# The Hypertext Corpus Initiative

# methods and tools for Social Sciences to build corpus from the web

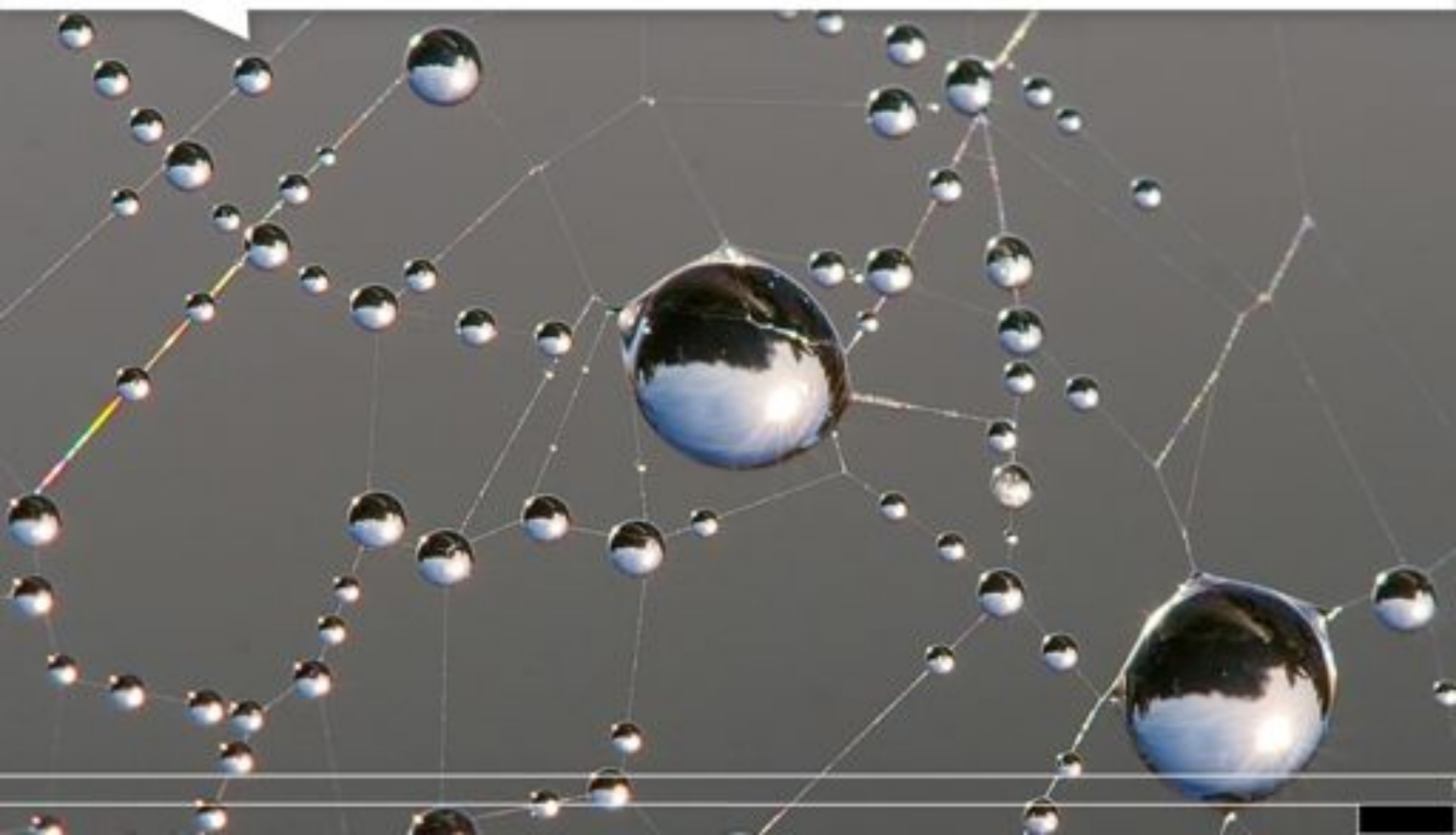Paul Girard, Mathieu Jacomy, Audrey Baneyx, Tommaso Venturini

# What is actually a website ?
Or why we need to define web entities ?

# From sources to corpus ?
a method to select and collect sources to build a corpus

# What for ?
Analysis opportunities and limitations
Interactions with web archive initiative
From tools to equipment

WWW : a network of ressources

# Uniform Ressource Locator
# Hypertext REFerence

*Nick Finck @ flickr*

# A network of references

# Hierarchies of addresses

*rfc3986*

« The URI syntax is organized hierarchically […]

It is **often** the case that a group or "tree" of documents has been constructed to serve a common purpose, wherein the vast majority of URI references in these documents point to resources within the tree rather than outside it.

Similarly, documents located at a particular site are **much more likely** to refer to other resources at that site than to resources at remote sites. »

# clusters of HTML pages

# hierarchical namespace

*sciences-po.fr domain*

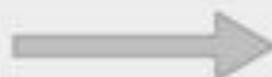# pointers in the URL hierarchy

RFC3986 : about URI

" The URI syntax is organized hierarchically, with components listed in order of **decreasing significance from left to right**."
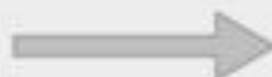
RFC882 : about domain names

"by convention, the labels that compose a domain name are read left to right, **from the most specific (lowest) to the least specific (highest)**."

# pointers in the URL hierarchy

# handle heterogeneity : aliases

slideshare.net/medialabSciencesPo

vimeo.com/medialab/videos

medialab.sciences-po.fr

medialab

ACTOR'S PRESENCE
ON THE WEB

sciences-po.fr

sciencespo.fr

sciences Po

ALIASES

# maintain complexity

médialab.
sciences-po

Sciences Po

recherche
sciences-po

# From sources to corpus

*sourcing*

## Define



WEB
ENTITY

## Select

**status**

✓ included
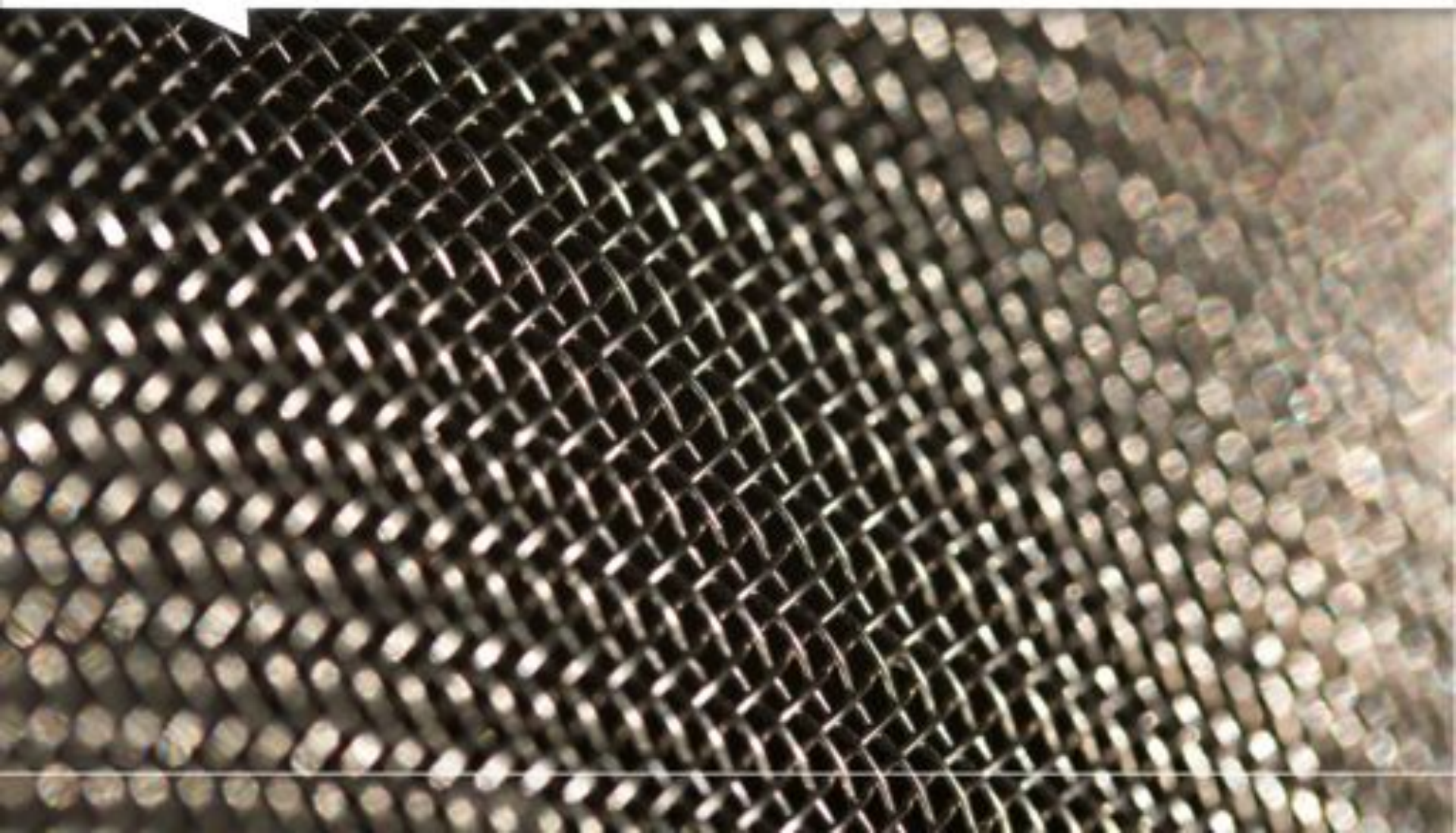
✗ excluded

? undefined

∅ unknown

## Describe

**tagging**

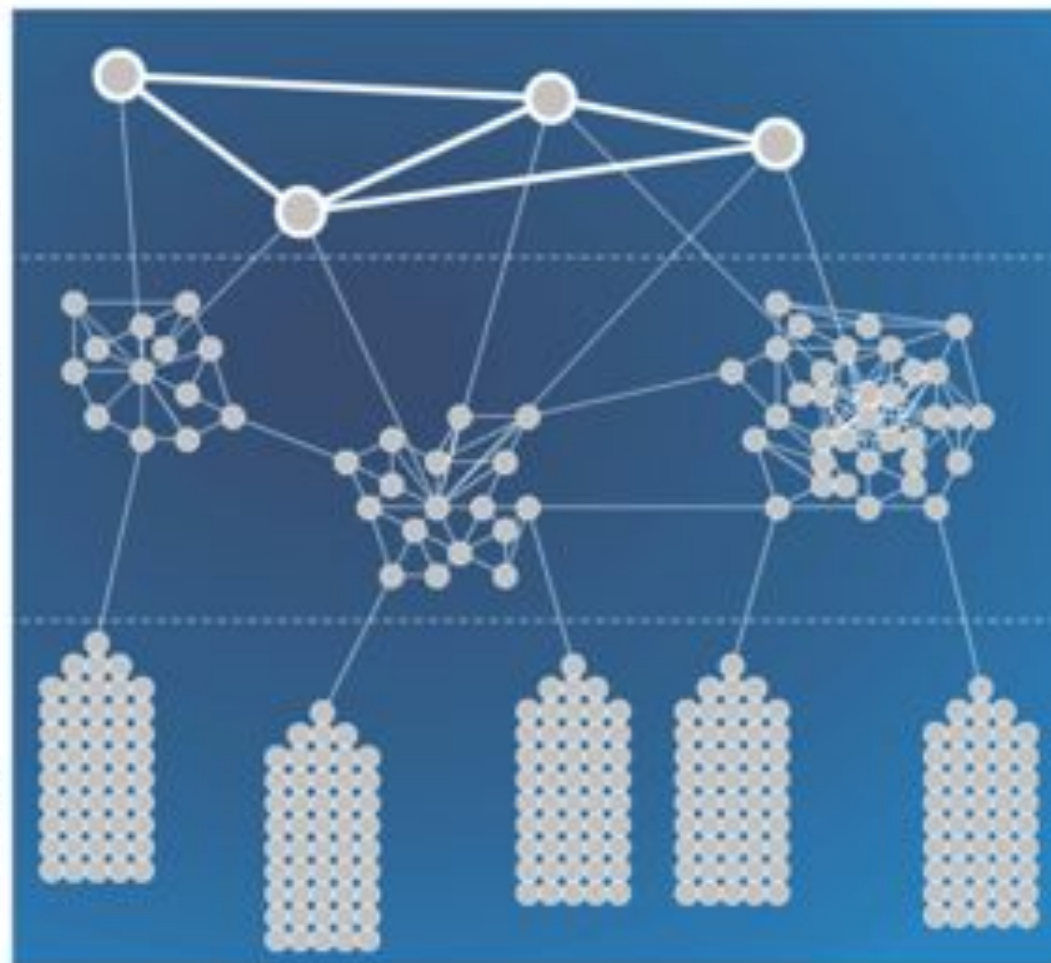Sourcing : how to sieve the web ?

© Mayhem Chaos

# WWWeb is an organized space



Couche **la plus visible** du web : Google, Amazon, Voilà, SNCF, etc...

Couche intermédiaire : **agrégats**, communautés en ligne

Couche profonde : **bases de données**

# Select : from sources to corpus

Sourcing : a qualitative task which gives entry points
- field enquiry, interviewing the actors
- use search engines and browse
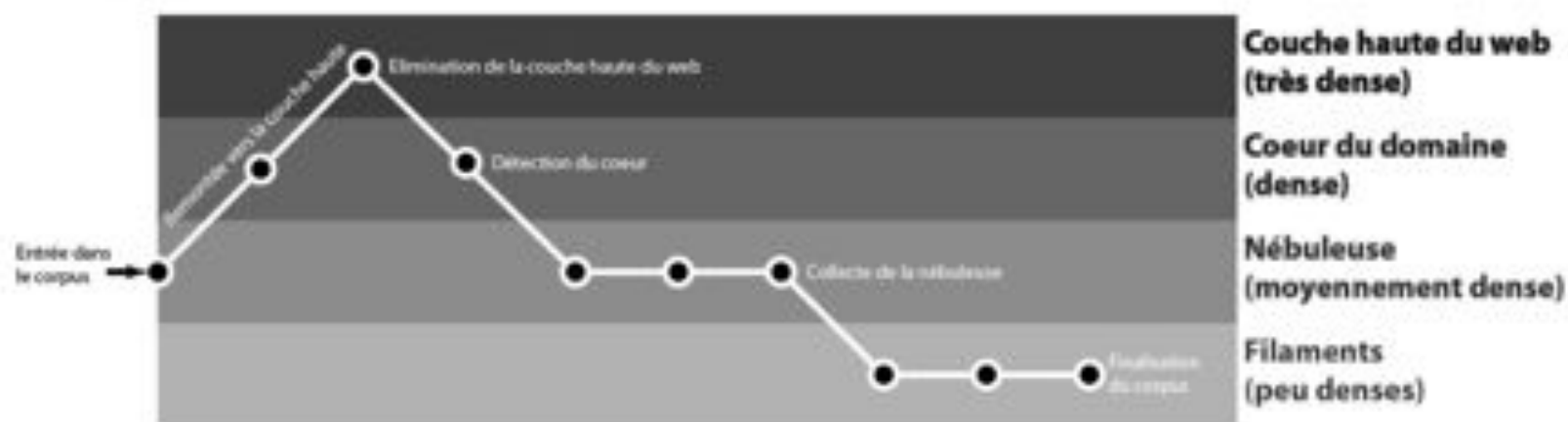
But how to construct and validate a web corpus ?
By following the medium...
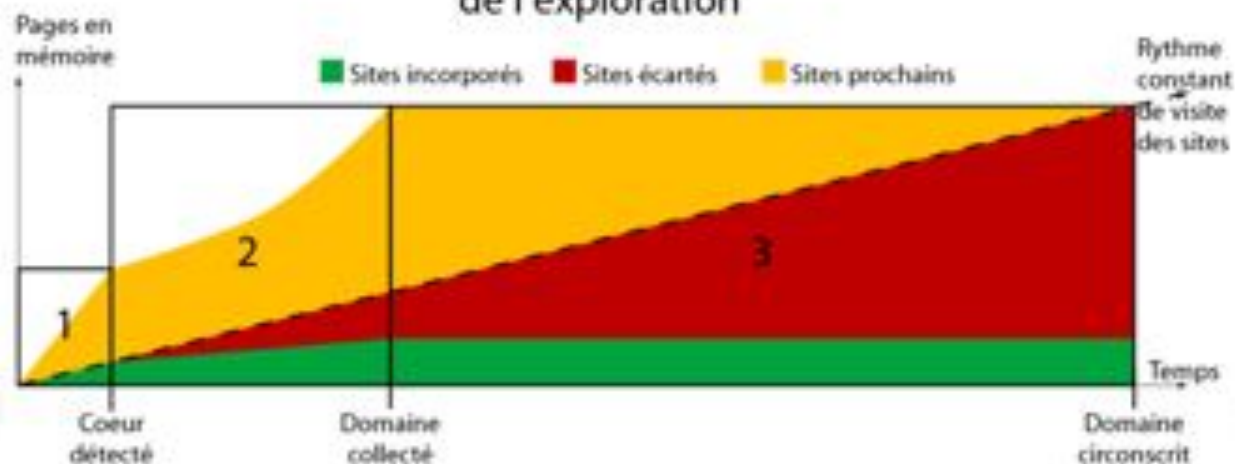Extending the corpus following links, peer to peer.
With prospective crawls.
Controlled by the researcher.

# Explore and limit the corpus



Courbe d'avancement
de l'exploration

# Which crawler ?



WebAtlas Navicrawler

navicrawler,
to build a corpus manually

**issue**crawler

issuecrawler,
to build a corpus with automatic crawling

# a difficult choice



tweezers

**?** or

caterpillars

# Research driven crawling

*topic focused corpus?*



SELECTION
OF WEB
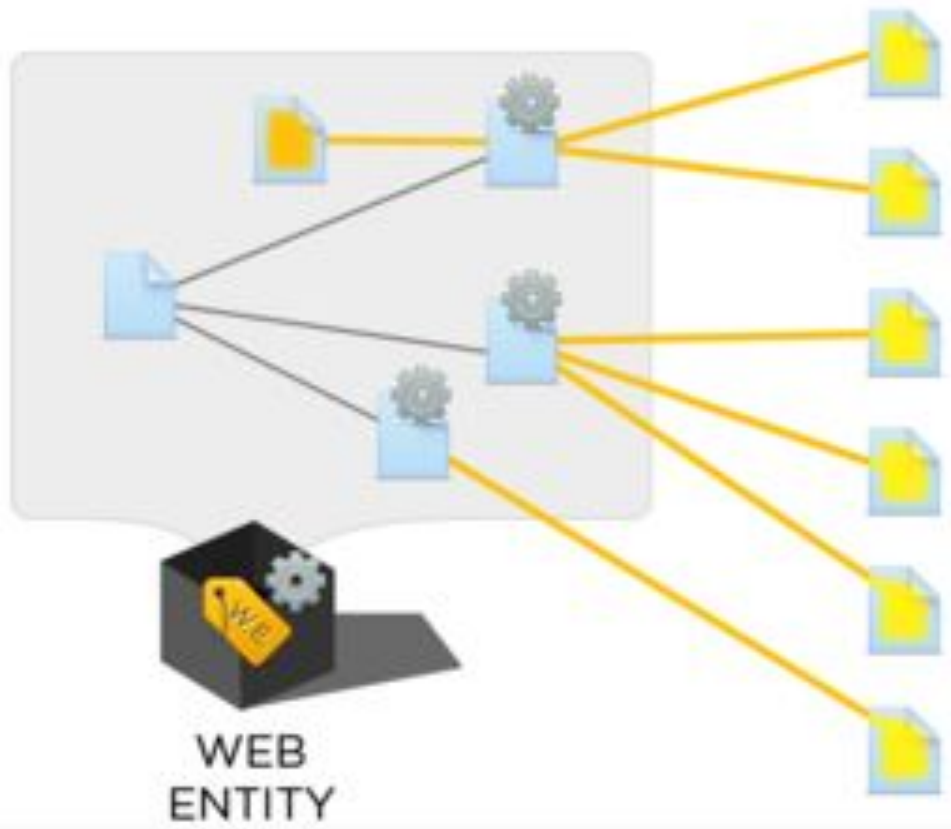ENTITIES

CRAWLING

URL
CANDIDATES

# From sources to corpus

*sourcing*

## Define



WEB
ENTITY

## Select

**status**

✔ included

✘ excluded

? undefined

∅ unknown

## Describe

**tagging**

# Links : topological analysis

*E-Diaspora Atlas*

**French Expatriates**

Politics

Welcoming and Integration

Expatriates' stories

Services for Expatriates

Institutions

# Qualification rules !

**French Expatriates**

Politics

Welcoming and
Integration

Expatriates'
stories

Services for
Expatriates

Institutions

# Web Archives

HCI → web archive : archive a corpus

- export a corpus to be archived by a web archive institution
- harvest rich content
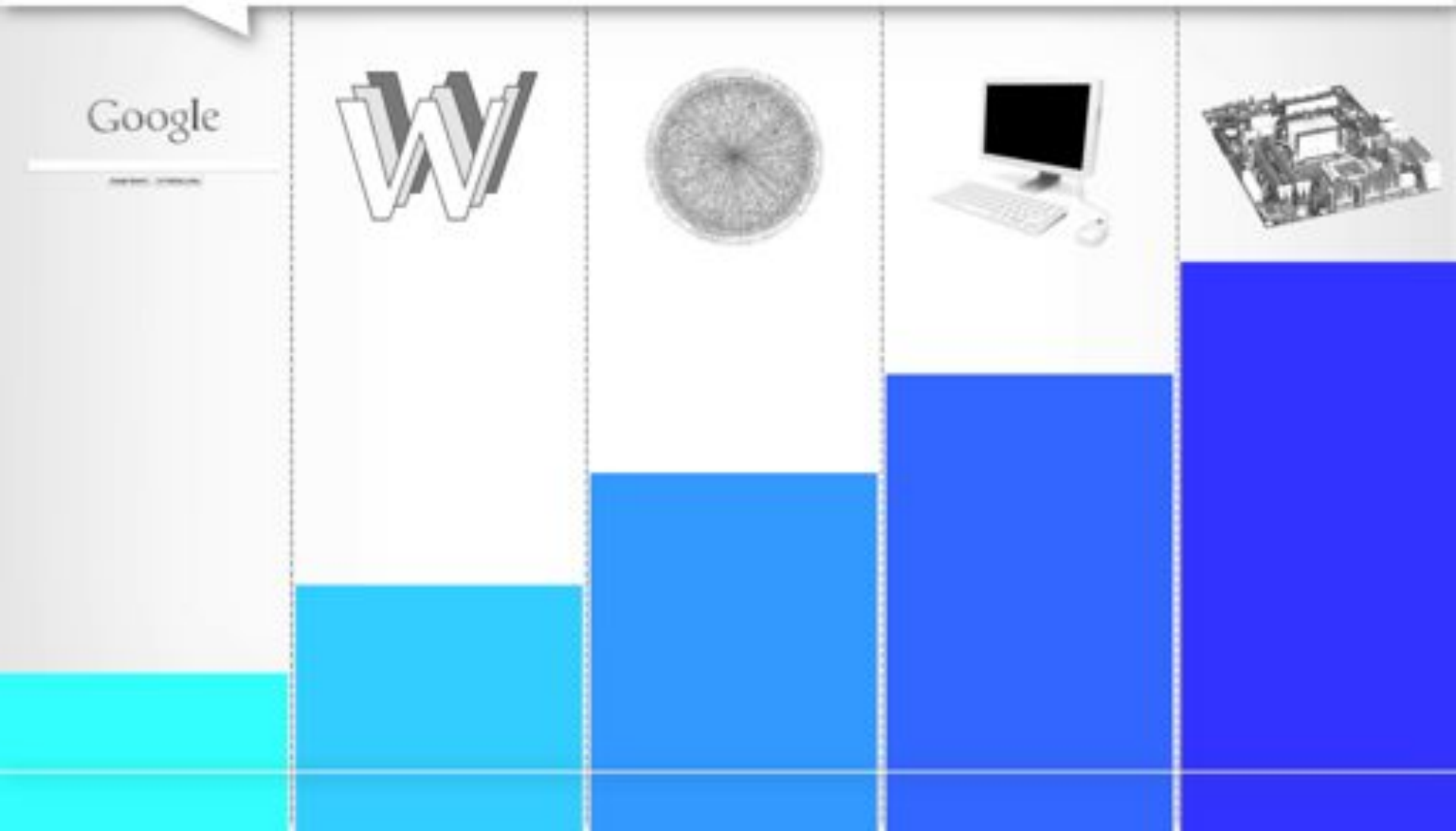- regular harvesting through time

Archive → HCI : crawl an archive

- use HCI to build a corpus from an archive and not the live web
- benefit from the anteriority of the archive

# Representativity ?

Google

# Web field as a carbon paper

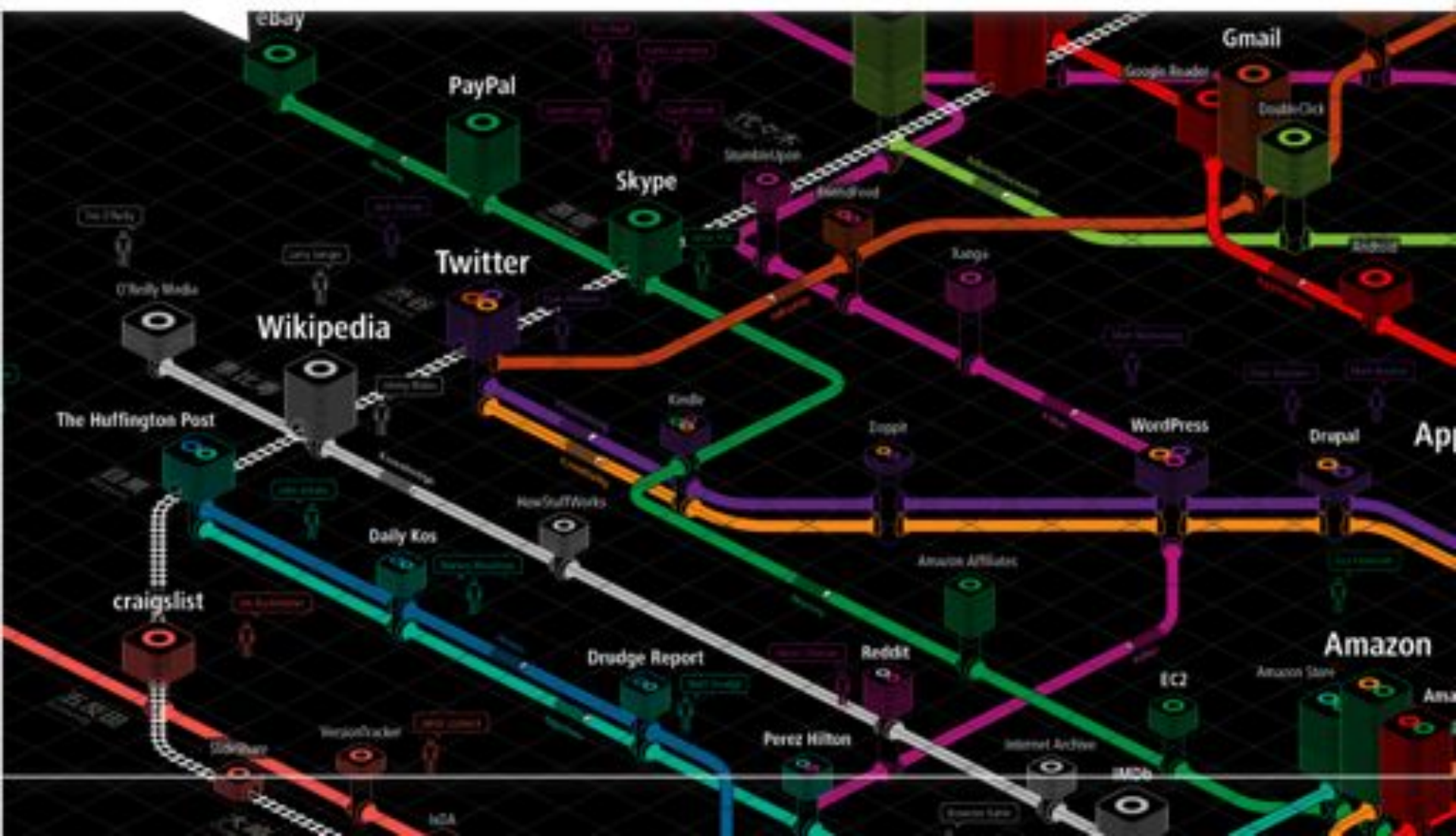*The world seen from Facebook*



facebook

The web is not everything
But don't forget it actually exists

Quanti : ELLIPS panel of 6000 french people equiped with tablets

Quali : BeQuali, a qualitative survey archive

Web : use the web as a survey field

- trainings and assistance : web corpus for Social Sciences
- tools and methods : Hypertext Corpus Initiative
- a technical architecture : storage and crawling servers

thank you

medialab.sciences-po.fr