# Two decades of web archiving: what's on the horizon?

Julien Masanès

Internet Memory (Research / Foundation)

j@internetmemory.net

@julienmasanes

LE JOURNAL    LES ARCHIVES    LA CARTE    VOS FAVORIS

INTERNATIONAL    POLITIQUE    SOCIÉTÉ    ÉCO    CULTURE    IDÉES    PLANÈTE    SPORT    SCIENCES    PIXELS    CAMPUS    LE MAG

# PIXELS
CHRONIQUES DES (R)ÉVOLUTIONS NUMÉRIQUES

VIE EN LIGNE    JEUX VIDÉO    BANC D'ESSAI    CULTURES WEB

## Vingt ans d'archivage du Web : les coulisses d'un projet titanesque

Depuis deux décennies, la fondation Internet Archive, avec d'autres institutions comme la BNF, consigne soigneusement la mémoire du Web pour les générations

---

# ARQUIVO.PT

Entrada    Notícias    Exemplos    Publicações ⌄    Colabore ⌄    Acerca ⌄    Aju

Entrada > Notícias > Investiga XXI > Evento Arquivo.pt no dia 8 de novembro: Inscreva-

## Evento Arquivo.pt no dia 8 de novembro: Inscreva-se!

Outubro 13, 2017

No dia **8 de novembro** vamos realizar o 1º Evento do Arquivo.pt no Pavilhão do Conhecimento em Lisboa e gostaríamos de convidá-lo a celebrar connosco os 10 anos do início do projecto!

## Palestras: porquê e como preservar a Web?

- José Pacheco Pereira, este historiador foi uma das primeiras individualidades portuguesas a destacar o problema da efemeridade da web
- Julien Masanès, é o "pai" do *web archiving* na Europa, editou o 1º livro no mundo acerca deste tema e actualmente lidera o Internet Memory Research
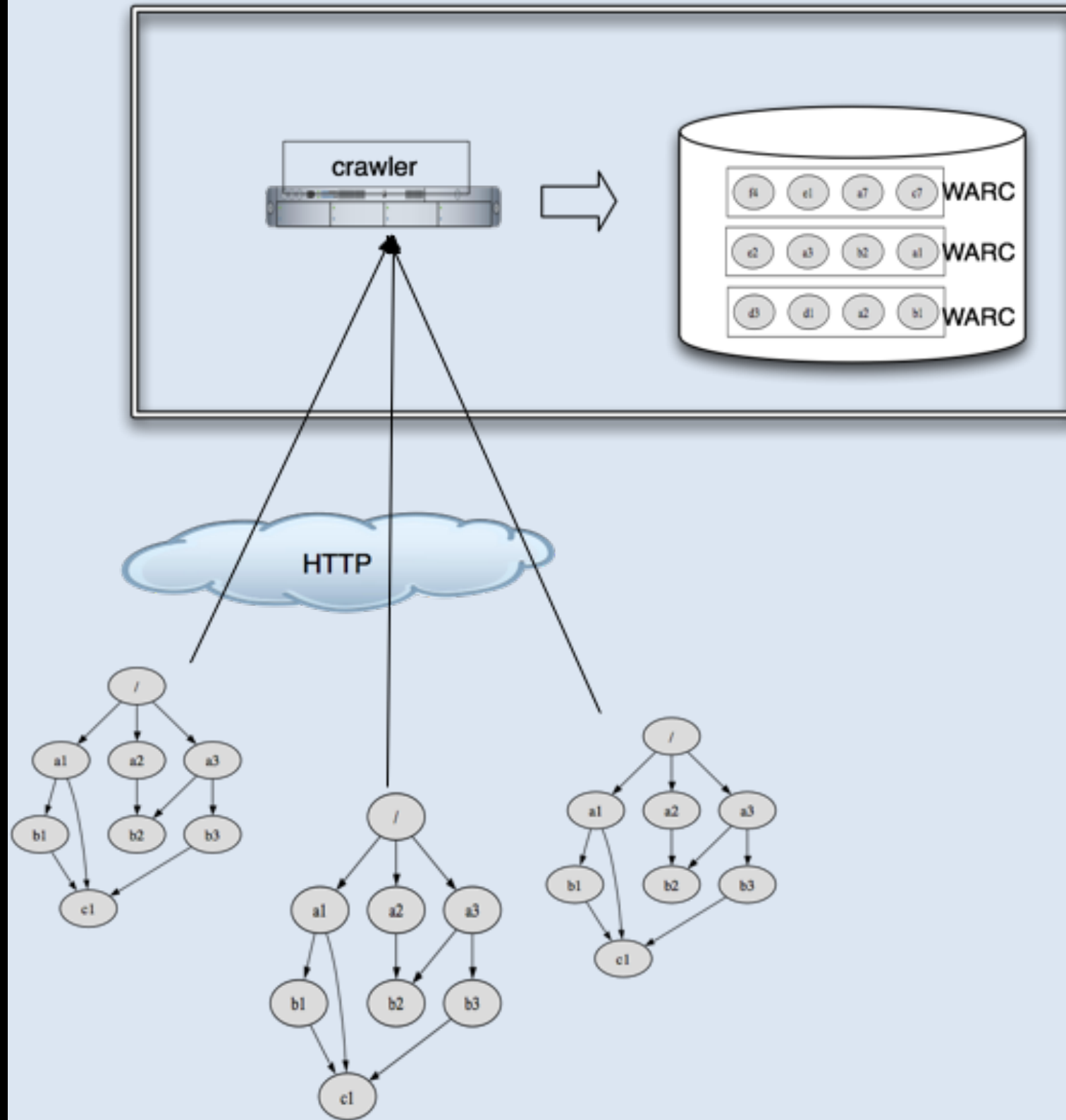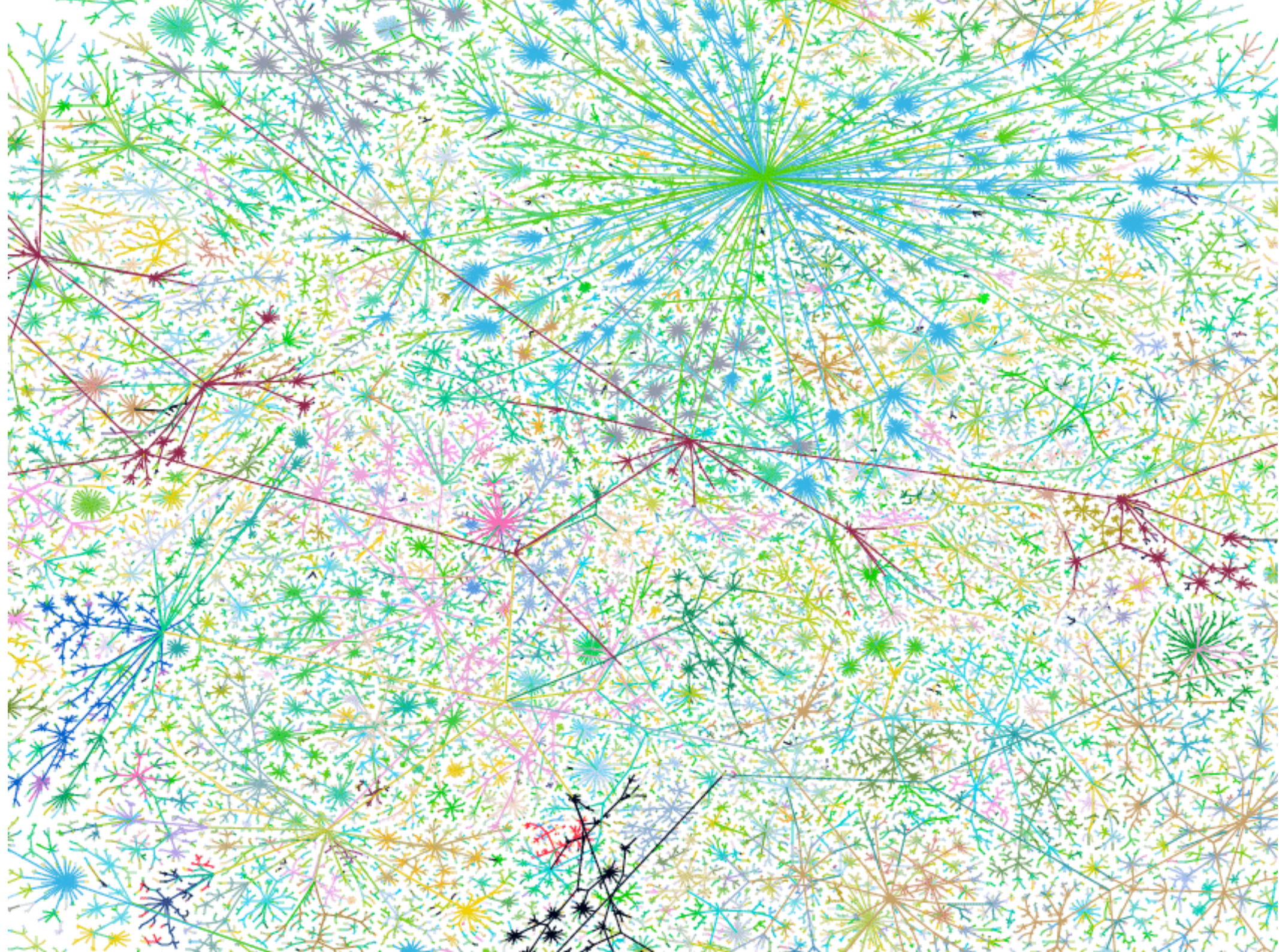
INTERNET ARCHIVE

# Two decades of web archiving: what's on the horizon?

Web archiving fundamentals

How platforms change the landscape

Web archives as analytical corpus

# Two decades of web archiving: what's on the horizon?

**Web archiving fundamentals**

How platforms change the landscape

Web archives as analytical corpus

crawler

WARC

WARC

WARC

HTTP

# Shaping the archive

1. Explicit limits
robots.txt

2. Implicit limits (no DoS)
Politeness

3. Technical limits (arm race)
Client-side code execution

# Cardinality of preserved object

Different for different institutions (museums, archives, libraries)

Cardinality of the incunabula

      20 M  (books) / 30 000 (editions) = 650

Consequence of high cardinality for preservation: redundancy and time

# What about the Web?

Virtually infinite number of copies

High dependancy on one single server

More like a museum than like libraries

# Active publishing system

- Web Information Systems

- Controlled by the producer

- Continuous publishing

Preserving: actively pursuing the exact opposite

# Internet Memory?

- Automated sampling of a virtually infinite information space

- Search and analysis specially leveraging the temporal dimension of the archive

- Has to be part of the internet

# Two decades of web archiving: what's on the horizon?

Web archiving fundamentals

How platforms change the landscape

Web archives as analytical corpus

**Search**    All Formats ▾    Search Loc.gov    **GO**

LIBRARY OF
CONGRESS BLOG

Search this blog    **GO**

**About This Blog**

**Categories**

Abraham Lincoln
American Folklife Center
Audiovisual
Blogging
Books
Capitol Hill
Cataloging
Cervantes
civil war
Collections
Concerts
Congress
Congress Blogs
Copyright

🖨 Print    📶 Subscribe    ♻ Share/Save

# Update on the Twitter Archive at the Library of Congress
January 4, 2013 by Erin Allen

*(The following is a guest post from the Library's Director of Communications, Gayle Osterberg.)*

An element of our mission at the Library of Congress is to collect the story of America and to acquire collections that will have research value. So when the Library had the opportunity to acquire an archive from the popular social media service Twitter, we decided this was a collection that should be here.

In April 2010, the Library and Twitter signed an agreement providing the Library the public tweets from the company's inception through the date of the agreement, an archive of tweets from 2006 through April 2010. Additionally, the Library and Twitter agreed that Twitter would provide all public tweets on an ongoing basis under the same terms.

The Library's first objectives were to acquire and preserve the 2006-10 archive; to establish a secure, sustainable process for receiving and preserving a daily, ongoing stream of tweets through the present day; and to create a structure for organizing the entire archive by date.

This month, all those objectives will be completed. We now have an archive of approximately 170 billion tweets and growing. The volume of tweets the Library receives each day has grown from 140 million beginning in February 2011 to

# Preserving the platforms?

- Who decides?
    - A handful of private company can now decide what will be preserved?

- Legal limits
    - They have technical capacity enable archiving, Updating regulation?

- Technical limits
    - We need more machines!
    - API access, how we deal with loss of presentation context?

http://webarchive.parliament.uk/*/http://twitter.com/
ukparliament

# Two decades of web archiving: what's on the horizon?

Web archiving fundamentals

How platforms change the landscape

**Web archives as analytical corpus**

# Macroanalysis

- Focus is no longer single document/page

- Extraction of data

- Indexing

- Analytics: statistical distribution , correlation, graph analyses

**Generic**

Sourcing

Crawl

Indexing

Processing

Mining

**Value Chain**

Analytics

Research

Visualisation

**Specific**

Research issues

22

# Working papers



Breton corpus    Chinese corpus    Egyptian corpus    French Expatriates    French Repatriates    Hindu – Hindutva

Hmong corpus    Indian corpus    Indian Real Estate    Italian corpus    Jewish corpus    Kerala corpus

Lebanese corpus    Macedonian corpus    Mexican corpus    Moroccans on FB    Moroccan corpus    Nepali corpus

*http://www.e-diasporas.fr/*

*Occupying the social Web: Moroccan students on the move. A method resulting from fieldwork*
*Sabrina Marchandise, 2012.*

# New possibilities to large scale analysis of web archives: Deep learning



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

# Deep learning to detect objects
## « luggages »

# Also to help building the archive!



ARQUIVO.PT

*Novas formas de viajar ao passado*

Type word or URL ✕

Search pages from the p

Millions of contents archived since 1996

Meet the service

# Better crawl

# Genre analysis

- Create model, language independent for publication genres:
  - News/e-commerce/Blogs/Forum/Porn/Company/
  - Detect webspam

- Evaluate how this can be use for large scale crawling priorization
  - Change revisit frequency depending on genre
  - Change crawling budget depending on genre
  - Blacklist spam

IMR INTERNET MEMORY RESEARCH

# Learning to detect genres

Site level

Shop

Page level

News

Shop

Forum

Shop

Vector
0.5, 0.3, 0.1, 0.2

Vector
0.2, 0.3, 0.1, 0.6

Vector
0.4 0.3, 0.7, 0.2

Vect
0.7, 0.3, (

# e-commerce sites

# Blog sites

# Obligado !

• 
: 
: 

**Julien Masanès**

**Internet Memory (Research / Foundation)**
**j@internetmemory.net**
**@julienmasanes**