Old Dominion University ODU Digital Commons

Computer Science Presentations

Computer Science

6-26-2014

Bits of Research

Michele C. Weigle Old Dominion University, mweigle@odu.edu

Follow this and additional works at: https://digitalcommons.odu.edu/ computerscience_presentations



Part of the <u>Archival Science Commons</u>, and the <u>Computer Sciences Commons</u>

Recommended Citation

Weigle, Michele C., "Bits of Research" (2014). Computer Science Presentations. 33. $https://digital commons.odu.edu/computerscience_presentations/33$

This Book is brought to you for free and open access by the Computer Science at ODU Digital Commons. It has been accepted for inclusion in Computer Science Presentations by an authorized administrator of ODU Digital Commons. For more information, please contact digitalcommons@odu.edu.



Bits of Research

Dr. Michele C. Weigle

Graduate Program Director
Web Sciences and Digital Libraries (WS-DL) Lab
Department of Computer Science
Old Dominion University

June 26, 2014

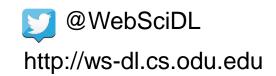


Outline

- Bits of Research
 - digital preservation and web archiving
 - information visualization

- Potential Summer Projects
 - Blue Button visualization (health informatics)
 - visualizing aggregate health data
 - exploring large document collections

ODU's WS-DL Group



Web Sciences and Digital Libraries

- digital preservation
- web archiving
- web science (social media analysis, web usage analysis)

Our recent work has been featured in the popular press











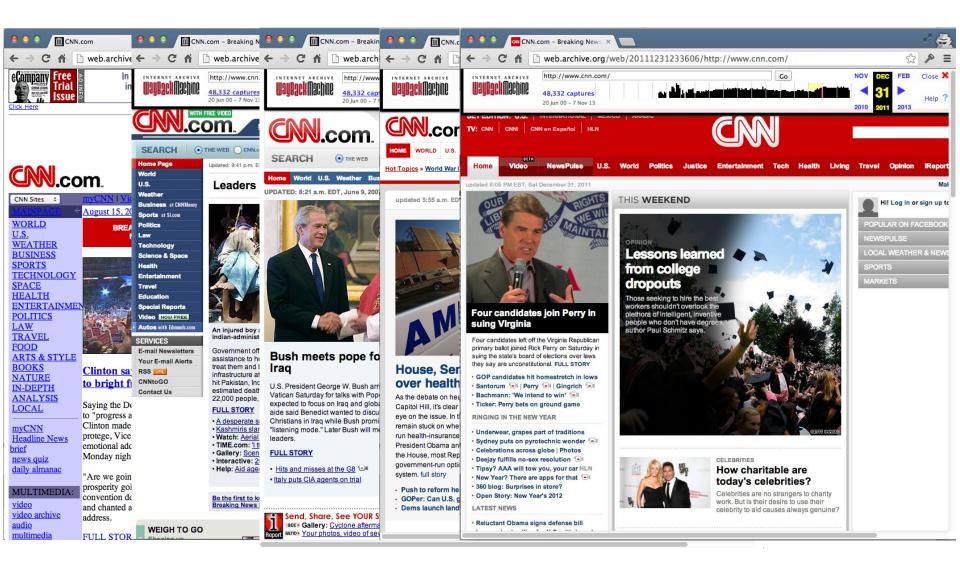








What is a web archive?





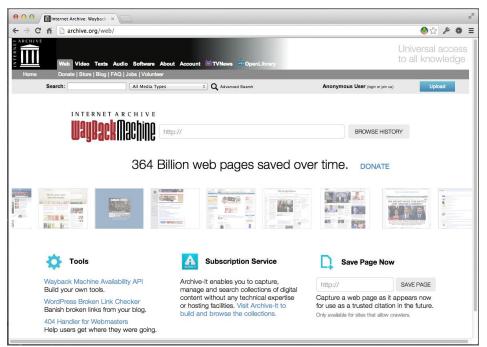
June 26, 2014

What are some web archives?





How can I access the archives?





Memento for Chrome

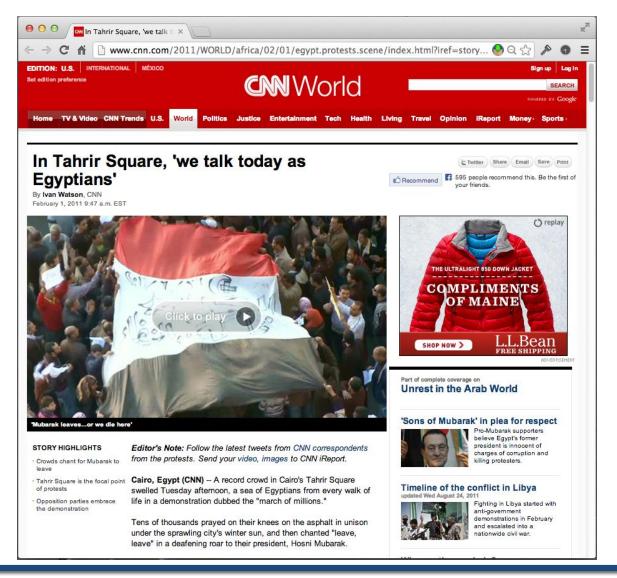


http://www.mementoweb.org/

http://ws-dl.blogspot.com/2010/03/2010-03-19-mementofox-add-on-released.html http://ws-dl.blogspot.com/2013/10/2013-10-14-right-click-to-past-memento.html



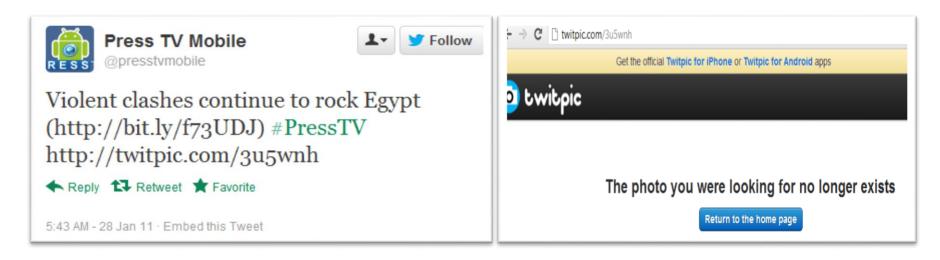
The Web holds our stories





June 26, 2014

But webpages can disappear



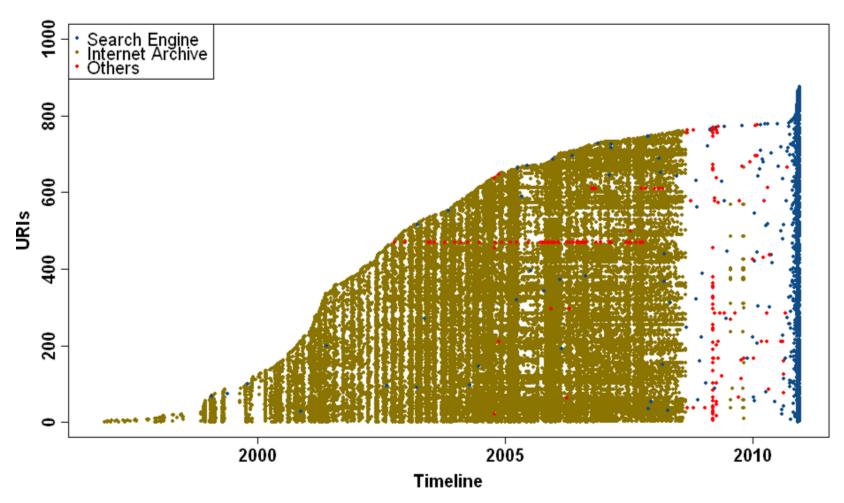
Average lifespan of a webpage - 50-100 days

• A year after publication, about 11% of content shared on social media will be gone.

SalahEldeen and Nelson, "Losing My Revolution: How Many Resources Shared on Social Media Have Been Lost?", TPDL 2012 http://ws-dl.blogspot.com/2012/02/2012-02-11-losing-my-revolution-year.html



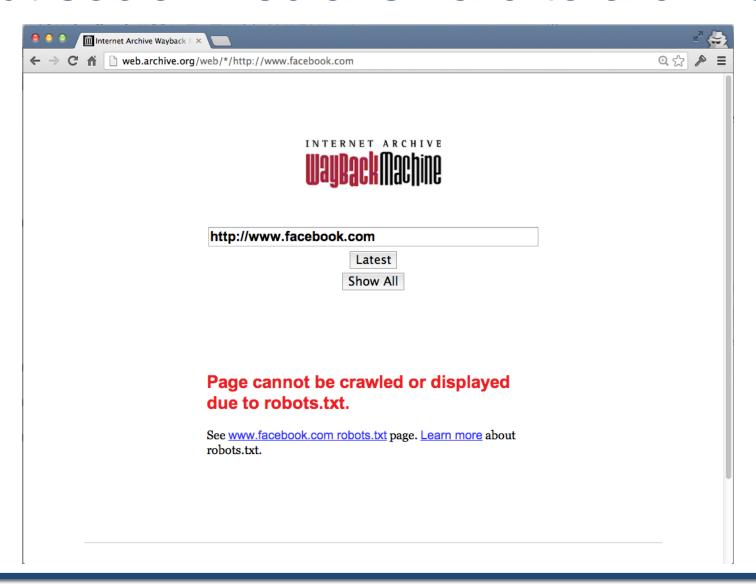
But maybe it's archived



Ainsworth, AlSum, SalahEldeen, Weigle, and Nelson, "How Much of the Web is Archived?", JCDL 2011 http://ws-dl.blogspot.com/2011/06/2011-06-23-how-much-of-web-is-archived.html



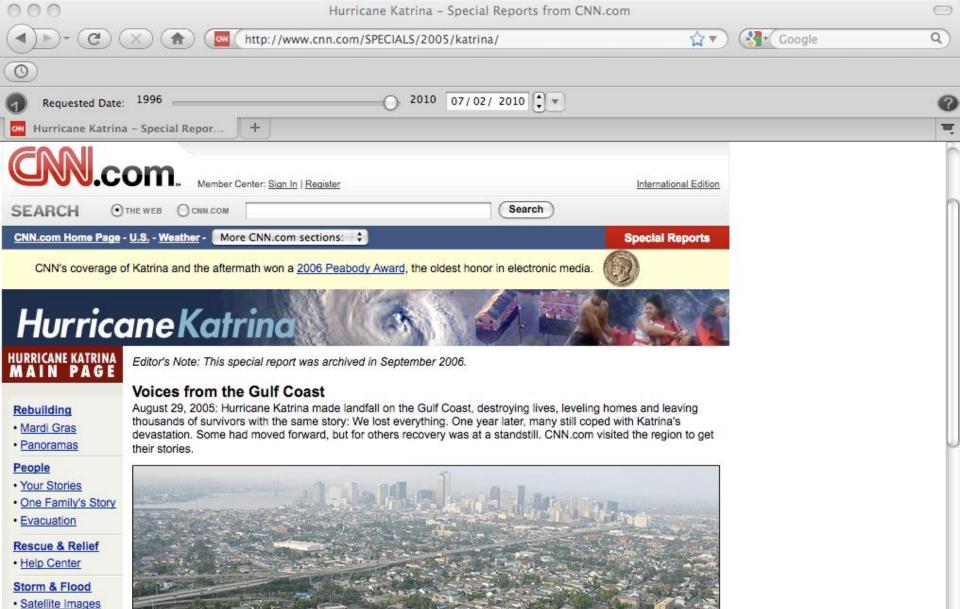
But social media is hard to archive





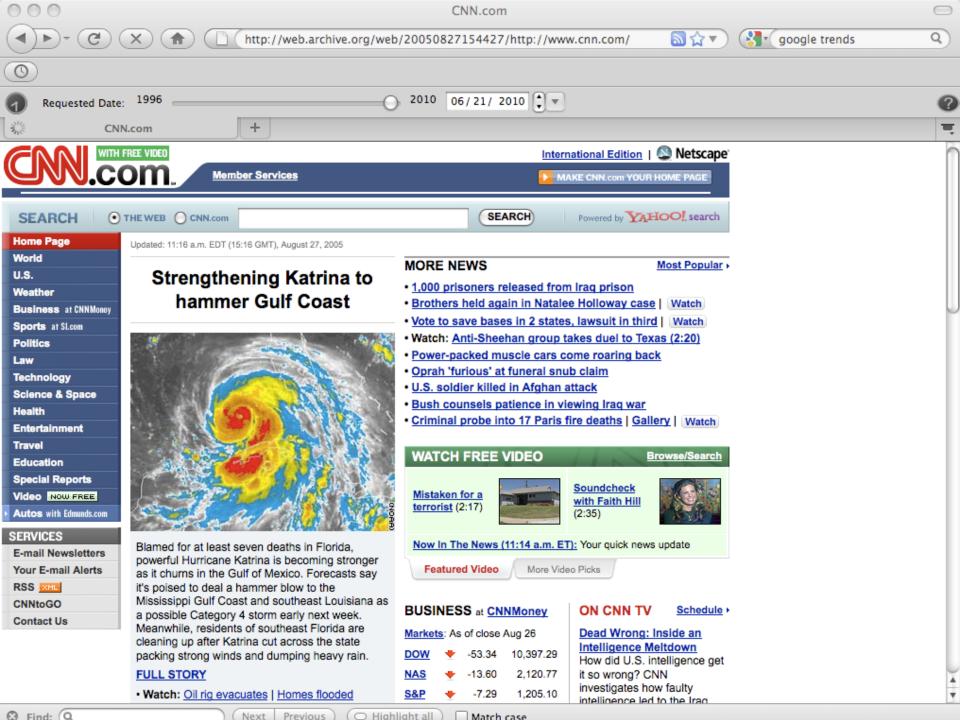
Our Research Group Goals

- We believe that web archives are valuable cultural resources, and we want everyone to know about them.
- We want to make it easy for people to bridge the gap between the live web and the archives.
- We believe that replaying the past is more compelling than reading a summary.



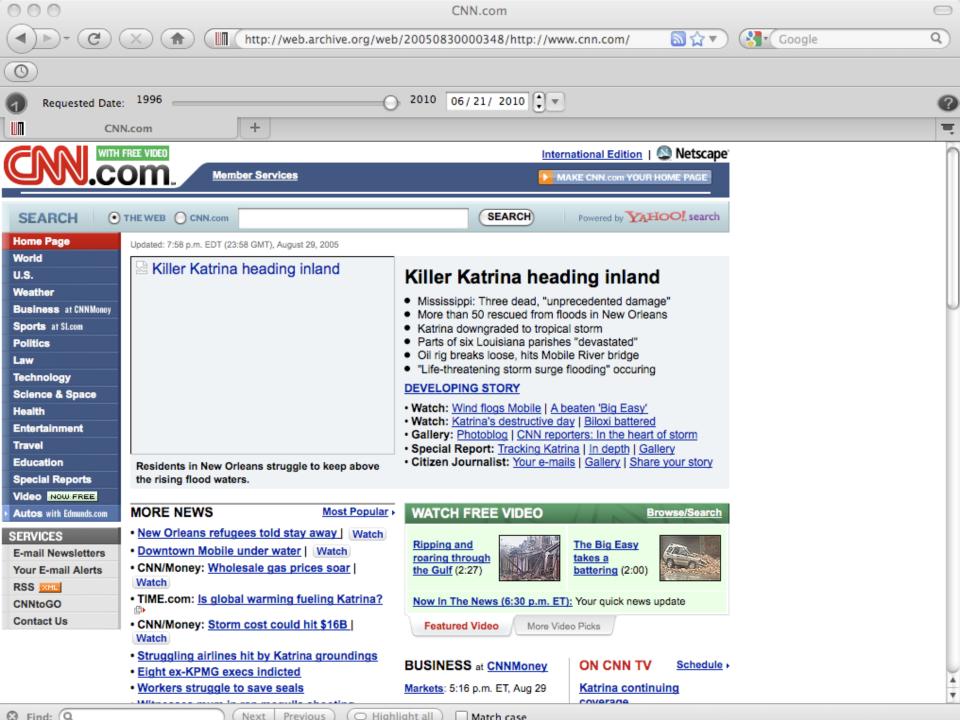
Hurricane Season Special Report VS.











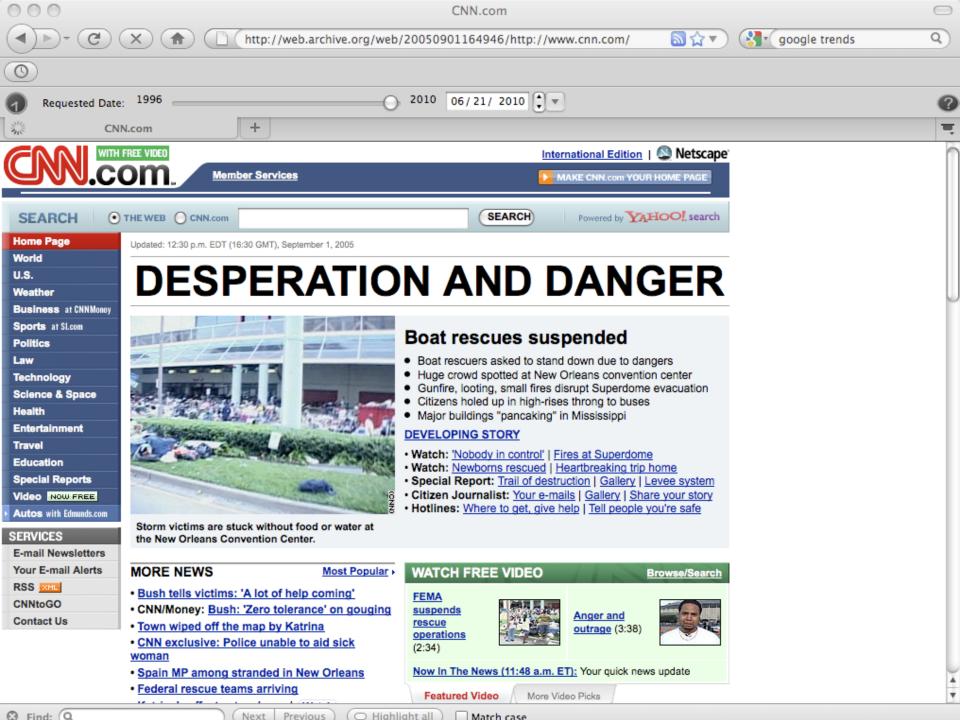


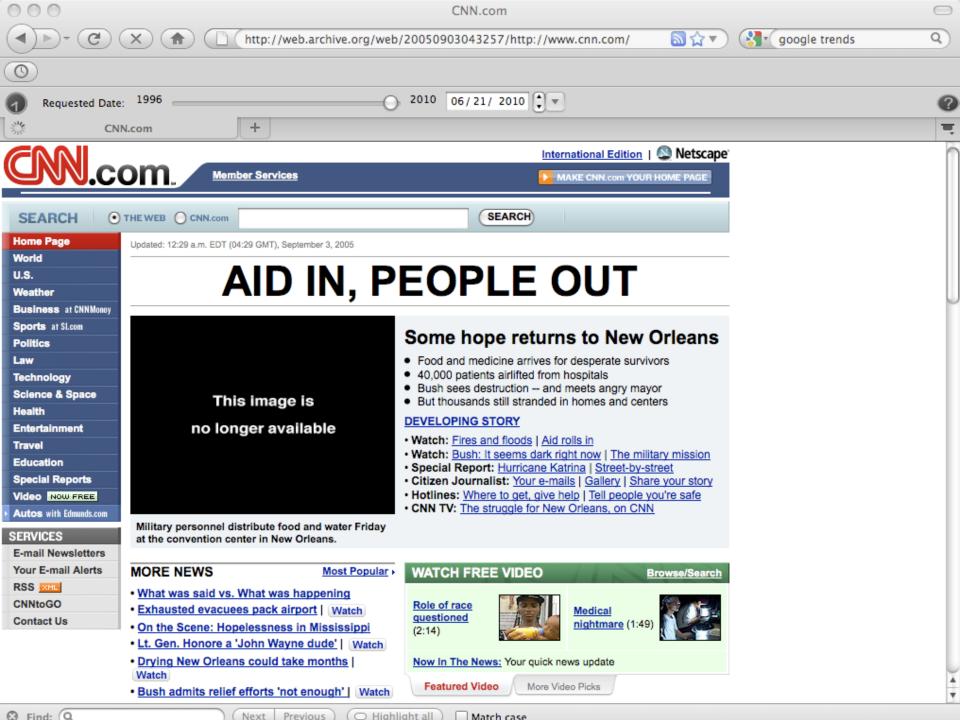


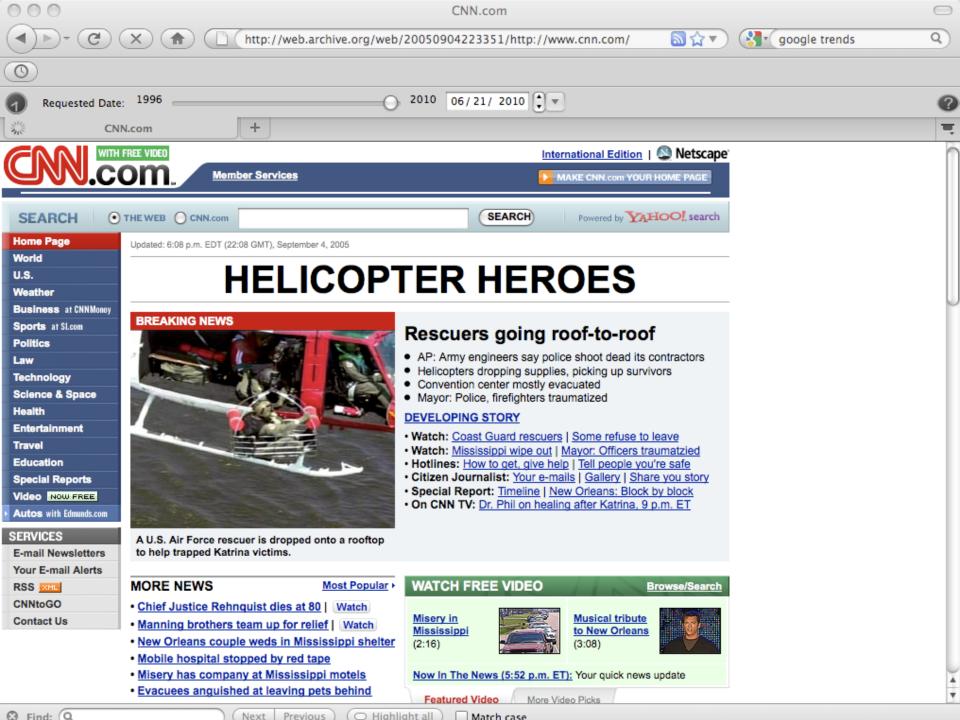












Replaying the past can be more compelling than just a summary

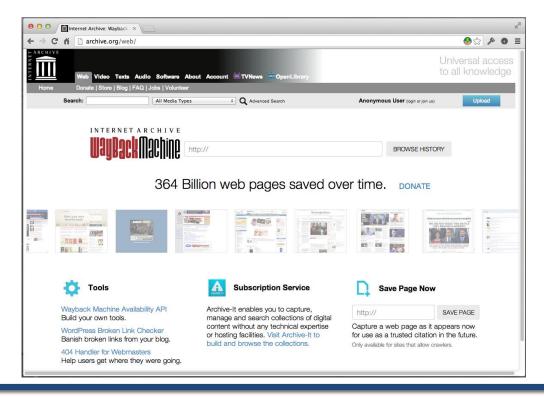
Overview of Current Projects

- What is archived?
 - how much of the web is archived?
 - how much of the Arabic web is archived?
 - what domains does each archive hold?
- How do people use web archives?
 - access patterns for humans and robots
 - what languages do they look for?
- How well are web pages archived?
 - archivability of webpages
 - damage of archived pages
 - tools to allow users to create their own archives
- Improving access for users
 - telling stories with web archives
 - archive visualization



How do people use web archives?

- We obtained a year's worth (2012) of requests to the Internet Archive's Wayback Machine
 - client IPs anonymized





How do people use web archives?

- First, there are a lot of robots (*aka* bots) who access the archive
 - 10 bot sessions for every 1 human session
 - maybe people don't know about the archive?
- Typical human sessions are pretty short
 - people aren't spending lots of time in the archive



AlNoamany, Weigle, and Nelson, "Access Patterns for Robots and Humans in Web Archives", JCDL 2013

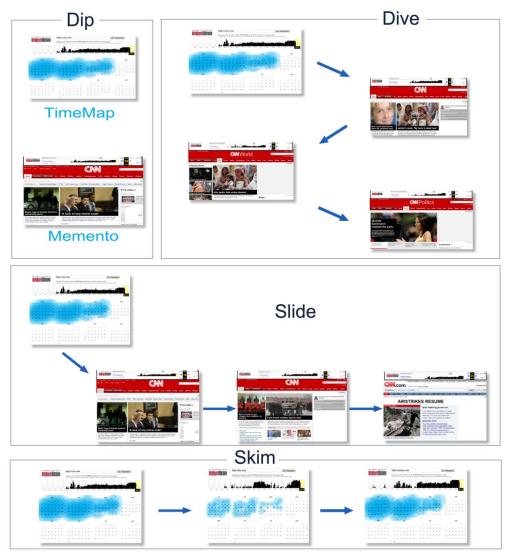
How do people use web archives?

• 65% of the requested archived pages no longer exist on the live web

• People use the archive because the pages they are interested in no longer exist

AlNoamany, AlSum, Weigle, and Nelson, "Who and What Links to the Internet Archive", IJDL, to appear, 2013

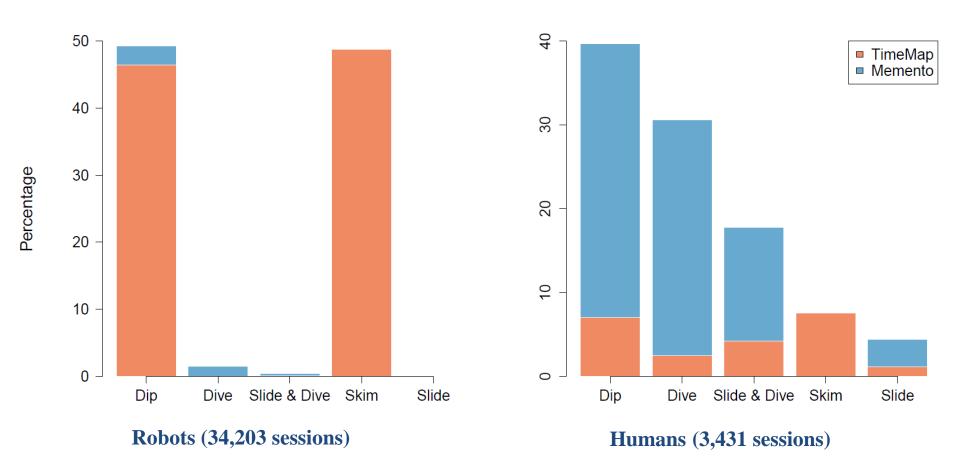
User Access Patterns



AlNoamany, Weigle, and Nelson, "Access Patterns for Robots and Humans in Web Archives", JCDL 2013



Everybody Dips, Humans Dive, Robots Skim



AlNoamany, Weigle, and Nelson, "Access Patterns for Robots and Humans in Web Archives", JCDL 2013



What domains does each archive hold?

	Archive Name	FullText	Website
		search	
$\mathbf{I}\mathbf{A}$	Internet Archive		web.archive.org
LoC	Library of Congress		www.loc.gov/lcwa
\mathbf{IC}	Icelandic Web Archive		vefsafn.is
CAN	Library & Archives Canada	x	www.collectionscanada.gc.ca
BL	British Library	x	www.webarchive.org.uk/ukwa
$\mathbf{U}\mathbf{K}$	UK Gov. Web Archive	x	webarchive.nationalarchives.gov.uk
PO	Portuguese Web Archive	x	arquivo.pt
CAT	Web Archive of Catalonia	x	www.padi.cat
CR	Croatian Web Archive	x	haw.nsk.hr
$\mathbf{C}\mathbf{Z}$	Archive of the Czech Web	x	webarchiv.cz
TW	National Taiwan University	x	webarchive.lib.ntu.edu.tw
AIT	Archive-It	x	www.archive-it.org

AlSum, Weigle, Nelson and Van de Sompel, "Profiling Web Archive Coverage for Top-Level Domain and Content Language," TPDL 2013.



What domains does each archive hold? ke uy uk ΖW net jр AIT CAT BL CAN edu eg de net us edu net es org gov orq ma com gov com com uk org ca cat 0.5 0.5 0.5 0.5 fг es it net uk uk CR CZ PO TW de bг com de ΓU net edu net cn net org org my рt com com org hг CZ com tw 0.5 0.5 0.5 0.5 net ca au ca pΙ gov pe gov fг edu eg iг ca ГU no jp it it it at UK IΑ IC LOC net to au de edu net cu uk net uk edu

AlSum, Weigle, Nelson and Van de Sompel, "Profiling Web Archive Coverage for Top-Level Domain and Content Language," TPDL 2013.

org

com

0.5

is



0.5

orq

de

com

orq

com

uk

June 26, 2014 34

gov

org

com

0.5

0.5

Sometimes the live web "leaks" into the archive



http://ws-dl.blogspot.com/2012/10/2012-10-10-zombies-in-archives.html



Archive What I See Now

- Standard web archiving tools are difficult for non IT experts. **Archive**
- "Save Page As" is not suitable for archiving purposes.
- Pages are behind authentication.
- Pages change quickly, but current state needs archiving.









http://bit.ly/wc-wail



WARCreate and WAIL

WARCreate



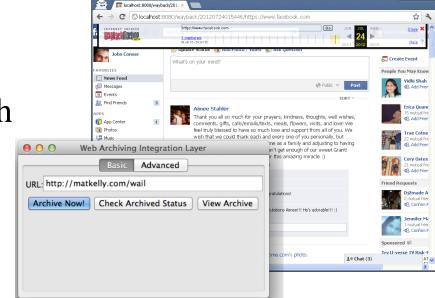
 generate archive of any webpage, right from your browser



• WAIL



- gathers institutional-strength archiving tools into simple package/GUI
- run your own WaybackMachine



http://bit.ly/wc-wail



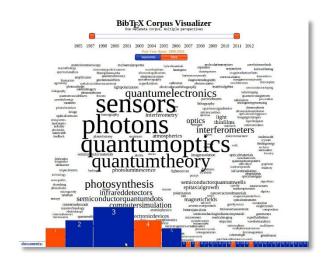
Outline

- Bits of Research
 - digital preservation and web archiving
 - information visualization

- Potential Summer Projects
 - Blue Button visualization (health informatics)
 - visualizing aggregate health data
 - exploring large document collections

What is Information Visualization?

- "The purpose of information visualization is insight, not pictures" -Ben Shneiderman
- Interactive representation of non-physically based data
- Concerned with how people think, human perception
- Has ties to cognitive science, psychology, data mining





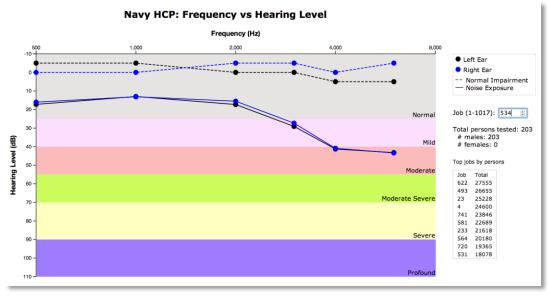
Info Vis Research Is Collaborative

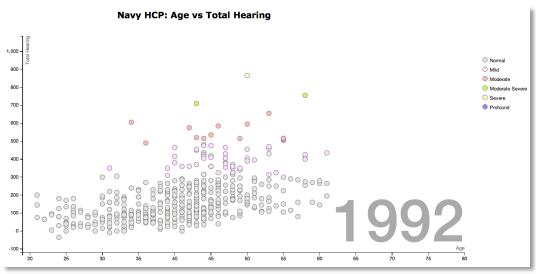
- Must have data and problems to solve
- Must work with experts in other domains
 - medical doctors hearing data analysis
 - social scientists welfare data record linkage
 - technology analysts trends in scientific research
 - web archivists viewing collections of web archives



Navy Hearing Conservation

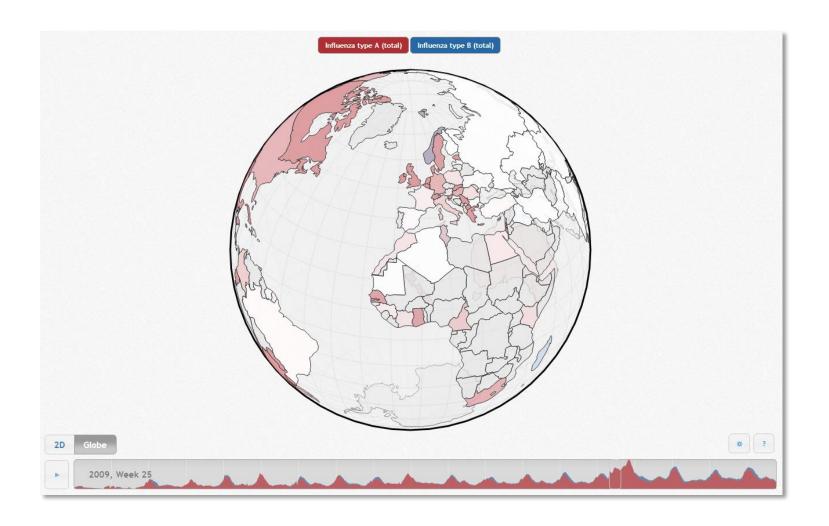
Program





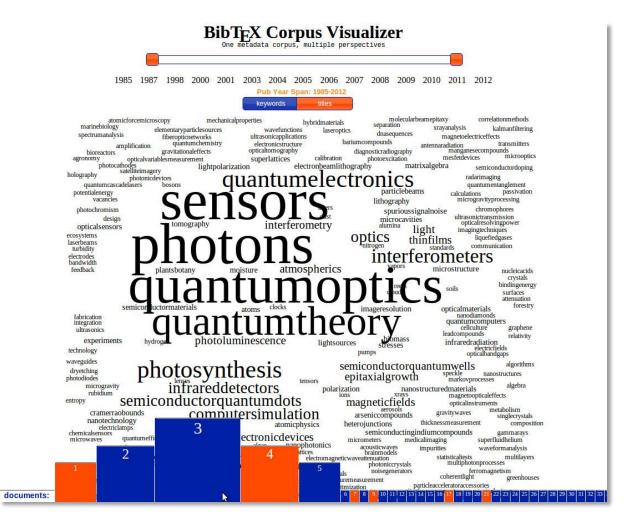


FluNet Visualization



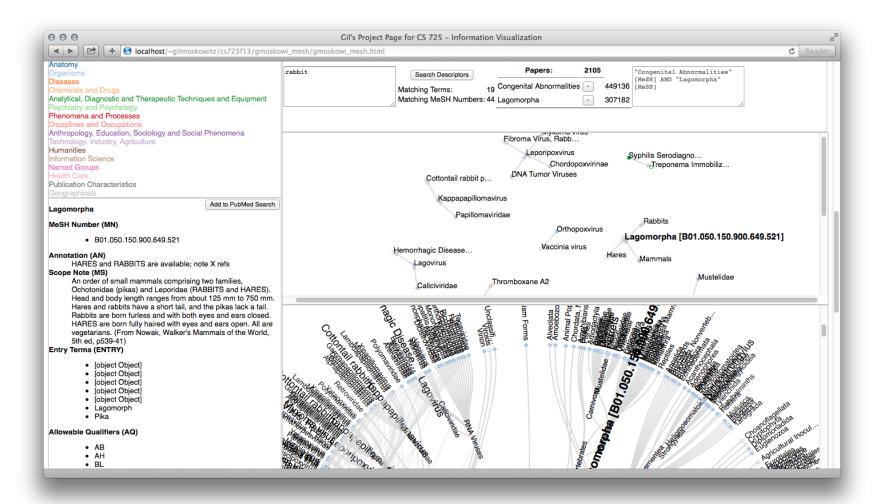


Academic Paper Browser





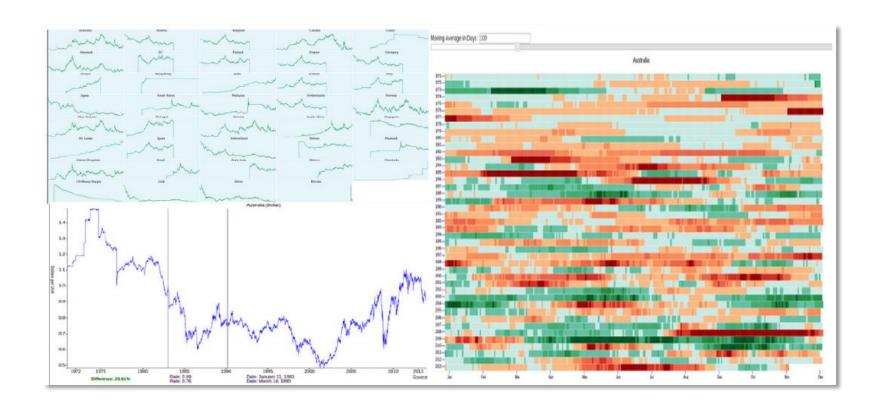
MeSH Viewer





June 26, 2014

Visualizing Currency Volatility



Outline

- Bits of Research
 - digital preservation and web archiving
 - information visualization

- Potential Summer Projects
 - Blue Button visualization (health informatics)
 - visualizing aggregate health data
 - exploring large document collections

Blue Button+ Visualization

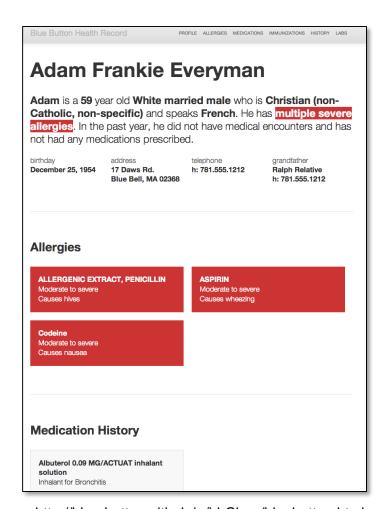
 Blue Button - allows patients to download their health data





• Goals:

- perform research into previous approaches to visualizing Blue Button data
- begin the development of a novel visualization of Blue Button data



http://blue-button.github.io/bbClear/bluebutton.html

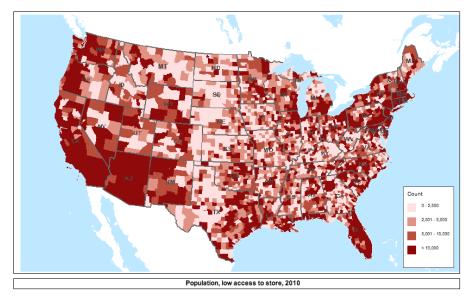


Visualizing Aggregate Health Data

- Aggregate health data
 - statistics for a large number of people (by country, US state, US county, etc.)

• Goal:

 become familiar with using public datasets and web-based visualization tools



http://www.ers.usda.gov/data-products/food-environment-atlas/goto-the-atlas.aspx#.U6hHpY1dUhs

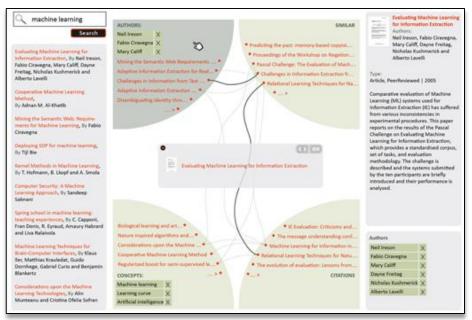
Exploring Large Document Collections

CORE

 API access to a large number of Open Access academic publication repositories

• Goals:

- investigate the use of the CORE API
- begin the implementation of tools that could be used in exploring large document collections



http://www.dlib.org/dlib/july12/herrmannova/07herrmannova.html



Outline

- Bits of Research
 - digital preservation and web archiving
 - information visualization

- Potential Summer Projects
 - Blue Button visualization (health informatics)
 - visualizing aggregate health data
 - exploring large document collections