

9-23-2013

# On the Change in Archivability of Websites Over Time

Mat Kelly

*Old Dominion University*

Justin F. Brunelle

Michele C. Weigle

*Old Dominion University, mweigle@odu.edu*

Michael L. Nelson

*Old Dominion University, mnelson@odu.edu*

Follow this and additional works at: [https://digitalcommons.odu.edu/computerscience\\_presentations](https://digitalcommons.odu.edu/computerscience_presentations)



Part of the [Archival Science Commons](#)

---

## Recommended Citation

Kelly, Mat; Brunelle, Justin F.; Weigle, Michele C.; and Nelson, Michael L., "On the Change in Archivability of Websites Over Time" (2013). *Computer Science Presentations*. 14.

[https://digitalcommons.odu.edu/computerscience\\_presentations/14](https://digitalcommons.odu.edu/computerscience_presentations/14)

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](mailto:digitalcommons@odu.edu).

# On the Change in Archivability of Websites Over Time

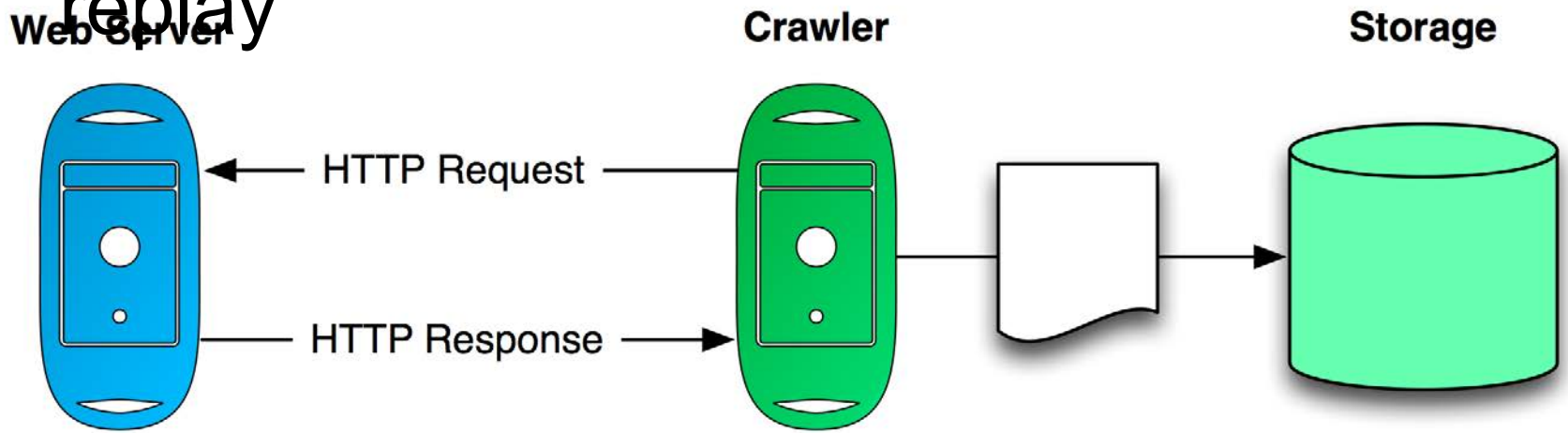
Mat Kelly, Justin F. Brunelle, Michele C. Weigle, and Michael L. Nelson

Old Dominion University

{mkelly,jbrunelle,mweigle,mIn}@cs.odu.edu

# Preserving Web Pages

- Identify page content needed for re-render
- Save page contents
  - HTML, Images, CSS, JavaScript
- Rewrite inter-resource references for  
replay



# Ease of Archiving

- **HTML** – text-based, references **other resources**
- Images
- CSS
- JavaScript

## EXAMPLE

---

```
<html><body>  
    
</body></html>
```

# Ease of Archiving

- HTML
- **Images** – binary data, no embedded URIs
- CSS
- JavaScript

## EXAMPLE

---



# Ease of Archiving

- HTML
- Images
- **CSS** – text-based, references other resources
- JavaScript

## EXAMPLE

---

```
body {  
  margin: 5px;  
  color: black;  
  background-image: url('outerSpace.png');  
}
```

# Ease of Archiving

- HTML
- Images
- CSS
- **JavaScript** – text-based, references other resources, URIs (not necessarily known until runtime)

## EXAMPLE

---

```
$.ajax({  
  url: "meatball-logo" + "." + "png"; //build logo URI at runtime  
});
```

# Ease of Archiving



- HTML
- Images
- CSS
- **JavaScript** – text-based, references other resources, URIs (not necessarily known until runtime)

---

## EXAMPLE

```
$.ajax({  
  url: "wormlogo" + "." + "png"; //build logo URI at runtime  
});
```

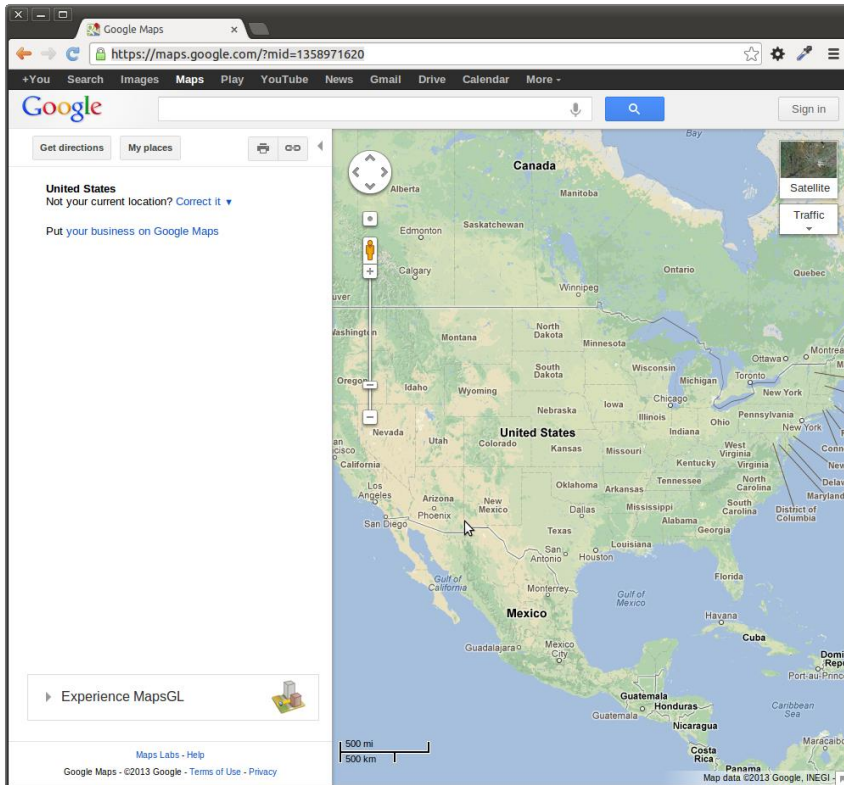


# BIG Problem

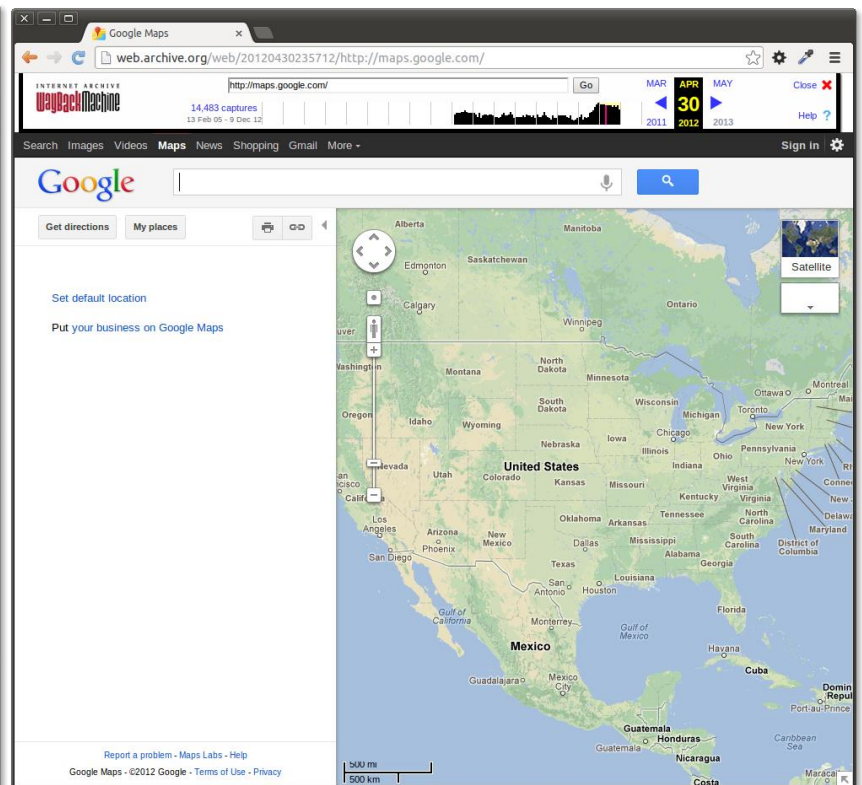
- JavaScript meant for browser
  - Capable of JS Execution
- Crawlers only recently became capable of exec
- Crawlers getting smarter at extracting URIs
  - Nowhere near perfect

# Identifying Missing Resources

- Sometimes missing resources are subtle



From Live Web

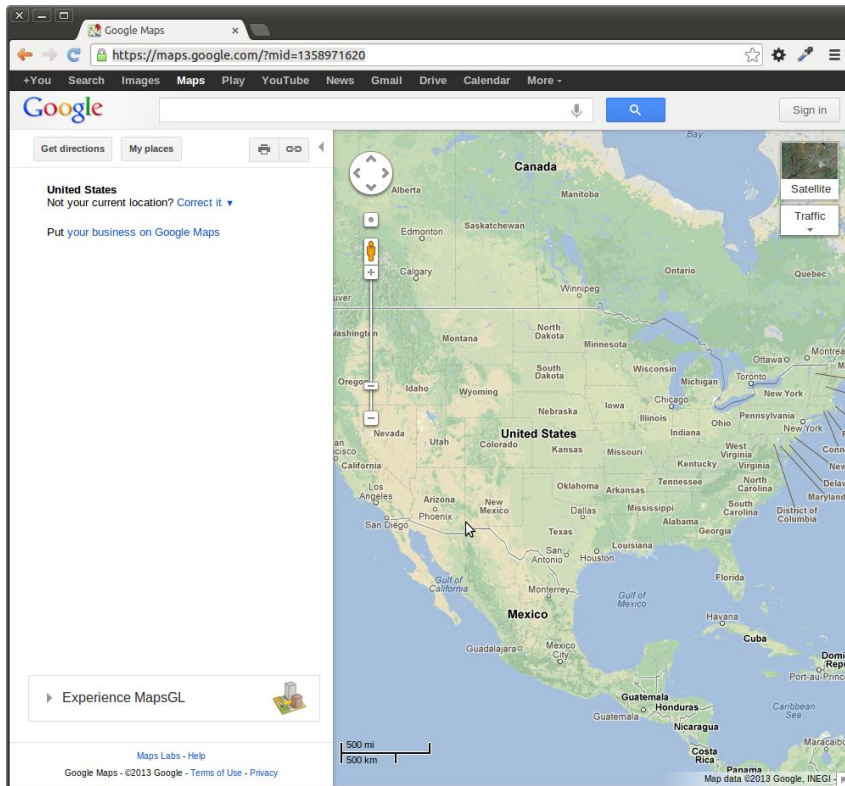


From Internet Archive Apr. 30, 2012

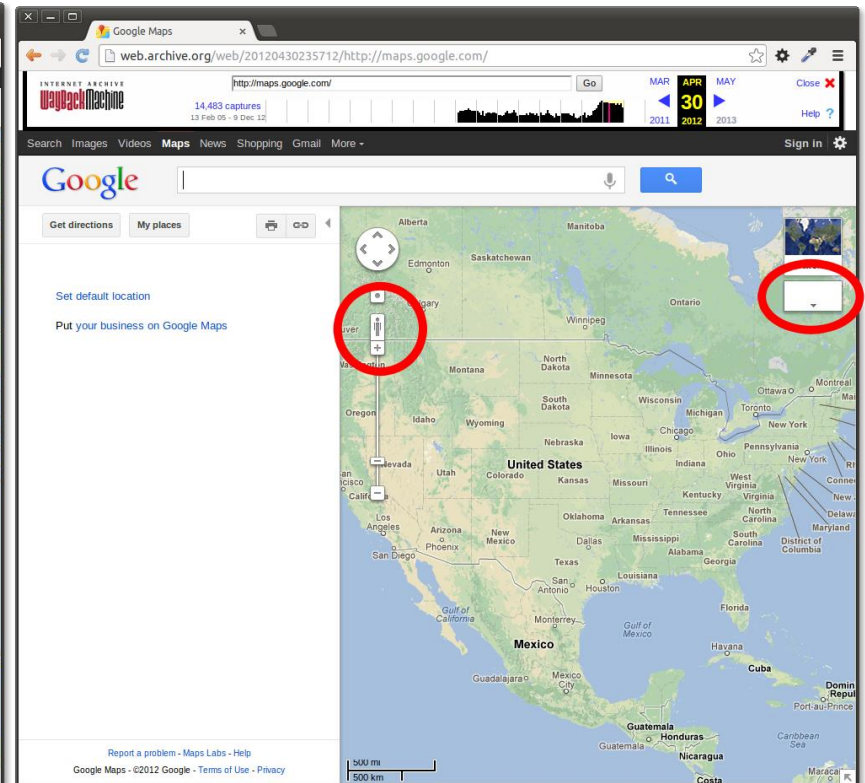
<http://web.archive.org/web/20120430235712/http://maps.google.com>

# Identifying Missing Resources

- Archived version not interactive & resources are missing



From Live Web

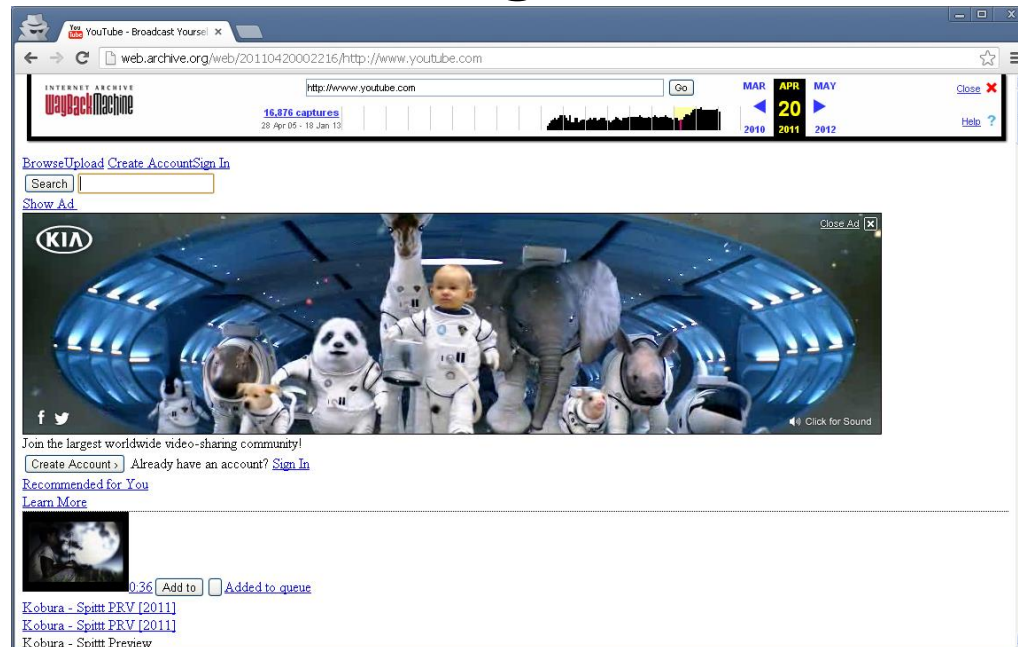


From Internet Archive Apr. 30, 2012

<http://web.archive.org/web/20120430235712/http://maps.google.com>

# Identifying Missing Resources

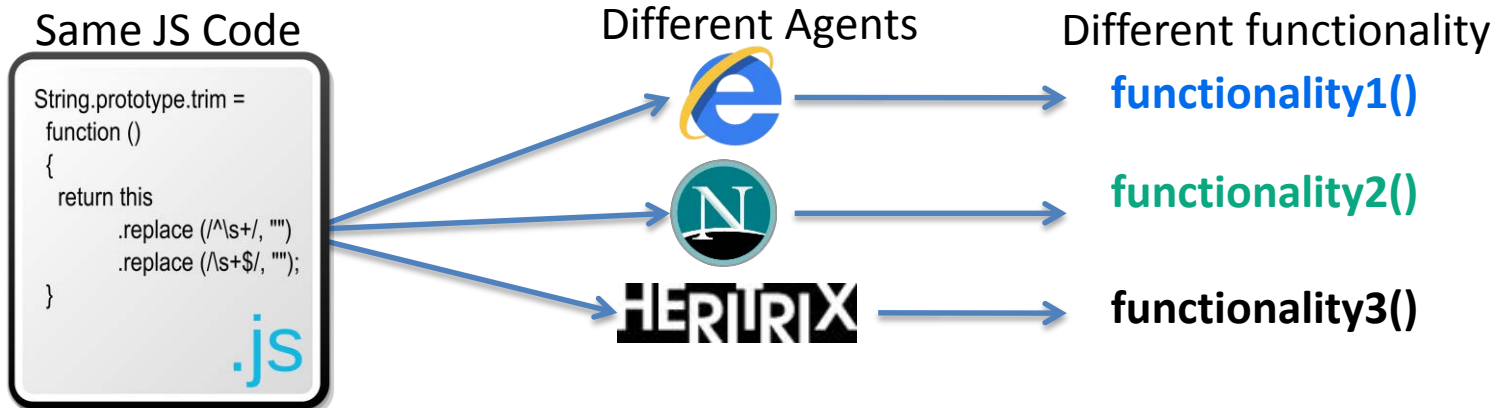
- Other times, a failed AJAX call prevents other resources from loading



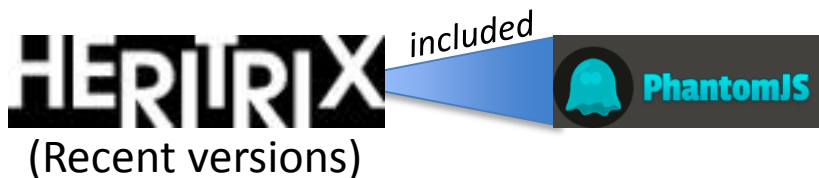
YouTube (2011) in archives with failed Ajax call

<http://web.archive.org/web/20110420002216/http://www.youtube.com/>

# Why JavaScript makes it difficult



- Archival crawlers don't interpret DOM, are made for capture and thus faster



- Interprets JavaScript
- Enabled pages w/ JS to be archived better!

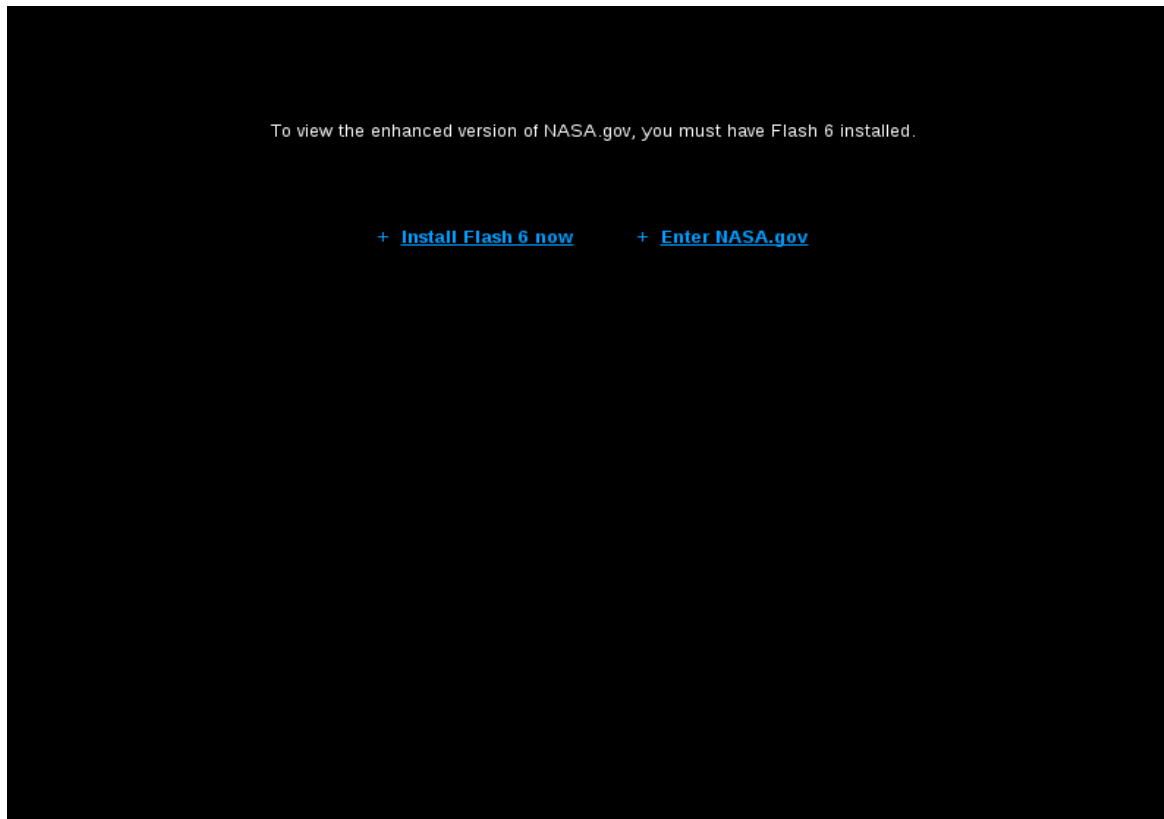
# JavaScript and Accessibility

- Content displaying should not be dependent on script execution



- U.S. Government
  - required to comply with accessibility standards
  - Content should be available to crawler w/o JS
  - Therefore, government sites are better preserved
    - Right?

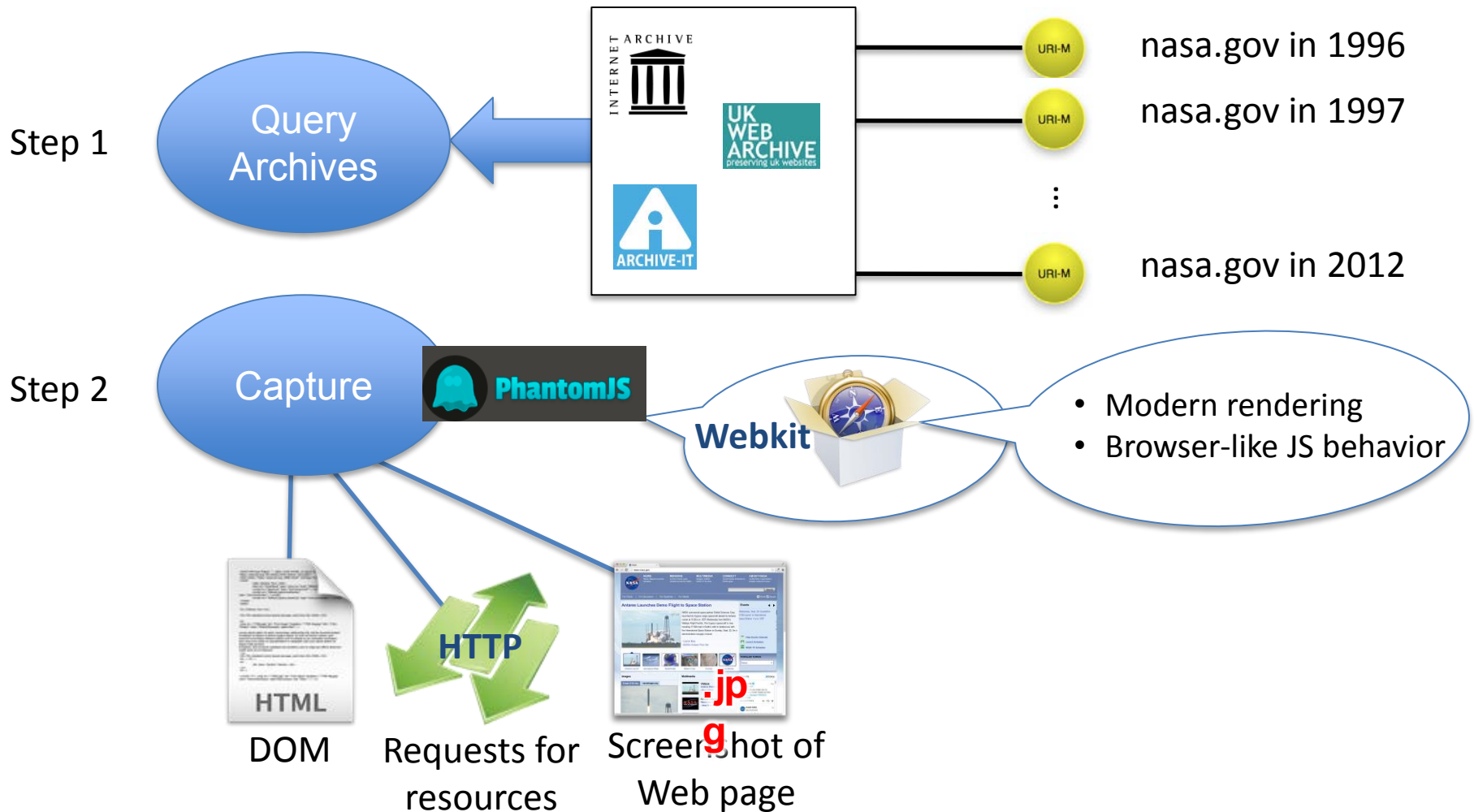
# If not Accessible then not Archivable



NASA.gov in Internet Archive, 2004

<http://web.archive.org/web/20041014205942/http://www.nasa.gov>

# Finding out where NASA went wrong





# NASA.gov over time



# NASA.gov over time



# Does this apply to popular sites?

- Alexa's Top 10 websites

Alexa Rank	Web Site Name
1	Facebook.com
2	Google.com
3	YouTube.com
4	Yahoo.com
5	Baidu.com
6	Wikipedia.org
7	Live.org
8	Amazon.com
9	QQ.com
10	Twitter.com

# Does this apply to popular sites?

- Alexa's Top 10 websites
- **No Mementos**
  - robots.txt exclusion prevents crawl

Alexa Rank	Web Site Name
1	Facebook.com
2	Google.com
3	YouTube.com
4	Yahoo.com
5	Baidu.com
6	Wikipedia.org
7	Live.org
8	Amazon.com
9	QQ.com
10	Twitter.com

# Does this apply to popular sites?

Alexa Rank	Web Site Name	Sampled Mementos
1	Facebook.com	No memento robots.txt exclusion
2	Google.com	15 mementos 1998 to 2012
3	YouTube.com	7 mementos 2006 to 2012
4	Yahoo.com	16 mementos 1997 to 2012
5	Baidu.com	No memento robots.txt exclusion
6	Wikipedia.org	12 mementos 2001 to 2012
7	Live.org	15 mementos 1999 to 2012
8	Amazon.com	14 mementos 1999 to 2012
9	QQ.com	15 mementos 1998 to 2012
10	Twitter.com	No memento robots.txt exclusion

all thumbnails at: [http://www.cs.odu.edu/~mkelly/semester/2013\\_spring/20130127alexatop10/](http://www.cs.odu.edu/~mkelly/semester/2013_spring/20130127alexatop10/)

# Case Study with Ajax emphasis: YouTube

Alexa Rank	Web Site Name	Sampled Mementos
1	Facebook.com	No memento robots.txt exclusion
2	Google.com	15 mementos 1998 to 2012
3	YouTube.com	7 mementos 2006 to 2012
4	Yahoo.com	16 mementos 1997 to 2012
5	Baidu.com	No memento robots.txt exclusion
6	Wikipedia.org	12 mementos 2001 to 2012
7	Live.org	15 mementos 1999 to 2012
8	Amazon.com	14 mementos 1999 to 2012
9	QQ.com	15 mementos 1998 to 2012
10	Twitter.com	No memento robots.txt exclusion

# Ajax's Effects on the Archivability of YouTube

- Recently Viewed content missing

YouTube Broadcast Yourself

Sign Up | Log In | Help

Videos Search

Home Videos Channels Groups Members Upload

My Videos | My Favorites | My Messages | My Subscriptions | My Playlists | My Groups | My Profile

Recently Viewed [1 - 5 of 12] >> More Recently Viewed...

Football skills 1 second ago

the police roxanne 3 seconds ago

World Trade Center 9/11 3 seconds ago

I dont wanna miss a t... 3 seconds ago

cuando calienta el sol 3 seconds ago

Today's Featured Videos... See More Videos

**F-16 crash**  
filmed from the outside  
Tags // cool  
Channels // Entertainment : Humor : People  
Added: 3 days ago by fatbob79  
Runtime: 00:48 | Views: 16009 | Comments: 81  
★★★★☆ (116 ratings)

**STANGERS WITH CANDY - Teaser Trailer - Exclusive First Look**  
ThinkFilm presents an Exclusive first look at the Strangers With Candy movie trailer. Keep an eye out, the film hits theaters June 28th.  
Tags // strangers : candy : amy : sedaris : david : gay : lesbian : humor : comedy : cult : classic : central : steven : colbert : funny  
Channels // Entertainment : Humor : Short Movies

**New Features**  
**Custom Profiles**  
Choose your colors and interact with comments and bulletin boards  
**Members**  
New Members tab—improved search, featured members, and more!  
**Explore YouTube**

**Pretty Girls Make Graves Video Contest Winners**

**Join the Cybersmack Group**

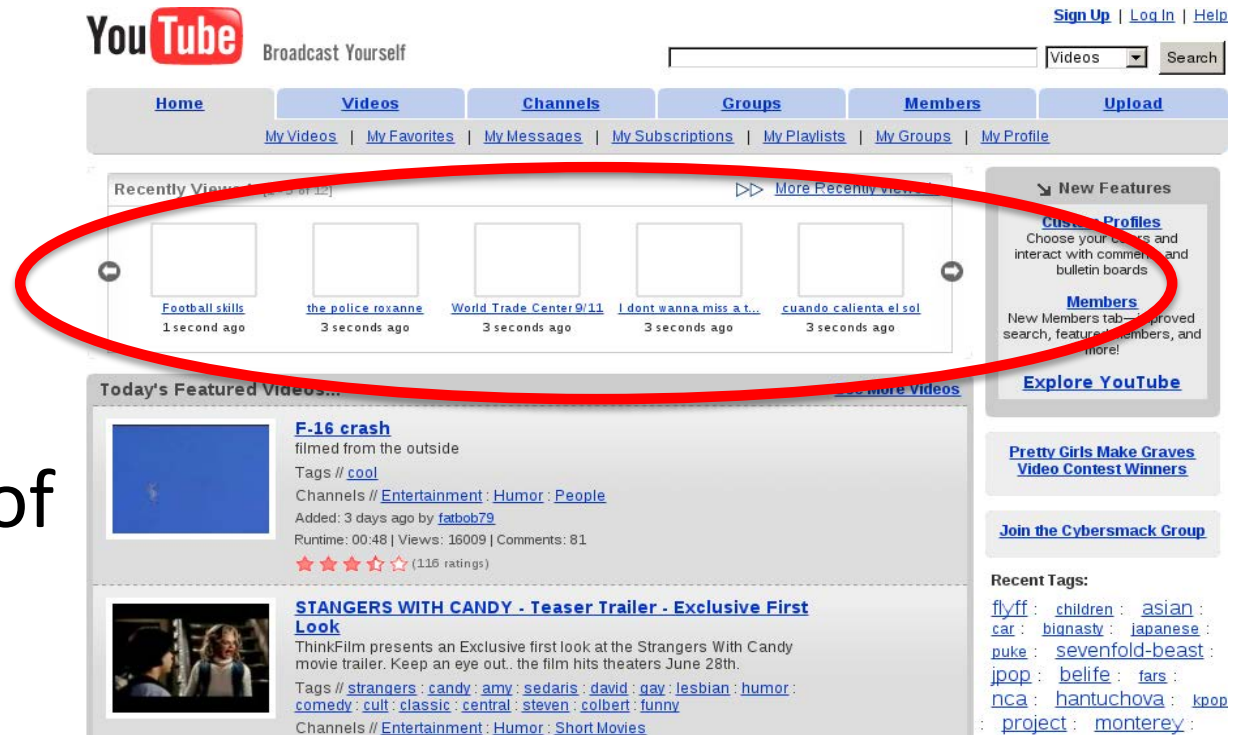
**Recent Tags:**  
flyff : children : asian : car : bignasty : japanese : puke : sevenfold-beast : jpop : belife : fars : nca : hantuchova : kpop : project : monterey :

2006 YouTube.com from Internet Archive

<http://web.archive.org/web/20060427213420/http://youtube.com>

# Ajax's Effects on the Archivability of YouTube

- Recently Viewed content missing
- How much of this is because of JavaScript?



2006 YouTube.com from Internet Archive

<http://web.archive.org/web/20060427213420/http://youtube.com>



# YouTube as captured without JavaScript capability

- Titles fetched with JavaScript
- Primary content (preview) still missing

The screenshot shows the YouTube homepage from 2006. At the top left is the YouTube logo with the tagline "Broadcast Yourself". To the right are links for "Sign Up", "Log In", and "Help". Below the logo is a search bar with a dropdown menu set to "Videos" and a "Search" button. A navigation bar contains tabs for "Home", "Videos", "Channels", "Groups", "Members", and "Upload". Underneath these tabs are links for "My Videos", "My Favorites", "My Messages", "My Subscriptions", "My Playlists", "My Groups", and "My Profile".

The main content area is divided into several sections:

- Recently Viewed:** A horizontal list of five video thumbnails, each with the text "loading..." below it. A "More Recently Viewed..." link is on the right.
- Today's Featured Videos...:** A section with two video entries. The first is "F-16 crash" with a blue thumbnail, tags like "cool", and a description. The second is "STANGERS WITH CANDY - Teaser Trailer - Exclusive First Look" with a movie trailer thumbnail, tags like "strangers", "candy", "amy", "sedaris", "david", "gay", "lesbian", "humor", "comedy", "cult", "classic", "central", "steven", "colbert", "funny", and a description.
- New Features:** A sidebar on the right with sections for "Custom Profiles", "Members", and "Explore YouTube".
- Recent Tags:** A list of tags such as "flyff", "children", "asian", "car", "bignasty", "japanese", "puke", "sevenfold-beast", "jpop", "belife", "fars", "nca", "hantuchova", "kpop", "project", "monterey".

2006 YouTube.com from Internet Archive

<http://web.archive.org/web/20060427213420/http://youtube.com>

# Dependent Loading of Resources

GET  
[http://web.archive.org/web/20121208145112cs\\_/http://s.ytimg.com/yt/cssbin/www-core-vfl\\_OJqFG.css](http://web.archive.org/web/20121208145112cs_/http://s.ytimg.com/yt/cssbin/www-core-vfl_OJqFG.css) 404 (Not Found) www.youtube.com:15

GET  
[http://web.archive.org/web/20121208145115js\\_/http://s.ytimg.com/yt/jsbin/www-core-vfl8PDcRe.js](http://web.archive.org/web/20121208145115js_/http://s.ytimg.com/yt/jsbin/www-core-vfl8PDcRe.js) 404 (Not Found) www.youtube.com:45

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:56

Uncaught TypeError: Cannot read property 'home' of undefined www.youtube.com:76

Uncaught TypeError: Cannot read property 'ajax' of undefined www.youtube.com:86

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:101

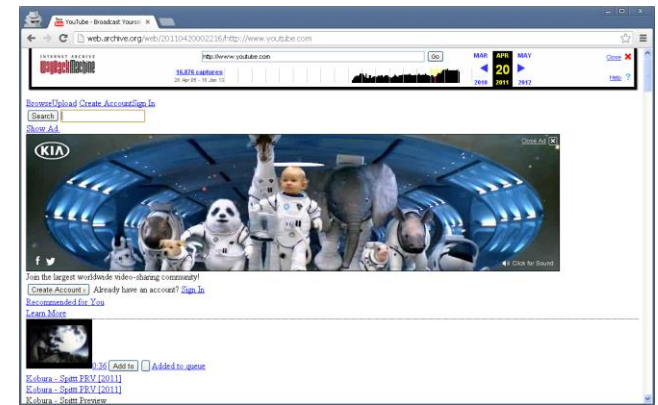
Uncaught ReferenceError: \_gel is not defined www.youtube.com:1784

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:1929

Uncaught TypeError: Cannot read property 'home' of undefined www.youtube.com:524

GET  
[http://web.archive.org/web/20130101024721im\\_/http://i2.ytimg.com/vi/1f7neSzDqvc/default.jpg](http://web.archive.org/web/20130101024721im_/http://i2.ytimg.com/vi/1f7neSzDqvc/default.jpg) 404 (Not Found)

YouTube 2011



<http://web.archive.org/web/20110420002216/http://www.youtube.com/>

# Dependent Loading of Resources

GET  
[http://web.archive.org/web/20121208145112cs\\_/http://s.ytimg.com/yt/cssbin/www-core-vfl\\_OJqFG.css](http://web.archive.org/web/20121208145112cs_/http://s.ytimg.com/yt/cssbin/www-core-vfl_OJqFG.css) 404 (Not Found) www.youtube.com:15

GET  
[http://web.archive.org/web/20121208145115js\\_/http://s.ytimg.com/yt/jsbin/www-core-vfl8PDcRe.js](http://web.archive.org/web/20121208145115js_/http://s.ytimg.com/yt/jsbin/www-core-vfl8PDcRe.js) 404 (Not Found) www.youtube.com:45

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:56

Uncaught TypeError: Cannot read property 'home' of undefined www.youtube.com:76

Uncaught TypeError: Cannot read property 'ajax' of undefined www.youtube.com:86

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:101

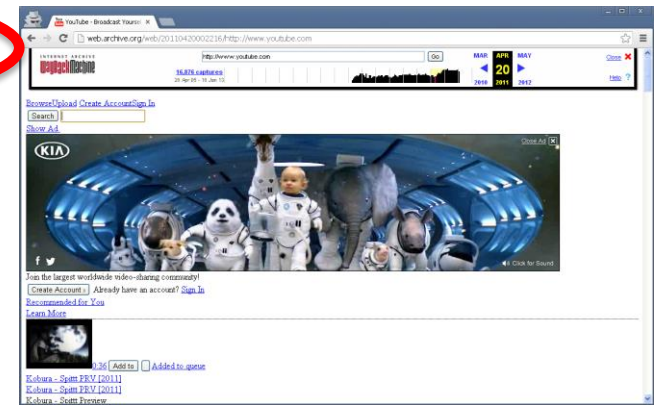
Uncaught ReferenceError: \_gel is not defined www.youtube.com:1784

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:1929

Uncaught TypeError: Cannot read property 'home' of undefined www.youtube.com:524

GET  
[http://web.archive.org/web/20130101024721im\\_/http://i2.ytimg.com/vi/1f7neSzDqvc/default.jpg](http://web.archive.org/web/20130101024721im_/http://i2.ytimg.com/vi/1f7neSzDqvc/default.jpg) 404 (Not Found)

YouTube 2011



<http://web.archive.org/web/20110420002216/http://www.youtube.com/>

**Missing CSS**

Missing JS

Missing JS-  
dependent  
resources

# Dependent Loading of Resources

YouTube 2011

GET  
[http://web.archive.org/web/20121208145112cs\\_/http://s.ytimg.com/yt/cssbin/www-core-vii\\_UJqFG.css](http://web.archive.org/web/20121208145112cs_/http://s.ytimg.com/yt/cssbin/www-core-vii_UJqFG.css) 404 (Not Found) www.youtube.com:15

GET  
[http://web.archive.org/web/20121208145115js\\_/http://s.ytimg.com/yt/jsbin/www-core-vf18PDcRe.js](http://web.archive.org/web/20121208145115js_/http://s.ytimg.com/yt/jsbin/www-core-vf18PDcRe.js) 404 (Not Found) www.youtube.com:45

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:56

Uncaught TypeError: Cannot read property 'home' of undefined www.youtube.com:76

Uncaught TypeError: Cannot read property 'ajax' of undefined www.youtube.com:86

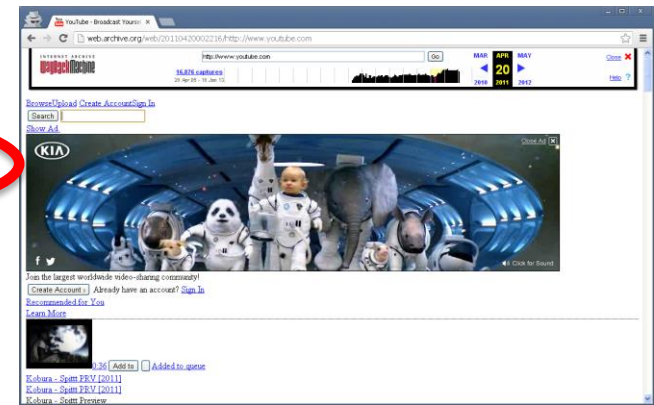
Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:101

Uncaught ReferenceError: \_gel is not defined www.youtube.com:1784

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:1929

Uncaught TypeError: Cannot read property 'home' of undefined www.youtube.com:524

GET  
[http://web.archive.org/web/20130101024721im\\_/http://i2.ytimg.com/vi/1f7neSzDqvc/default.jpg](http://web.archive.org/web/20130101024721im_/http://i2.ytimg.com/vi/1f7neSzDqvc/default.jpg) 404 (Not Found)



<http://web.archive.org/web/20110420002216/http://www.youtube.com/>

Missing CSS  
**Missing JS**  
Missing JS-  
dependent  
resources



# Dependent Loading of Resources

YouTube 2011

GET  
http://web.archive.org/web/20121208145112cs\_/http://s.ytimg.com/yt/cssbin/www-core-vfl\_OJqFG.css 404 (Not Found) www.youtube.com:15

GET  
http://web.archive.org/web/20121208145115js\_/http://s.ytimg.com/yt/jsbin/www-core-vfl8PDcRe.js 404 (Not Found) www.youtube.com:45

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:56

Uncaught TypeError: Cannot read property 'home' of undefined www.youtube.com:76

Uncaught TypeError: Cannot read property 'ajax' of undefined www.youtube.com:86

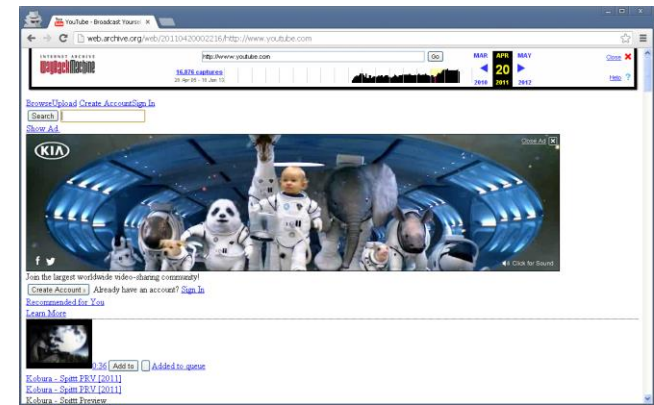
Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:101

Uncaught ReferenceError: \_gel is not defined www.youtube.com:1784

Uncaught TypeError: Object #<Object> has no method 'setConfig' www.youtube.com:1929

Uncaught TypeError: Cannot read property 'home' of undefined www.youtube.com:524

GET  
http://web.archive.org/web/20130101024721im\_/http://i2.ytimg.com/vi/1f7neszDqvc/default.jpg 404 (Not Found)



http://web.archive.org/web/20110420002216/http://www.youtube.com/

Missing CSS

Missing JS

Missing JS-  
dependent  
resources

Address	Status
▼ YouTube - Broadcast Yourself.	
▶ about:blank	53 items
▶ http://web.archive.org/liveweb/http://ad-g.double...eative_1;tile=1;dcopt=ist;ord=7629290234443648	1 item
▶ http://web.archive.org/web/20110420002216/http://ad-creative_3;tile=3;ord=9792874818943184	7 items
http://i1.ytimg.com/vi/4ALrBtVxn0A/default.jpg	2.9 KB
http://i1.ytimg.com/vi/I55jZclGDsl/default.jpg	not found
http://i2.ytimg.com/vi/-CrcLolPzow/default.jpg	4.0 KB
http://i2.ytimg.com/vi/Q2G4ETZvjU/default.jpg	2.5 KB
http://i3.ytimg.com/vi/2zw8SmsovjC/default.jpg	4.6 KB
http://i3.ytimg.com/vi/6ckaW1WWmZo/default.jpg	3.3 KB
http://i3.ytimg.com/vi/bb6cBKE3WzQ/default.jpg	3.4 KB
http://i3.ytimg.com/vi/jrN41AUXKo4/default.jpg	3.8 KB
http://i4.ytimg.com/vi/C9jghLeYufQ/default.jpg	3.0 KB
http://i4.ytimg.com/vi/KqCjC6cvjec/default.jpg	2.9 KB
http://s.ytimg.com/yt/swf/masthead_child-vflRMMO6_swf	2.7 KB
http://web.archive.org/static/images/toolbar/transp-red-pixel.png	0.1 KB
http://web.archive.org/static/images/toolbar/transp-yellow-pixel.png	70 bytes
http://web.archive.org/static/images/toolbar/wayback-toolbar-logo.png	7.4 KB
http://web.archive.org/static/images/toolbar/wm_tb_bk_trns.png	0.1 KB
http://web.archive.org/static/images/toolbar/wm_tb_close.png	0.1 KB
http://web.archive.org/static/images/toolbar/wm_tb_help.png	0.1 KB
http://web.archive.org/static/images/toolbar/wm_tb_nxt_on.png	0.1 KB
http://web.archive.org/static/images/toolbar/wm_tb_prv_on.png	0.1 KB
http://web.archive.org/static/jquery/jquery.min.js	70 KB
http://web.archive.org/static/js/disclaim-element.js	0.1 KB
http://web.archive.org/static/js/graph-calc.js	0.1 KB
http://web.archive.org/static/js/jwplayer/jwplayer.js	0.1 KB
http://web.archive.org/static/js/video-embed-rewriter.js	0.1 KB
http://web.archive.org/web/20110419231106/http://s.ytimg.com/yt/imgbin/www-master-vflX0FKfU.png	0.1 KB
http://web.archive.org/web/20110419231135/http://pccs.9145.hpjs.10598.ct.10598.ol.93461.aft.93461	0.1 KB
http://web.archive.org/web/20110420002216/http://www.youtube.com/	0.1 KB
http://web.archive.org/web/20110420002216cs_/http://s.ytimg.com/yt/cssbin/www-core-vfl_OjQFG.css	202.1 KB
http://web.archive.org/web/20110420002216im_/http://i1.ytimg.com/vi/DF7Gd0YShio/default.jpg	4.6 KB
http://web.archive.org/web/20110420002216im_/http://i1.ytimg.com/vi/dPqhb8ouToY/default.jpg	4.5 KB
http://web.archive.org/web/20110420002216im_/http://i1.ytimg.com/vi/pfDEaciayoc/default.jpg	4.5 KB
http://web.archive.org/web/20110420002216im_/http://i1.ytimg.com/vi/phASMLiFRR8/default.jpg	3.0 KB
http://web.archive.org/web/20110420002216im_/http://i1.ytimg.com/vi/Xio61D2On14/default.jpg	not found
http://web.archive.org/web/20110420002216im_/http://i2.ytimg.com/vi/1f7neSzDqvc/default.jpg	2.5 KB
http://web.archive.org/web/20110420002216im_/http://i2.ytimg.com/vi/QYPgD6wlxco/default.jpg	4.2 KB
http://web.archive.org/web/20110420002216im_/http://i2.ytimg.com/vi/UgwJfoZ-12c/default.jpg	4.4 KB
http://web.archive.org/web/20110420002216im_/http://i3.ytimg.com/vi/BdPjocfSoll/default.jpg	4.7 KB
http://web.archive.org/web/20110420002216im_/http://i3.ytimg.com/vi/bMGatrWkG2c/default.jpg	4.4 KB
http://web.archive.org/web/20110420002216im_/http://i3.ytimg.com/vi/ftASuLqTYg/default.jpg	4.7 KB
http://web.archive.org/web/20110420002216im_/http://i3.ytimg.com/vi/NJOIR0sIR_c/default.jpg	3.3 KB
http://web.archive.org/web/20110420002216im_/http://i3.ytimg.com/vi/NjZ-4qQyw_s/default.jpg	4.7 KB
http://web.archive.org/web/20110420002216im_/http://i3.ytimg.com/vi/Z5U3QvKfUG8/default.jpg	3.2 KB
http://web.archive.org/web/20110420002216im_/http://i4.ytimg.com/vi/7G0SaC05V3Q/default.jpg	4.6 KB

# Live Web Leaks Into Archive Via Javascript





Sept 3, 2008

updated 4:24 p.m. EDT, Wed September 3, 2008

Make CNN Your Home Page



### Palin's path from city hall to governor's mansion

When she played basketball in high school, Sarah Palin, the soon-to-be Republican vice presidential nominee, earned the nickname "Sarah Barracuda" for her fierce competitiveness. In the 21 months she has served as governor of Alaska, no one is suggesting she's lost her fighting edge. full story

- Palin ready to lead? | Analysts: Her task
- Palin to call for reform in RNC speech
- Will her gender sway women to Palin?

HEALTH MAGAZINE

Secrets to a healthy heart

### Latest News

- Poll measures race in three key states 19 min
- Dems blast Lieberman, say he lied to delegates
- Ticker: McCain greets Palin's daughter, fiancé
- Martin: Keep politics out of Bristol Palin issue
- Biographer: America, meet Sarah Palin
- iReport.com: Can a mom of 5 be VP?
- Cafferty: Should McCain consider different VP?
- CNNMoney: Auto sales plunge
- Gunman's bloody trail leads to six bodies
- Caylee's grandma disputes evidence
- Nagin: Residents won't be turned away
- CNNMoney: Gulf oil production scrambling back
- Hanna roams Caribbean, may hit U.S.
- Pakistan protests 'coalition air attack'
- The cartoonist beloved by GIs and regular guys
- 8 dogs die in hot truck during lunch break
- Review: Google's Chrome needs more polish
- New Marilyn Monroe film scenes found
- Beyonce's sister aims for a path of her own
- CNN Wire: Paraguay coup plot alleged

all news from the past 24hrs » | all most popular »

### Republican National Convention »

- Bernstein: Democrats better take note
- Bush, McCain still an uneasy alliance
- Interactive: How conventions work
- McCain arrives in Minnesota for convention
- Key players | Trivia | The Forum | In Depth

Viewer's Guide

### Video »



Gustav recovery briefing

LIVE



Dogs die during lunch break 1:39



The RNC in St. Paul

LIVE

LIVE: Watch coverage of the RNC & Tropical Storm Hanna



2012

### CNN TV »



Palin in prime-time

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



# Stepping Back, NASA 1999

- Few embedded resources
- Little to no JavaScript
  - If some, not Ajax
- No resources dependent on JavaScript

Nov 9, 1999

"NASA is deeply committed to spreading the unique knowledge that flows from its aeronautics and space research..."

Read NASA Administrator Daniel S. Goldin's [welcome letter](#), to and [speeches](#).

[Welcome to NASA Web](#)

Do you dream of exploring space or working for NASA? If so, avoid black holes and drugs. [You decide.](#)

**Navigating NASA's Strategic Enterprises**

- [Aero-Space Technology](#)
- [Human Exploration and Development of Space](#)
- [Earth Science](#)
- [Space Science](#)

**NASA for Kids**

**More About NASA:**

- [Doing Business with NASA](#)
- [Educational Resources](#)
- [Freedom of Information Act](#)
- [History](#)
- [Jobs and Internships](#)
- [News and Information](#)
- [Organization and Subject Index](#)
- [Project Home Pages](#)
- [Research Opportunities](#)
- [Scientific and Technical Information](#)
- [See a Launch](#)
- [Launch Schedule](#)
- [Spinoffs and Commercial Technology](#)
- [Visiting NASA](#)

**NASA**

**SEARCH**  
the NASA web

**NASA Technology to Make Air Travel Safer, Reduce Delays**

Interested in the latest information NASA has to offer? Then take a look at [today@nasa.gov](mailto:today@nasa.gov). This on-line newsletter, updated daily, contains the latest news about NASA science and technology.

- [NASA Technology to Make Air Travel Safer, Reduce Delays](#)
- [A Different Route to High Altitude Research](#)
- [NASA Steps Up on Aviation Safety](#)

NASA and an industry partner have developed technology that could ease travelers' frustration with flight delays caused by bad weather. The Airborne Information for Lateral Spacing system and the Closely Spaced Parallel Approaches system expand on existing communication and navigation technology to allow airplanes to land safely in bad weather on parallel runways spaced as closely as 2,500 feet apart. Currently, the minimum runway separation during low visibility is 4,300 feet, which means that some of the nation's busiest airports have to shut down one of their closely spaced runways when weather conditions deteriorate. Some of the airports where this new technology could improve on-time arrivals are Detroit, Seattle, Minneapolis and Memphis. ([Full Story](#)) (11/8/99)

**Cool NASA Websites**

- [Galileo](#)
- [Kids' issues in Space](#)
- [Spaceflight](#)
- [Women & NASA](#)

**Other Cool NASA Websites**

- [\[Frequently Asked Questions\]](#) [\[Hot Topics\]](#) [\[Multimedia Gallery\]](#)
- [\[NASA Television\]](#) [\[Textonly Version\]](#) [\[NASA Privacy Statement\]](#) [\[Site Map\]](#)

Author: Brian Dunbar  
Curator: Sudha V. Chudamani  
[Comments and Questions](#)  
Last Updated: Nov 9, 1999

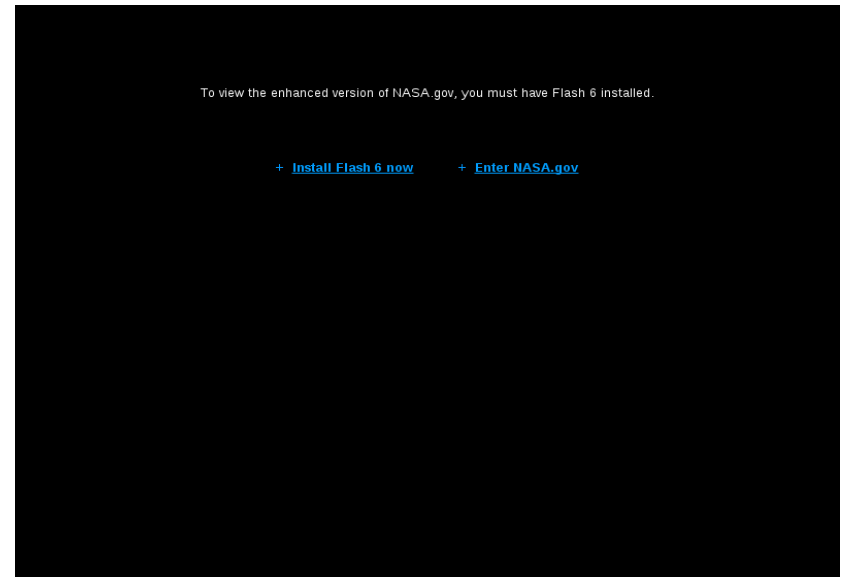
<http://web.archive.org/web/19991109105105/http://www.nasa.gov>



# Stepping Back, NASA 2003

<http://web.archive.org/web/20041014205942/http://www.nasa.gov>

- Content requires JS
- Accessibility goes awry
- Archivability decreases

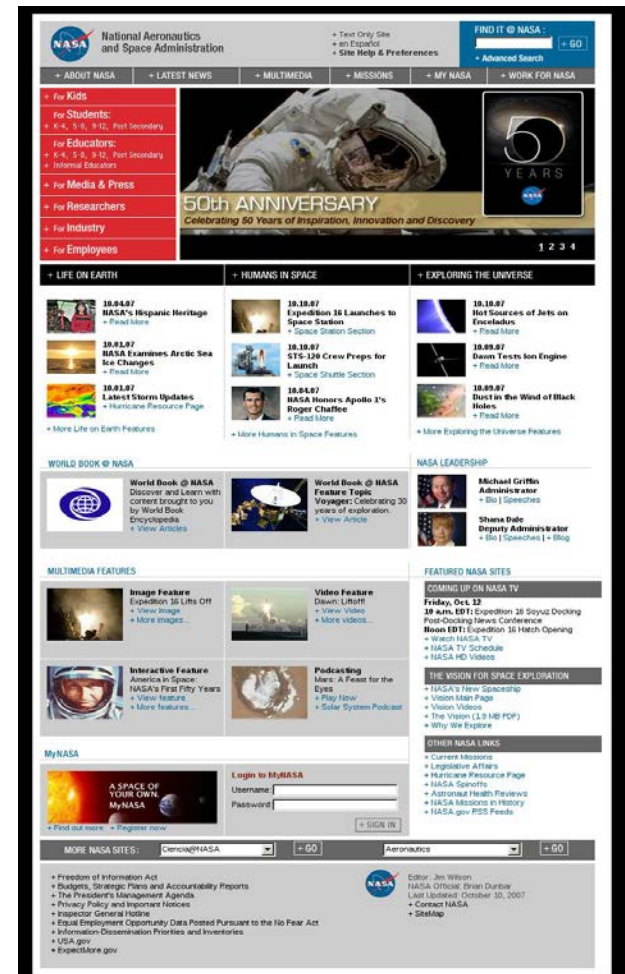


```
var fstr = "";
if(hasFlash(6)) {
    fstr+="
```

# Stepping Back, NASA 2007

<http://web.archive.org/web/20071011055607/http://www.nasa.gov>

- JS check removed
- Content is Accessible
- Archivability Goes Up
- Accessibility → Archivability



# Reinforcing Case: Wikipedia



2001

2002

2003

2004

2005

2006

2007

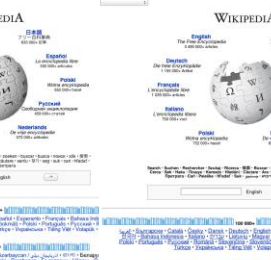
2008

2009

2010

2011

2012



TPDL 2  
Septer

# Summary

- Javascript: good for interaction, bad for archiving
  - crawlers miss URIs at crawl time
  - rendering yesterday's pages causes them to reach into today's web
- Different trends of archivability over time:
  - YouTube: bad-->worse
  - NASA: good-->bad-->good
  - Wikipedia: good
  - see all: [http://www.cs.odu.edu/~mkelly/semester/2013\\_spring/20130127alexatop10/](http://www.cs.odu.edu/~mkelly/semester/2013_spring/20130127alexatop10/)
- Overall, archivability is getting worse
  - 24% increase in missing embedded resources from 2006-2010 due to Javascript