# How is the web written (and how to read it)

SICSS

CREST

INSTITUT POLYTECHNIQUE DE PARIS

# Outline

**1. The WWW is not the Internet**
- A Brief History of the Internet
- Birth of the WWW
- Internet, HTML and Browsers

**2. The WWW, Back & Front**
- How is the web written?
- The structure of a simple webpage
- The www *millefeuille*

**3. How-to in R**

**4. A note on APIs**

**5. Headless browsers**

2

# The WWW is not the Internet

## 1. A Brief History of the Internet

# The WWW is not the Internet

## 1. A Brief History of the Internet

- *Competing Histories*
  - A "Cold War Technology"
    - DARPA and the race for the future of technology
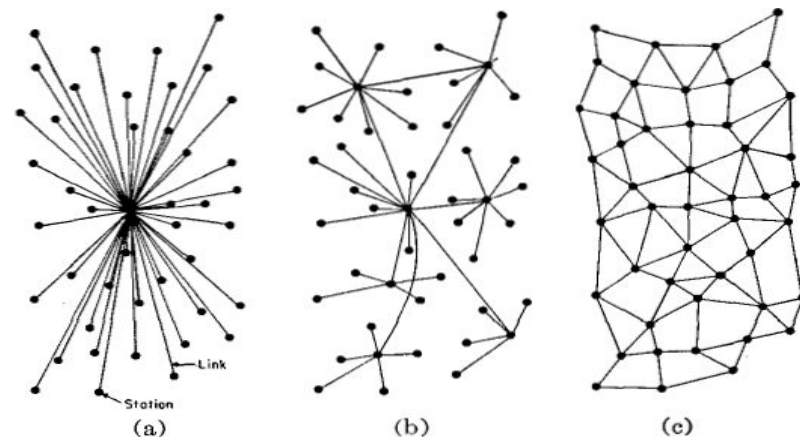    - DARPA & ARPA-NET (1958-1969)



Безопасное путешествие лайка



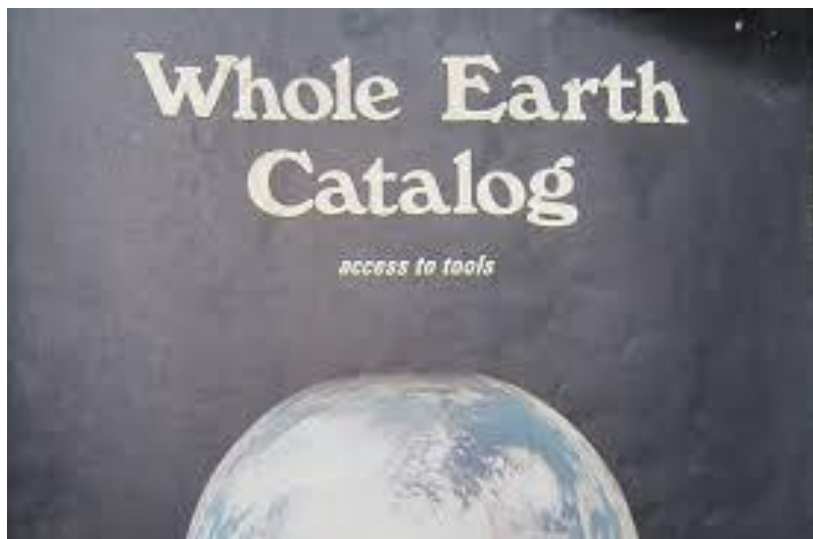Fig. 1—(a) Centralized. (b) Decentralized. (c) Distributed networks.

Paul Baran, 1962

# The WWW is not the Internet

## 1. A Brief History of the Internet

- *Competing Histories*
  - "The Hippies did it": Libertarian origins
    - Augmenting individuals: D. Engelbart against technoscience
    - Collective collaboration and the rise of hackers: science (and connected computers) for the people.

# The WWW is not the Internet

## 1. A Brief History of the Internet

- *Technological & Political Struggles*
  - From circuit switching to packet switching
  - A flurry of networks → TCP/IP (1978-1983)
  - Whose Technology?
  - 1980s: NSF grant to connect US universities
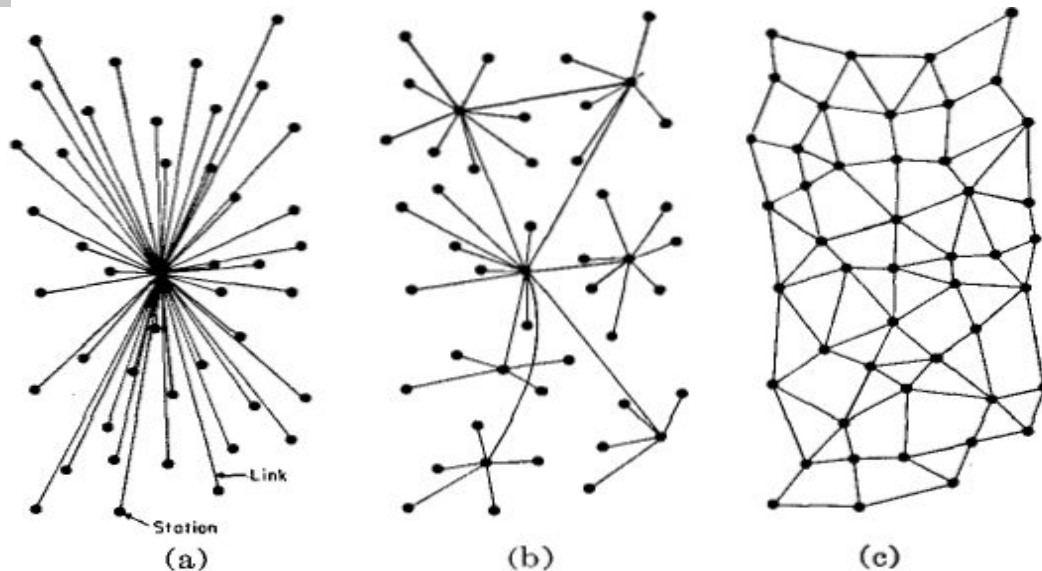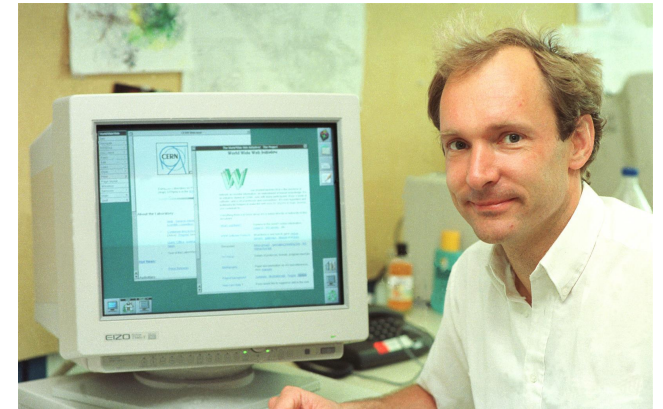
The French Contender (the Minitel)



Fig. 1—(a) Centralized. (b) Decentralized. (c) Distributed networks.

# The WWW is not the Internet

**2. Birth of the WWW**

- *A network that uses the internet*
  - Tim Berners-Lee (1989)

- *A decentralized sorting system*
  - Child of the previous evolution
  - The hyperlink at its core

- *HTML & HTTP*
  - HTML: Hyper Text Markup Language
  - HTTP: Hyper Text Transfer Protocol

# The WWW is not the Internet

## 2. Birth of the WWW

- *The Internet as a Product*
    - IMDB 1990; Amazon, Ebay, Craiglist 1995;
    - Hotmail 1996; Yahoo, Google, Paypal 1998; 2001 Wiki;etc

| Results Summary | | | |
|---|---|---|---|
| Month | # of Web sites | % .com sites | Hosts* per Web server |
| 6/93 | 130 | 1.5 | 13,000 (3,846) |
| 12/93 | 623 | 4.6 | 3,475 (963) |
| 6/94 | 2,738 | 13.5 | 1,095 (255) |
| 12/94 | 10,022 | 18.3 | 451 (99) |
| 6/95 | 23,500 | 31.3 | 270 (46) |
| 1/96 | 100,000 | 50.0 | 94 (17) |
| 6/96 | 230,000 (est) | 68.0 | 41 |
| 1/97 | 650,000 (est) | 62.6 | NA |



Avenue Q: The internet is for...
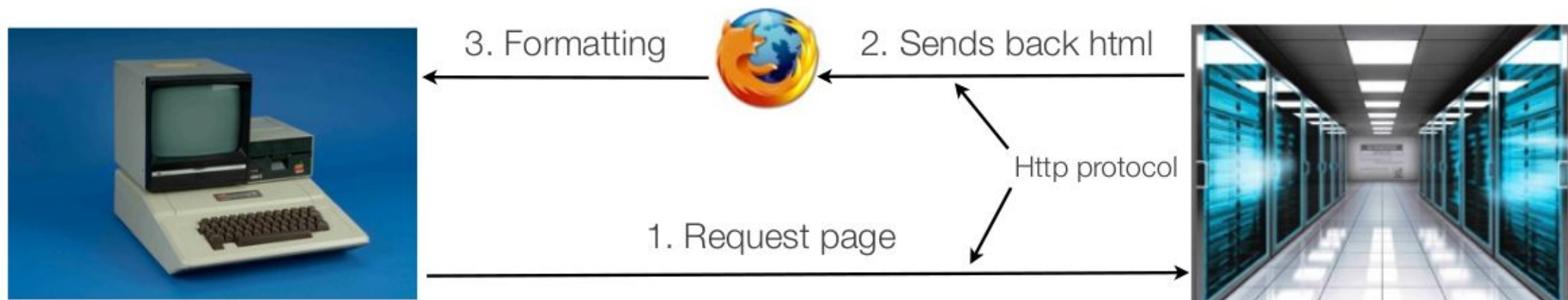
# The WWW is not the Internet

**3. The Internet, the WWW and Browsers**

- *Browsers were essential to popularizing the WWW*

# The WWW is not the Internet

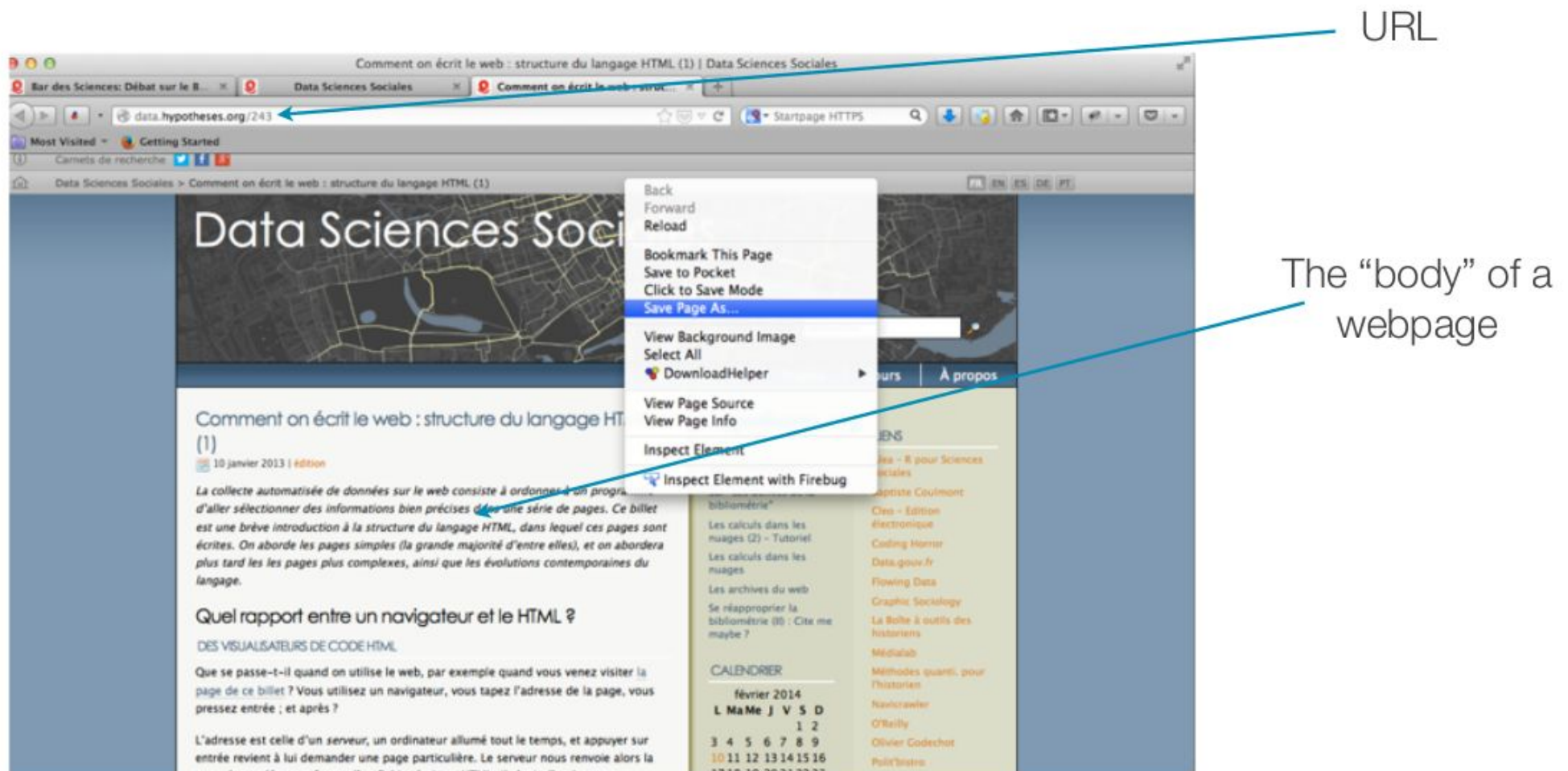## 3. The Internet, the WWW and Browsers

- *Browsers send, receive and interpret code*
- The web relies on the circulation of text in HTML
  - The **www** is based on communication between computers via a protocol called **HTTP**
  - Computers & pages are identified by their address, called a **URL** (Uniform Resource Locator)
  - **HTML** files are transferred and subsequently formatted into a legible format.



3. Formatting    2. Sends back html

Http protocol

1. Request page

# The WWW is not the Internet

## 3. The Internet, the WWW and Browsers

- *Browsers visualize code*

# The WWW is not the Internet

## 3. The Internet, the WWW and Browsers

- *Browsers visualize code*

# The WWW, Back & Front

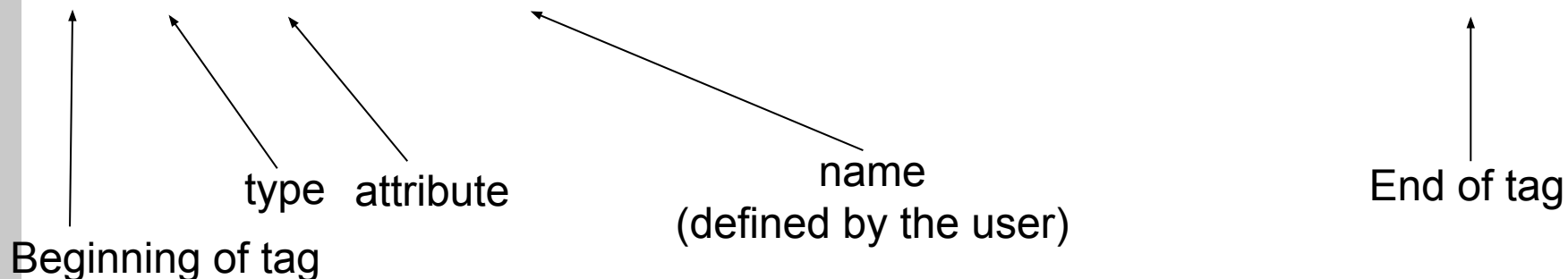# The WWW, Back & Front

**1. How is the Web Written?**
- *The Web is Premised on HTML*
  - A "Rich Language"
  - A Structuring Principle: tags

  - HTML works with tags
    - The text displayed is surrounded by extra information, contained in these containers.

  - Ex. <TAG> This is very interesting </TAG>

14

# The WWW, Back & Front

## 1. How is the Web Written?

- *The Web is Premised on HTML*
  - A "Rich Language" (more than meets the eye)
  - A Structuring Principle: **tags**
  - Tags have a **type** and an **attribute**
    - *Types* are fix (a, span, li, div). They have a limited number of attributes.
    - For more, see <u>this page</u>.

&lt;span id="cheers"&gt; This is very interesting &lt;/span&gt;

Beginning of tag

type    attribute

name
(defined by the user)

End of tag

# The WWW, Back & Front

## 1. How is the Web Written?
- *Head & Body*

```
<!DOCTYPE html ...>  ←——————————————————— Type of document (declares format)
<html>
    <head> ←————————————————————————————— Head (meta informations
        <meta ... >                              about the document)
        <link ...>
        <script ...>
        ...
    </head>
    <body> ←————————————————————————————— Body of text (where
        ...                                   most of the action
    </body>                                  happens for us)
</html>
```

# The WWW, Back & Front

**1. How is the Web Written?**

- *A few common tags*
    - \<div> : block of text
    - \<p> : paragraph
    - \<a> : hypertext link
    - \<h1> : (resp. h2, h3, h4, h5) titles
    - \<!...> : Comment

# The WWW, Back & Front

## 2. The Structure of a Simple Webpage

- *An HTML File has a Tree Structure*

# The WWW, Back & Front

## 2. The Structure of a Simple Webpage

- *What happens in the code is visible on the page*



Display source code

19

# The WWW, Back & Front

## 2. The Structure of a Simple Webpage

- *Tools in your browsers help you "inspect elements"*

# The WWW, Back & Front

**3. The www *millefeuille***

- A webpage is built on HTML

- And it includes other types of files
  - Pictures
  - Content Style Sheet (CSS)
  - Javascript

> Increasingly, the web has become a *millefeuille*

# The WWW, Back & Front

## 3. The www *millefeuille*

> Increasingly, the web has become a *millefeuille*

This has consequences for scraping. But keep in mind that a regularity on the screen means regularity in the code. We are going to use this

# How-to in R

**1. Finding your way around**
There is a wealth of dedicated libraries

This page maintains a list of all that there is at moment (and it's plenty)

Scraping: *httr or **rvest***

Selecting in HTML (or XML): *XML or **rvest***
Selection in json: *rjson*, *rjsonio, jsonlite...*

# How-to in R

**2. Basic instructions**

**read_html()** will read the page and transform it into an XML document.

Thus,

read_html("https://sicss.io/2022/paris/schedule") will output the source code of the schedule page for the SICSS-Paris program.

# How-to in R

**2. Basic instructions**

# How-to in R

**2. Basic instructions**

Yes, in ~70% of the cases, all you need to do to scrape a page is to do

**read_html("*PAGE*")**

# How-to in R

Sometimes that won't be enough, **and the website will see the evil crawler in you**.

You will need to dress-up like an honest browser

This goes through the user_agent command (*httr*)

user_agent("Mozilla/5.0 (Macinstosh;U; Intel MacOS X 10.6; en-US")

user_agent("roger.rabbit@gmail.com")

"On the Internet, nobody knows you're a dog."

# How-to in R

Sometimes that won't be enough, **because you'll need cookies**.

With rvest, you will have to create a **session**, which stores the said cookies and allows you to navigate from there.

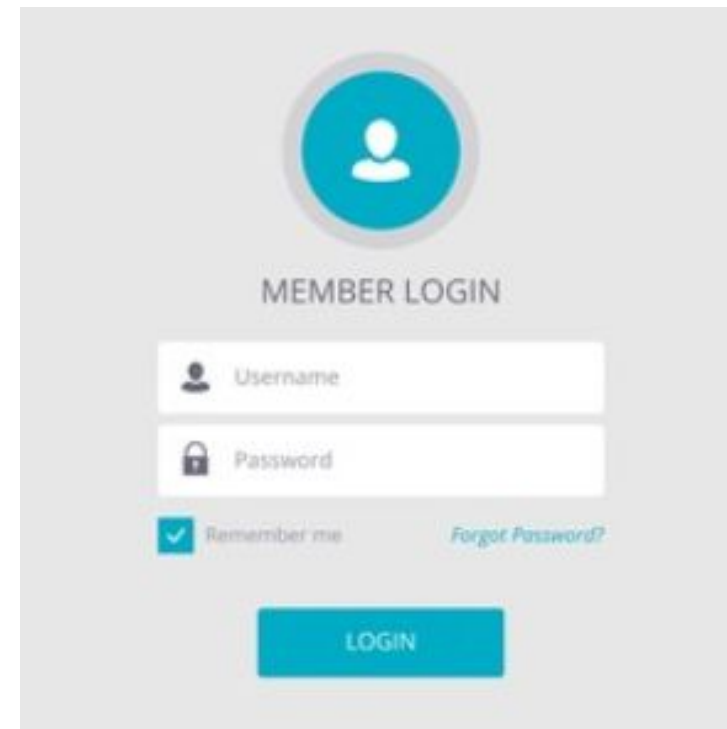And then you'll use **read_html()**

# How-to in R

Sometimes that won't be enough, **because you'll need to log in**.

You'll use a function called **html_form(),**

And then you'll use **read_html()**

# How-to in R

Once you have done that:

**Great news: we are back to square 1!**

(You could have copied and pasted the source code in your console, couldn't you?)

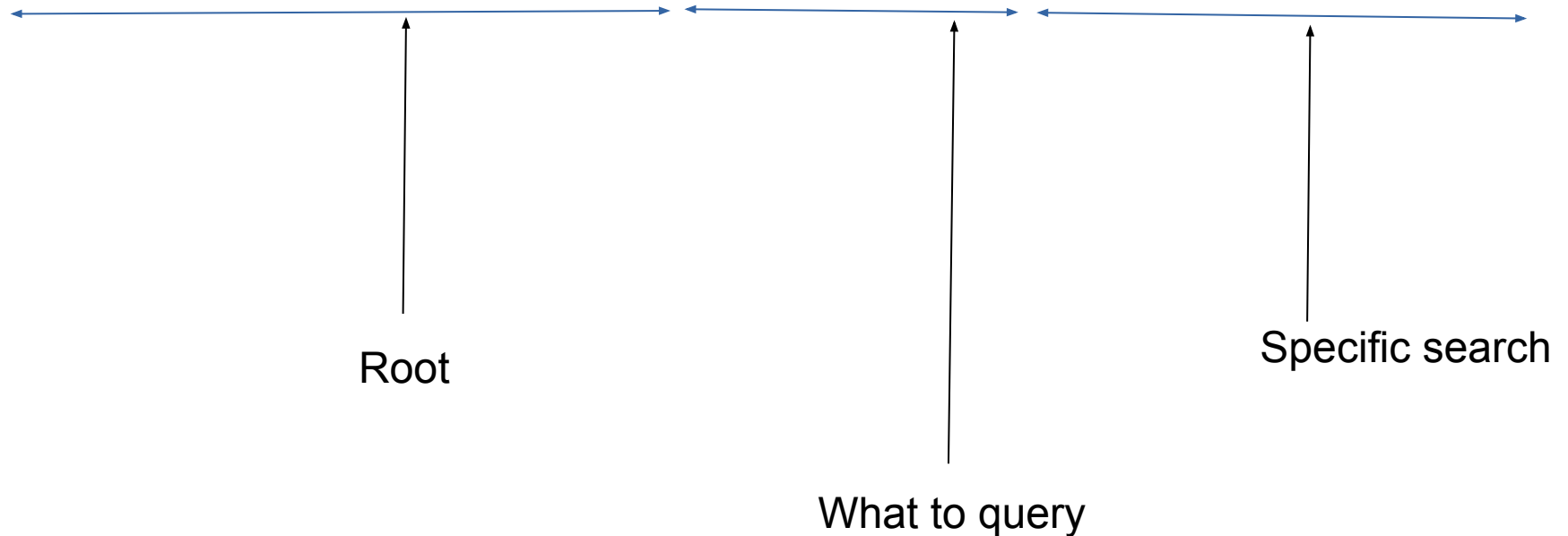Except that now, it is in R, and you can use it

**30**

# A quick note on APIs

# APIs

Application Programmers Interface
- A feature of the web 2.0
- Legal, and often easier
- Different forms: login, with rate limits, or just open

# APIs
## A simple yet elaborate API

- https://www.openstreetmap.org/search?query=ecole%20polytechnique%20paris

Root

Specific search

What to query

# APIs

Most APIs require registration

See the now ~~over~~often-used Twitter API
for academics



1. Register for a project
https://apps.twitter.com
2. Wait for approval
3. Get your credentials
4. Make requests

See Chris Bail's underlined_detailed page on
APIs

34  See academicTwitteR webpage

## APIs

Back to legal and deontological matters

- If there is an API, use it
- If there is no API, ask yourself: are you doing something illegal?

- Sure, scraping Twitter is legal, and its content is public. But what do you learn from individuals? And how should you protect them?

# A Quick Note on Headless Browsers

A growing trend in the web industry is to have websites that respond to your behavior (scrolling, clicking, etc).

# A Quick Note on Headless Browsers

For behavior-based websites

# A Quick Note on Headless Browsers

To do this, you will need to use Javascript in order to create a "**headless browser**", i.e. a browser piloted from your command line.

# A Quick Note on Headless Browsers

To do this, you will need to use Javascript in order to create a "**headless browser**", i.e. a browser piloted from your command line.

This is slightly more complex as you need to install other software, but we'll see an example later.

This is also an easy way to avoid some classic headaches.

# A Quick Note on Headless Browsers

In R, this is often done using "Selenium".
To do so, install "RSelenium"

For me, it worked better installing Docker too (…)

And you will need to type, in the command line, a few lines of code. See this explanation by <u>Chris Bail</u>.

# Conclusion



But all you need to know is **read_html()**
(and where to look for)