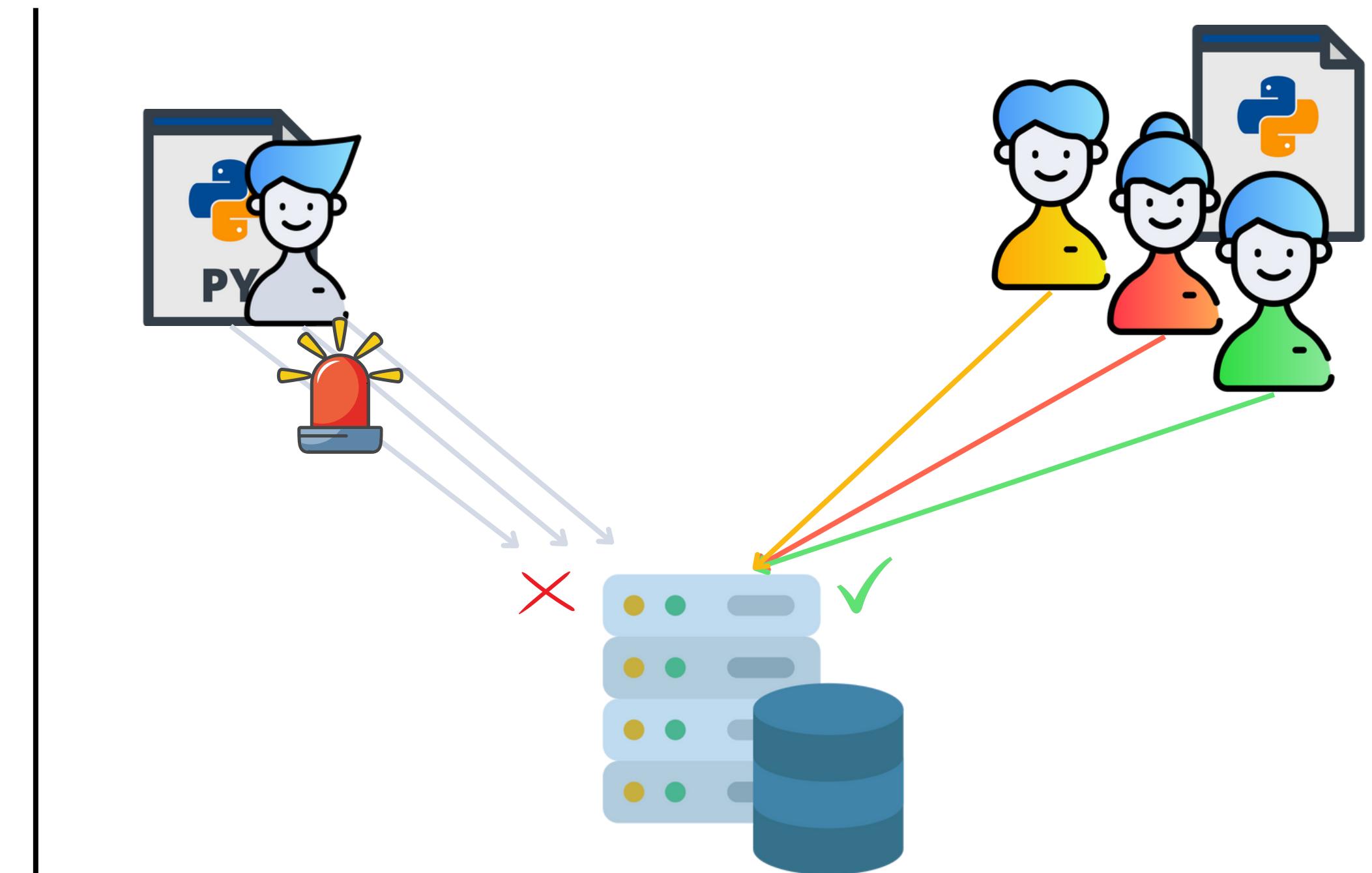
User-Agents Rotation

Why ?

Repeated requests from a **single** user-agent **can trigger server alerts**.

Rotating user-agents makes requests **appear to come from different sources**.

This helps **delay or avoid detection** by server security measures.



User-Agents Rotation



DEMO IN PYTHON : IMDB



EXERCICE

Let's
practice!

[HTTPS://WWW.SCRAPETHISITE.COM/PAGES/ADVANCED/?GOTCHA=HEADERS](https://www.scrapethissite.com/pages/advanced/?gotcha=headers)

**SCRAPE THIS WEBSITE BY SPOOFING
HEADERS**

IP/Proxys limitations

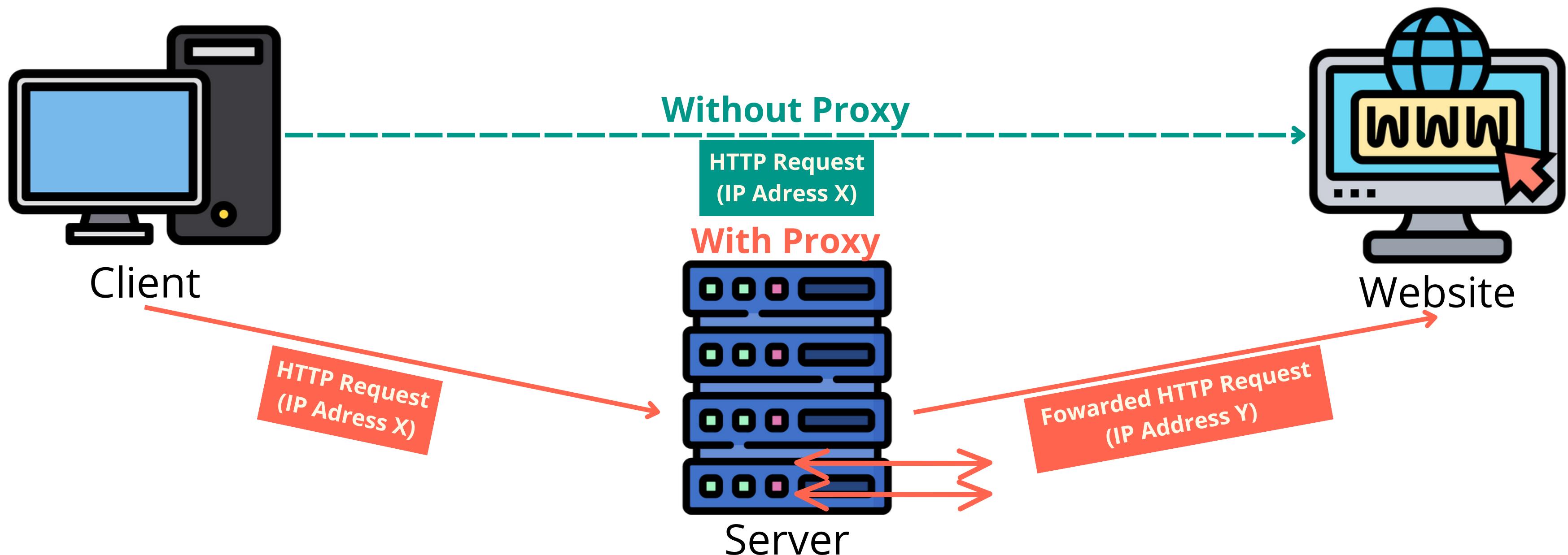
An **Internet Protocol address** (IP address) is a numerical label such as 192.0.2.1 that is assigned to a device connected to a computer network that uses the Internet Protocol for communication.

IP addresses serve two main functions: **network interface identification**, and **location addressing**.



IP/Proxys limitations

A proxy is an **intermediary server** that **forwards requests and responses** between a client and the internet, providing various functions such as **anonymity and content filtering**.



How to get a list of proxies ?

option1



option2



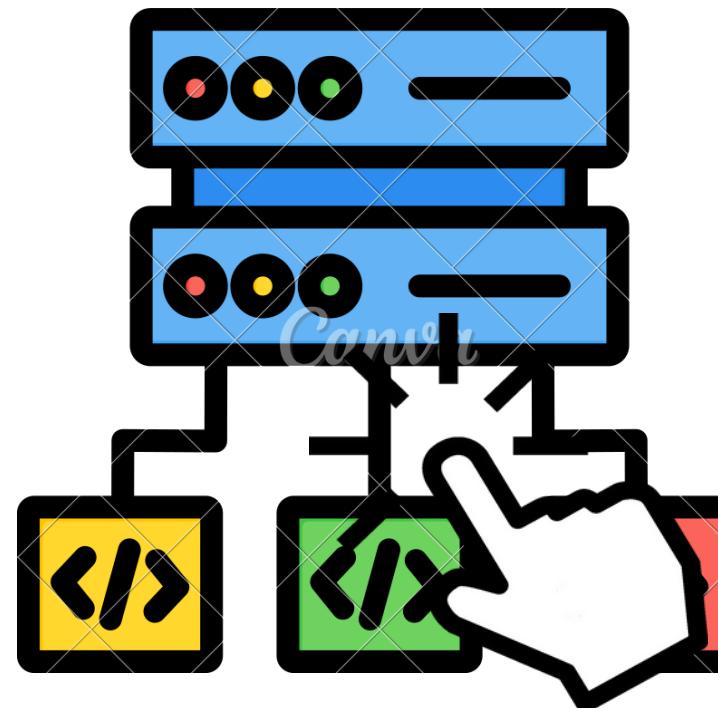
option3



IP address rotation



DEMO IN PYTHON



Legality: Do's and Don'ts

- **Terms of Service (ToS):** Scraping may violate a website's ToS, leading to potential legal consequences.
- **Public vs. Private Data:** Scraping public data is usually more acceptable, but private/protected data can lead to legal issues.
- **Copyright:** Scraping copyrighted content without permission can infringe intellectual property rights.
- **GDPR (Europe):** Scraping personal data must comply with GDPR regulations.
- **Court Rulings:** Legal cases have led to mixed results; public data scraping has sometimes been ruled legal (e.g., HiQ Labs v. LinkedIn).

Good Practice:

- Always check and follow a **website's ToS**.
- Respect **robots.txt files** and avoid scraping **private or protected data**.
- Ensure **compliance with GDPR** or other privacy laws **if scraping personal data**.
- Use **scraping ethically and avoid overloading websites** to prevent harm or misuse.

La violation des droits d'auteurs est constitutive du délit de contrefaçon puni d'une peine de 300 000 euros d'amende et de 3 ans d'emprisonnement (CPI, art. L. 335-2 s.).



Ministère de la Culture

<https://www.culture.gouv.fr> › Thematiques › Files PDF :

La protection par le droit d'auteur fiche N° 1

**END OF
1ST
LESSON**

Introduction to scrapy



- What is **Scrapy**? → A **fast, open-source web crawling framework** for Python.
- **Key Features:**
 - **Efficient:** Handles large-scale scraping with built-in mechanisms for speed and data integrity.
 - **Robust:** Built-in support for handling retries, redirects, and handling of cookies.
 - **Asynchronous:** Uses Twisted for non-blocking requests, making it capable of scraping multiple pages concurrently.

- Core Concepts:
 - **Spiders:** Define how to crawl and extract data from websites.
 - **Items:** Represent the data scraped from the pages.
 - **Pipelines:** Post-process and store scraped data (e.g., cleaning, saving to databases).
 - **Middlewares:** Customize request and response processing (e.g., handling proxies, user agents).
 - **Settings:** Control the behavior of Scrapy (e.g., request delays, concurrency limits).

Introduction to scrapy



DEMO IN PYTHON : SCRAPY PROJECT

Structure of scrapy projects

Script	Function	Example of use
<code>scrapy.cfg</code>	Global configuration of the project	Define the default settings module, configure deployment
<code>settings.py</code>	Project configuration parameters	Set bot name, configure robots.txt obedience, set delays between requests
<code>items.py</code>	Definition of data structures for scraped items	Create classes for products with fields like title, price, description
<code>pipelines.py</code>	Processing of extracted data	Clean data, convert types (e.g., string to float for prices), store in a database
<code>middlewares.py</code>	Customization of request/response behavior	User-Agent rotation, proxy management, modification of HTTP headers
<code>spiders/ (folder)</code>	Contains spiders that define how to scrape sites	Create a spider to extract data from an e-commerce site, define navigation and extraction rules



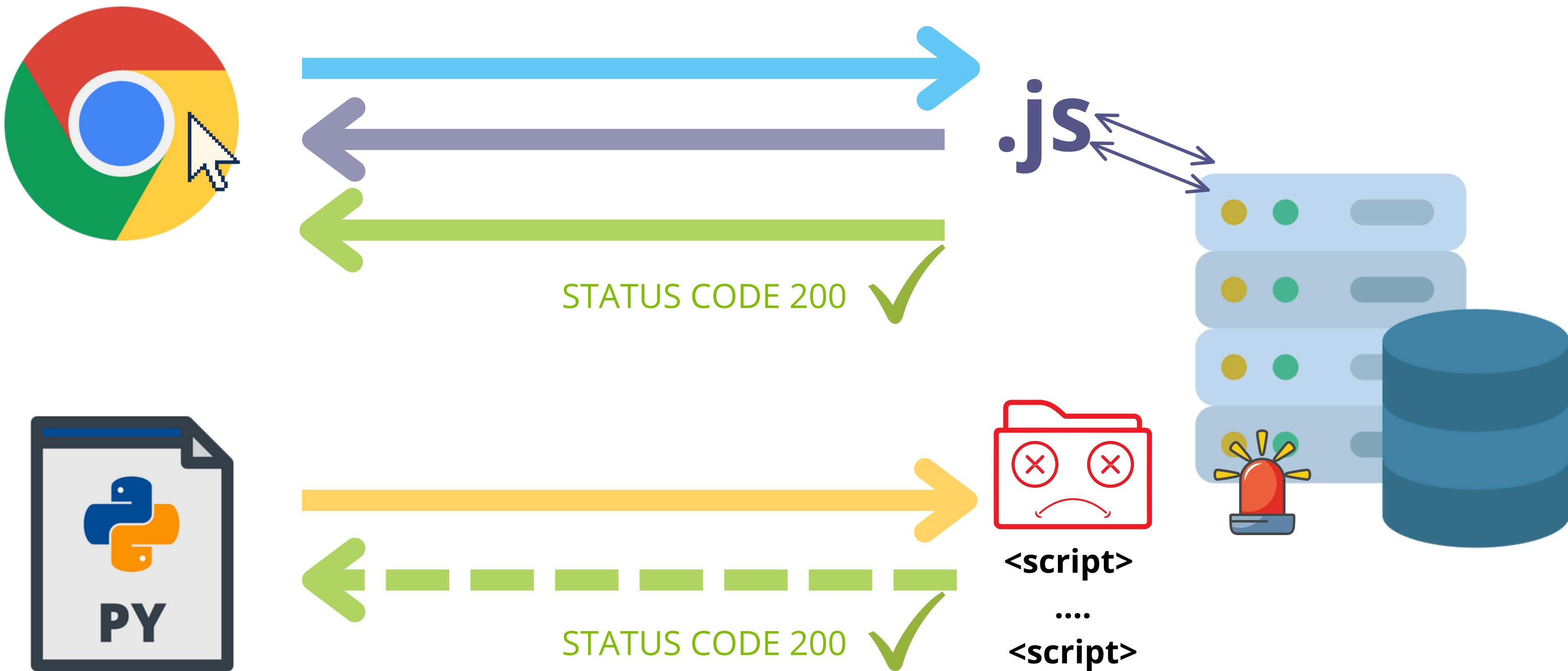
[HTTPS://WWW.SCRAPETHISITE.COM/PAGES/AJAX-JAVASCRIPT/](https://www.scrapethissite.com/pages/ajax-javascript/)
SCRAPE THIS WEBSITE

Principal terminal commands you'll need

Command	Description
<code>scrapy startproject <nom_projet></code>	Create a new scrapy project
<code>scrapy genspider <nom_spider> <domaine></code>	Create a new spider
<code>scrapy crawl <nom_spider></code>	Execute a spider
<code>scrapy shell <url></code>	Open an interactive Scrapy console to do some tests
<code>scrapy list</code>	List all the available spiders in your project
<code>scrapy fetch <url></code>	Fetch the HTML of a page and print it in terminal
<code>scrapy view <url></code>	Open the page in the navigator as Scrapy sees it
<code>scrapy parse <url></code>	Parse the HTML page and print results
<code>scrapy version</code>	Print Scrapy's verison
<code>scrapy bench</code>	Small benchmark to evaluate performances

**END OF
2ND
LESSON**

Javascript problems

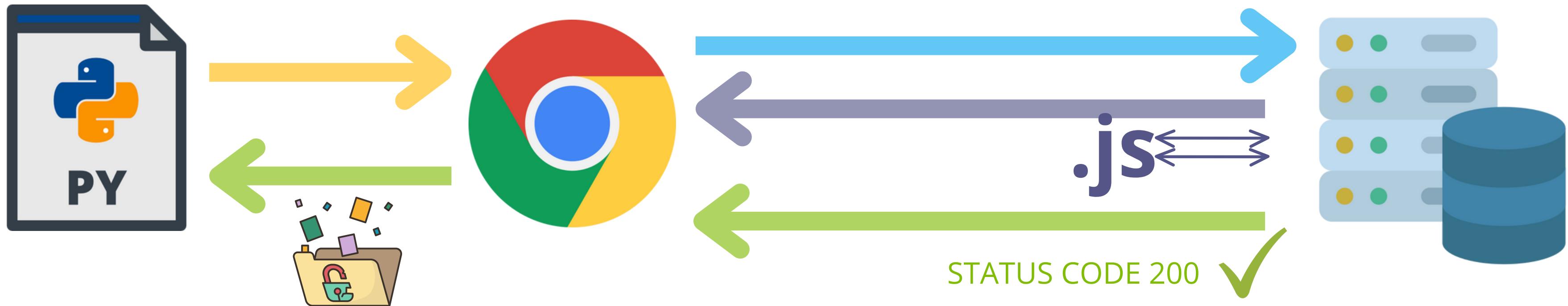


Introduction to selenium

- **What is Selenium?** → A powerful **web browser automation framework** that **simulates real user interactions**.
- Key Features:
 - **Browser Automation:** Controls web browsers programmatically across all major browsers.
 - **Dynamic Content:** Handles JavaScript-rendered content and dynamic page updates.
 - **Cross-Platform:** Works on multiple operating systems and browsers.
 - **Language Support:** Available in multiple programming languages (Python, Java, C#, etc.).

- **Core Concepts:**
 - **WebDriver:** The core component that controls browser actions and interactions.
 - **Elements:** Represent HTML elements on the page that can be interacted with.
 - **Locators:** Methods to find elements (ID, CSS, XPath, etc.).
 - **Waits:** Mechanisms to handle dynamic content loading:
 - **Implicit Waits:** Global timeout for finding elements
 - **Explicit Waits:** Conditional waiting for specific elements/states

Selenium against JS problems



- **Executes JavaScript**, loading dynamic content.
- Captures full content **after JS is rendered**.
- Mimics **browser behavior**, avoiding bot detection.
- **Simulates user interactions** (mouse, scroll, clicks).

What can you make a webdriver do automatically?

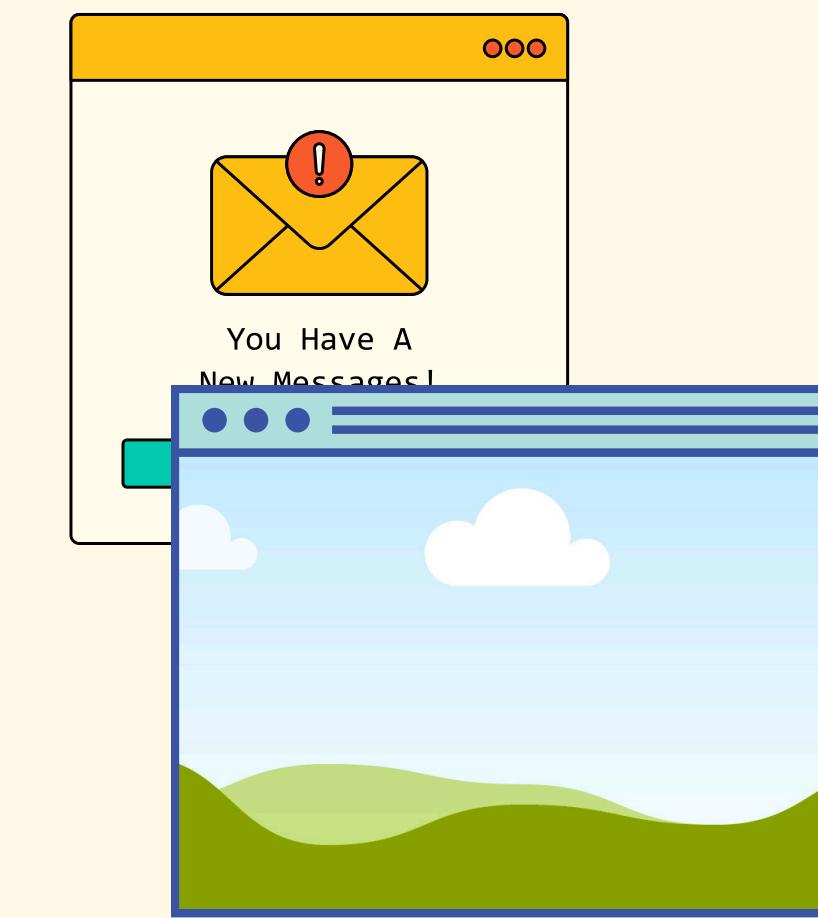
Open URLs, refresh,
back, forward...



Click, type, scroll ...



Handle alerts,
frames, windows...



And
others...

Selenium Introduction



DEMO IN PYTHON





[HTTPS://GITHUB.COM/SARRABENYAHIA/TUTO-WEBSRAPING](https://github.com/sarrabenyahia/tuto-webscraping)

WEBSCRAP THE README OF THE OUR REPO

**END OF
LESSON !**

Sources

Github repository :



GitHub



benyahiasarra9@gmail.com

**Complementary course
material:**

