4-11-2016

# Combining Heritrix and PhantomJS for Better Crawling of Pages with Javascript

Justin F. Brunelle
*Old Dominion University*, jbrun008@odu.edu

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

Michael L. Nelson
*Old Dominion University*, mnelson@odu.edu

# Combining Heritrix and PhantomJS for Better Crawling of Pages with Javascript

Justin F. Brunelle
Michele C. Weigle
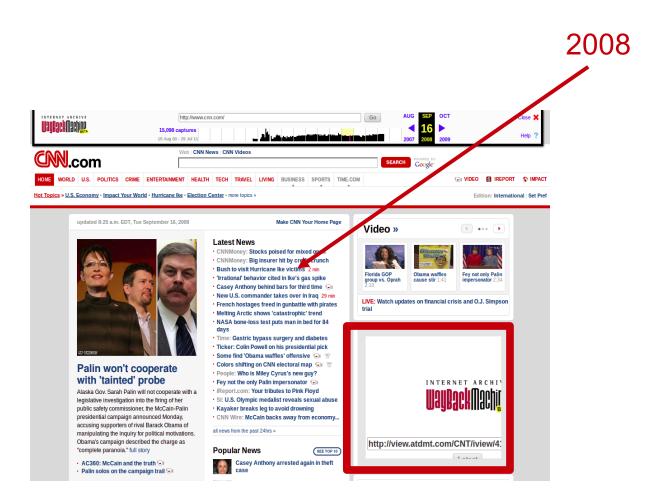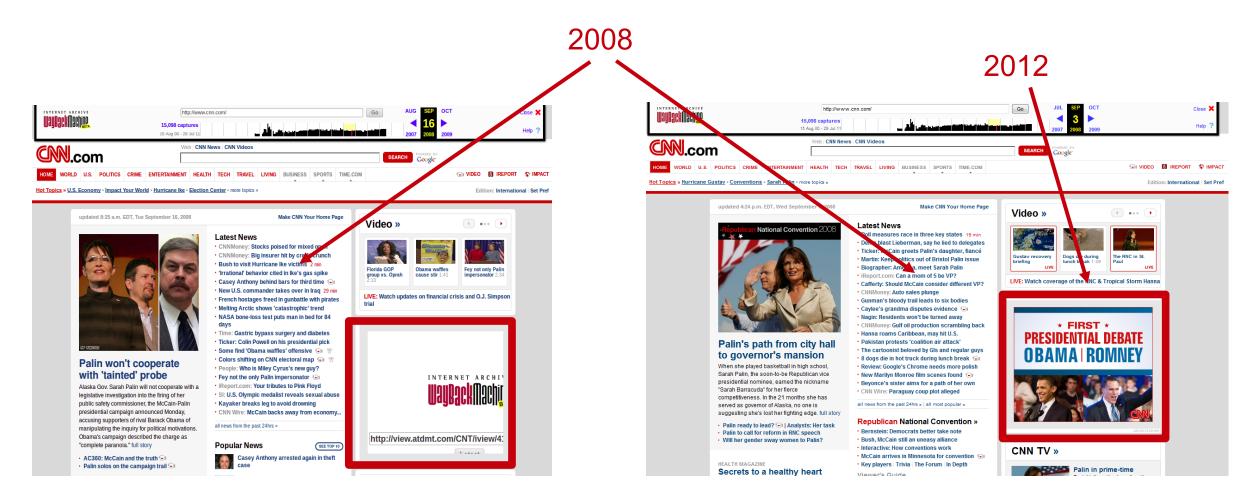**Michael L. Nelson**
Web Science and Digital Libraries Research Group
Old Dominion University
@WebSciDL

IIPC 2016
Reykjavik, Iceland, April 11, 2016

# Javascript can create missing resources (bad)

2008

2

# Javascript can create missing resources (bad) or Temporal violations (worse)

2008

2012

# Old ads are interesting

# New ads are annoying…for now.



"Why are your parents wrestling?"

# Today's ads are missing from the archives

http://adserver.adtechus.com/addyn/3.0/5399.1/2394397/0/-1/QUANTCAST;;size=300x250;target=_blank;alias=p36-17b4f9us2qmzc8bn;kvp36=p36-17b4f9us2qmzc8bn;sub1=p-4UZr_j7rCm_Aj;kvl=172802;kvc=794676;kvs=300x250;kvi=c052a803d0b5476f0bd2f2043ef237e27cd48019;kva=p-4UZr_j7rCm_Aj;rdclick=http://exch.quantserve.com/r?a=p-4UZr_j7rCm_Aj;labels=_qc.clk,_click.adserver.rtb,_click.rand.85854;rtbip=192.184.64.144;rtbdata2=EAQaFUhSQmxvY2tfMjAxNlRheFNlYXNvbiCZiRcogsYKMLTAMDoSaHR0cDovL3d3dy5jbm4uY29tWih UUEhwYlUzM3ZqeFFU5LTA1SGZEMk1SXzE0anBVcGU0d0dxTG10STFUdUs2IECAAb_JicoFoAEBqAGhy7YCugEoVFBIcGJVMzN2anhVOS0wNUhmRDJNUl8xNGpwVXBlNHdHcUxtdEkxVMAB3ed3yAGUp7GUqSraAShjMDUyYTgwM2QwYjU0NzZmMGJkMmYyMDQzZWYyMzdlMjdjZDQ4MDE55QHvEWs-6AFkmAK2wQqoAgWoAgawAgi6AgTAuECQwAICyAIA0ALe9baMj4Cos-oB

# JavaScript is hard to replay

What happens when an event is completely lost?

http://ws-dl.blogspot.com/2013/11/2013-11-28-replaying-sopa-protest.html

# SOPA: Historically significant, archivally difficult



https://en.wikipedia.org/wiki/Stop_Online_Piracy_Act
https://en.wikipedia.org/wiki/Protests_against_SOPA_and_PIPA

http://en.wikipedia.org/wiki/Main_Page     January 18th, 2012

http://web.archive.org/web/20120118110520/http://en.wikipedia.org/wiki/Main_Page  January 18th, 2012

# Problem!



The archives contain the Web as
seen by crawlers

# Why archive?

The Internet Archive has everything!

Why didn't you back it up?

Participating institutions can hand over their databases.

# Crimean Conflict

Russian troops captured the Crimean Center for Investigative Journalism

**Gunman: "We will try to agree on the correct truthful coverage of events."**

# Archive-It to the rescue!

# How well is it archived?

- Masked gunman have your servers
  - anything onsite is gone or altered

Threat models: http://blog.dshr.org/2011/01/threats-to-preservation.html
Automating assessment of crawl quality: http://www.cs.odu.edu/~mln/pubs/jcdl-2014/jcdl-2014-brunelle-damage.pdf

# Any future discussion of the 21$^{st}$ century will involve the web and the web archives

# Any future discussion of the 21$^{st}$ century will involve the web and the web archives

But JavaScript is hard to archive, resulting in archives of content as seen by crawlers rather than as seen by users

# Any future discussion of the 21$^{st}$ century will involve the web and the web archives

But JavaScript is hard to archive, resulting in archives of content as seen by crawlers rather than as seen by users

Goal: Mitigate the impact of JavaScript on the archives by making crawlers behave like users

# W3C Web Architecture



**Dereference** a URI, get a representation

# JavaScript Impact on the Web Architecture



http://maps.google.com

Identifies

Represents

HTTP GET

HTTP Response (200 OK)

Delivered to

JavaScript

# JavaScript makes requests for new resources after the initial page load

http://maps.google.com

Identifies

HTTP GET

HTTP Response (200 OK)

Delivered to

Represents

JavaScript

*Deferred Representation*

Live: JavaScript    PhantomJS: JavaScript    Heritrix: *No JavaScript*

Live Resource

PhantomJS Crawled

Heritrix Crawled, Wayback replayed

9

# JavaScript != Deferred

# Web Browsing Process



Archival Tools stop here

HTTP GET Request for Resource R

HTTP 200 OK Response: R Content

Browser renders and displays R

JavaScript requests embedded resources

Server returns embedded resources

R updates its representation

# Web Browsing Process



HTTP GET Request for Resource R

HTTP 200 OK Response: R content

HTTP GET Request: URI-R

HTTP 200 OK: Representation

Browser renders and displays R

JavaScript in R requests embedded resource X from the server (XMLHttpRequest)

Server returns X as a XMLHttpRequest

*Deferred representations*

R updates its representation

# Web Browsing Process

Archival Tools stop here

HTTP GET Request for Resource R

HTTP 200 OK Response: R Content

Browser renders and displays R

JavaScript requests embedded resources

Server returns embedded resources

R updates its representation

26

# Web Browsing Process



Archival Tools stop here

HTTP GET Request for Resource R

HTTP 200 OK Response: R Content

Browser renders and displays R

JavaScript requests embedded resources
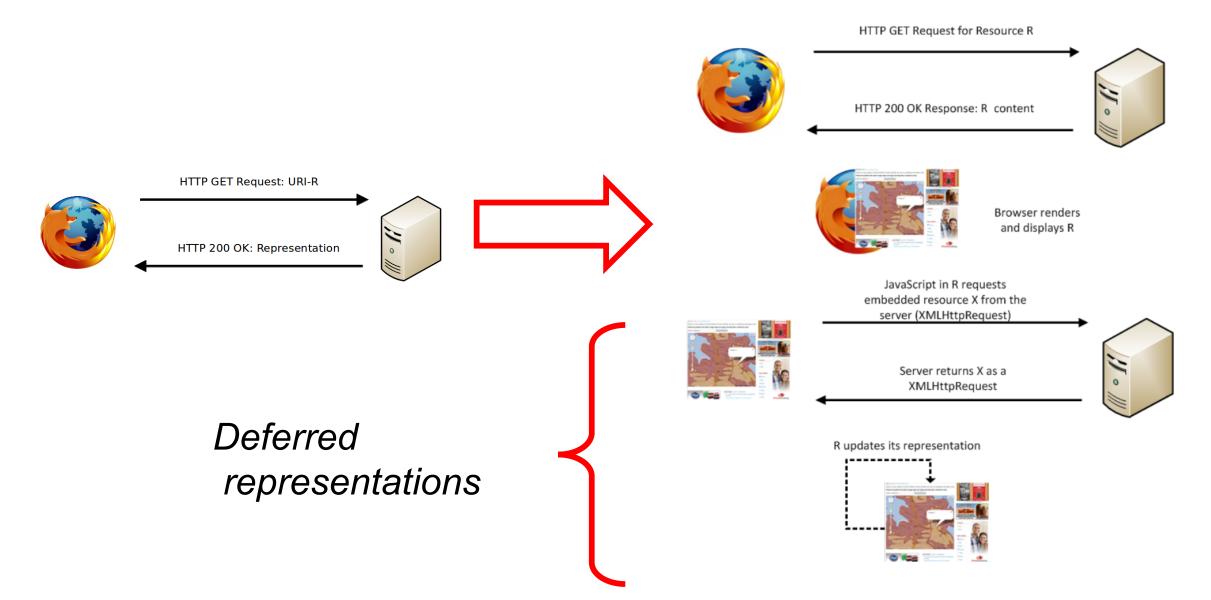
Server returns embedded resources

R updates its representation

Archival approach not defined!

# Current Workflow

- Dereference URI-Rs
- Archive representation
- Extract embedded
- URI-Rs
- Repeat

**Crawl Frontier**

$U_1$

$U_2$

$U_3$

$U_4$

$U_5$

1. Crawl $U_1$

4. Repeat $U_{2...n}$

2. Extract embedded resources and link $U_1$

3. Add $U_4$, $U_5$ to frontier

$U_4$

$U_5$

# Two-Tiered Crawling

**Crawl Frontier**

1. Crawl & archive $U_1$

4. Repeat $U_{2...n}$

2'. Classify representation
- Non-Deferred
- Deferred

3. Extract embedded resources and links from $U_1$

3'. Extract interactive elements in $U_1$

4. Add $U_3$, $U_4$ to frontier

4. Add embedded links and resources from $U_{1.1}$, $U_{1.2}$, $U_{1.3}$ to frontier

XMLHttpRequest

3". Interact with elements & archive $U_{1.1}$, $U_{1.2}$, $U_{1.3}$

**Interaction Frontier of $U_1$**

"**Archiving Deferred Representations Using a Two-Tiered Crawling Approach**", *iPRES2015*

"**Adapting the Hypercube Model to Archive Deferred Representations at Web-Scale**", *Technical Report, arXiv:1601.05142, 2016*
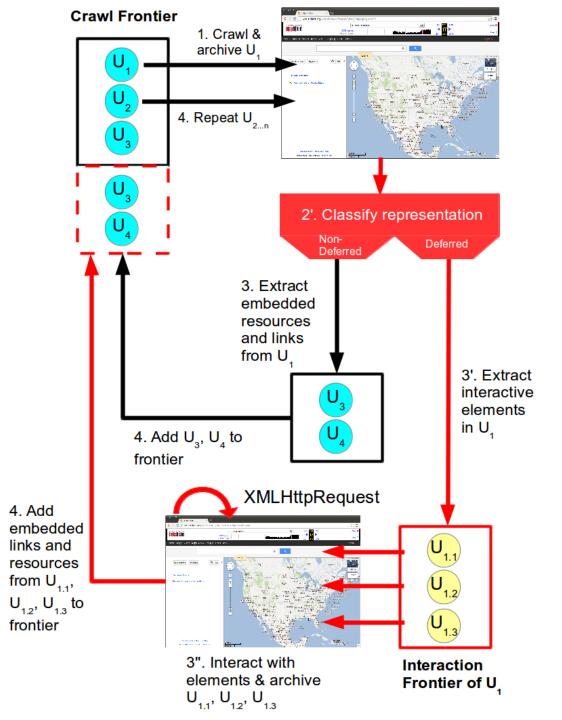
29

**Crawl Frontier**



1. Crawl & archive $U_1$

4. Repeat $U_{2...n}$

2'. Classify representation

Non-Deferred | Deferred

3. Extract embedded resources and links from $U_1$

3'. Extract interactive elements in $U_1$

4. Add $U_3$, $U_4$ to frontier

4. Add embedded links and resources from $U_{1.1}$, $U_{1.2}$, $U_{1.3}$ to frontier

XMLHttpRequest

3". Interact with elements & archive $U_{1.1}$, $U_{1.2}$, $U_{1.3}$
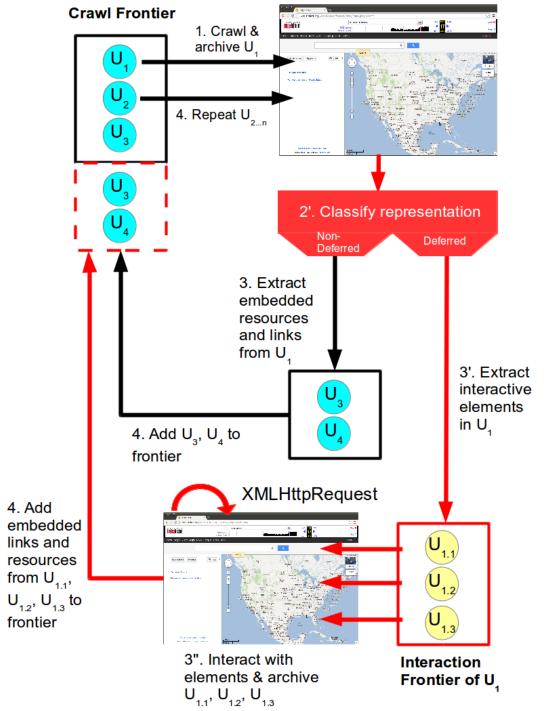
**Interaction Frontier of $U_1$**

rrent workflow not suitable for deferred resentations

<script> tags alone are not indicative of a deferred epresentation. JavaScript can be played back in the rchives!

Two-tiered crawling approach to optimize performance

Use PhantomJS to run JavaScript, interact with the epresentation

30

**Crawl Frontier**

1. Crawl & archive $U_1$

4. Repeat $U_{2...n}$

2'. Classify representation

Non-Deferred | Deferred

3. Extract embedded resources and links from $U_1$

3'. Extract interactive elements in $U_1$

4. Add $U_3$, $U_4$ to frontier

4. Add embedded links and resources from $U_{1.1}$, $U_{1.2}$, $U_{1.3}$ to frontier

XMLHttpRequest

3". Interact with elements & archive $U_{1.1}$, $U_{1.2}$, $U_{1.3}$

**Interaction Frontier of $U_1$**

More URI-Rs in the crawl frontier

Runs more slowly but more deeply

rrent workflow not suitable for deferred presentations

<script> tags alone are not indicative of a deferred representation. JavaScript can be played back in the archives!

Two-tiered crawling approach to optimize performance

Use PhantomJS to run JavaScript, interact with the representation

31

# Classifying deferred representations

- Manually classify 440 URIs (generated from random bitlys) as deferred or non-deferred; build classifier based on 12 different features (8 DOM-based, 4 resource-based)

- On a 10,000 URI set (random bitlys, including 440 from before) compare crawl speed & discovered frontier size with and without classifier
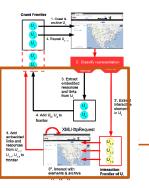
- 

- Data set & code available at:
  - https://github.com/jbrunelle/classifyDeferred/
  - https://github.com/jbrunelle/DataSets

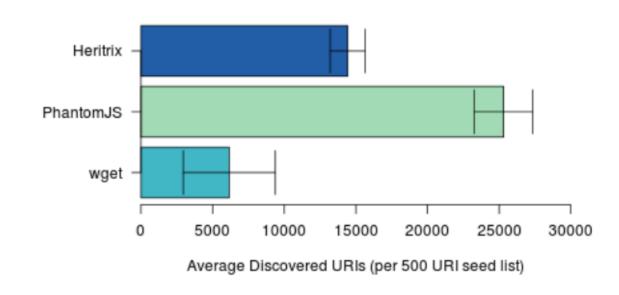# Classifier accuracy improved slightly when monitoring HTTP requests

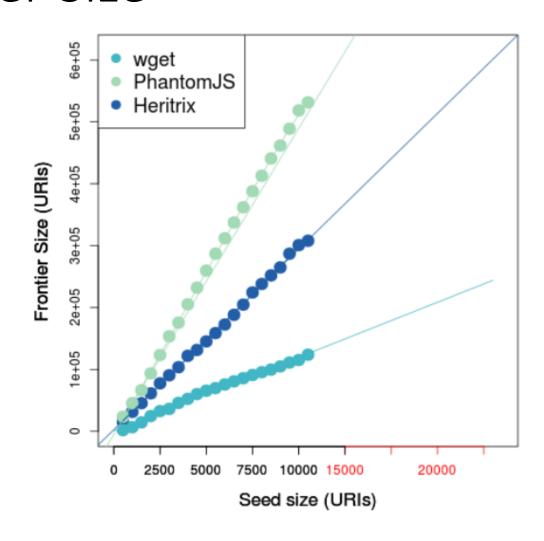| Features | Classification | Accuracy | F-measure | Precision | Recall |
|---|---|---|---|---|---|
| DOM Features Only | Deferred | 79% | 79% | 78% | 81% |
| | Non-deferred | | | 76% | 80% |
| DOM & Resource Features | Deferred | 81% | 82% | 79% | 81% |
| | Non-deferred | | | 90% | 80% |

Table 8: Classification success statistics for DOM-only and DOM and Resource feature sets.

# Performance: Frontier Size



**Average Frontier Size by Tool**

PhantomJS creates a 1.5x larger crawl frontier than Heritrix

# Are all those URIs the same?

| Trim Policy | Original URI-R | Trimmed URI-R |
|---|---|---|
| No Trim | `http://example.com/folder/index.html?param=value` | `http://example.com/folder/index.html?param=value` |
| Origin Trim | `http://example.com/folder/index.html?callback=cs.odu.edu` | `http://example.com/folder/index.html` |
| Base Trim | `http://example.com/folder/index.html?param=value` | `http://example.com/folder/index.html` |
| Session Trim | `http://example.com/folder/index.html?param=value&sessionid=12345` | `http://example.com/folder/index.html?param=value` |
| HTTP Trim | `http://example.com/folder/index.html?param=value&httpParam=http://www.test.com/` | `http://example.com/folder/index.html?param=value` |

**Table 3: Examples of the URI trim policies.**

| Trim Type | URI Duplicates | URI and Entity Duplicates | Accuracy |
|---|---|---|---|
| No Trim | 6,469 | 4,684 | 0.68 |
| Origin Trim | 7,078 | 4,749 | 0.68 |
| Base Trim | 10,359 | 5,191 | 0.56 |
| Session Trim | 8,159 | 4,921 | 0.64 |
| HTTP Trim | 7,315 | 4,868 | 0.67 |

**Table 4: Detected duplicate URIs, entity bodies, and the overlap between the two using the five URI string trimming policies.**
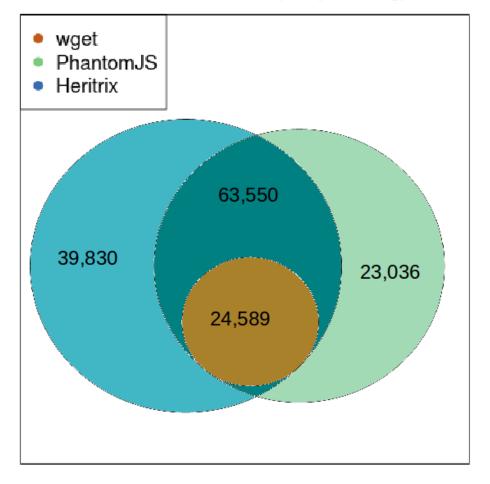
TP = URIs match & entities match
TN = neither URI nor entity matches
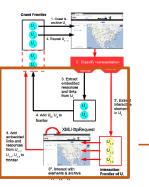P + N = 19,522

# Trimming shrinks the PhantomJS Frontier



**Unions and Intersections (String Matching)**

- wget
- PhantomJS
- Heritrix

39,830
63,550
194,818
24,589

**Unions and Intersections (Fuzzy Matching)**

- wget
- PhantomJS
- Heritrix
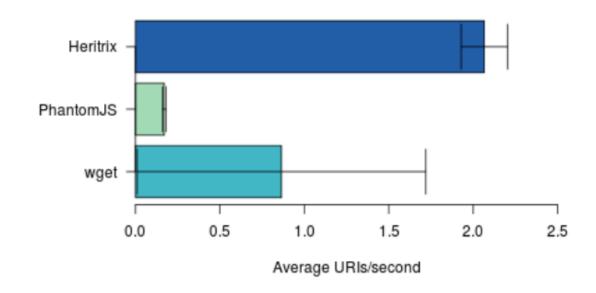
63,550
39,830
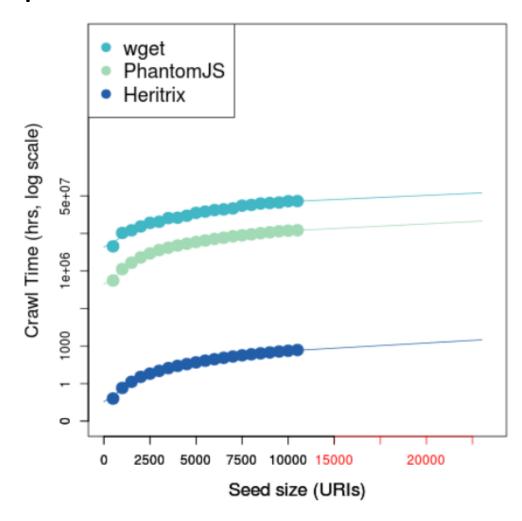23,036
24,589

(Base policy shown)

# Performance: Crawl Speed



Average Crawl Rate by Tool



Heritrix:        ~2 URIs/second

PhantomJS: ~4 seconds/URI
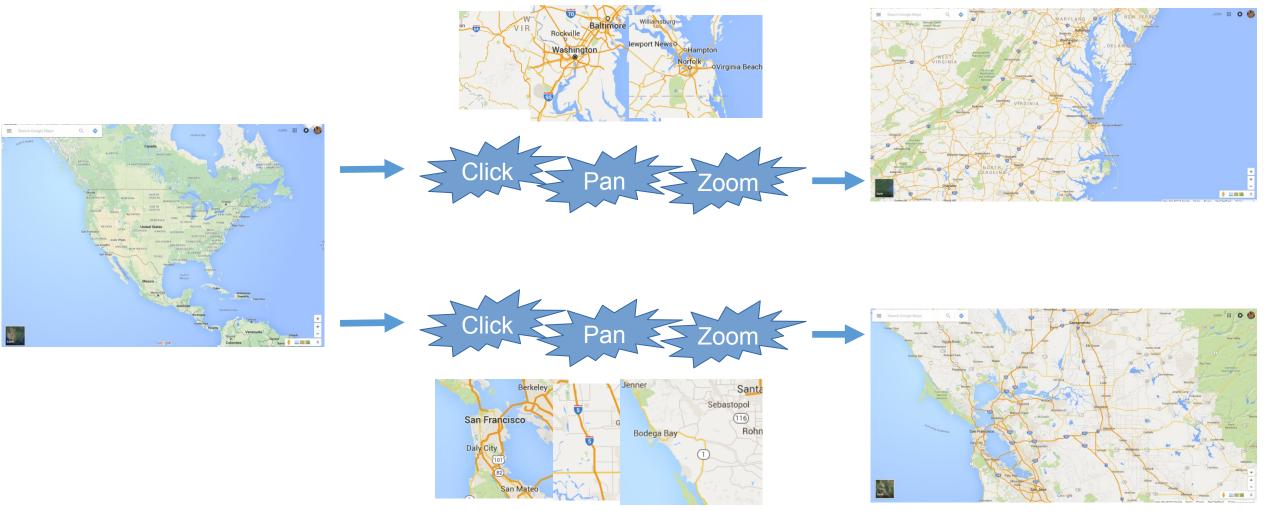
# How long would it take to crawl *everything*?

| Crawl Strategy | Crawl Time (hrs) | Crawl Rate ($t_{URI}$) | Frontier Size ($|F|$) |
|---|---|---|---|
| wget | 416.16 | 0.864 | 129,443 |
| Heritrix | 407.53 | 2.065 | 302,961 |
| PhantomJS | 8,684.38 | 0.170 | 531,484 |
| Heritrix + PhantomJS | 9,100.54 | 0.152 | 537,609 |
| Heritrix + PhantomJS with Classifier | 6,495.23 | 0.196 | 458,815 |

nearly a year!

(obviously parallelization would help)

Table 9: A summary of *extrapolated* performance (based on our calculations) of single- and two-tiered crawling approaches.

# Descendants = States of deferred representations reached through client-side events
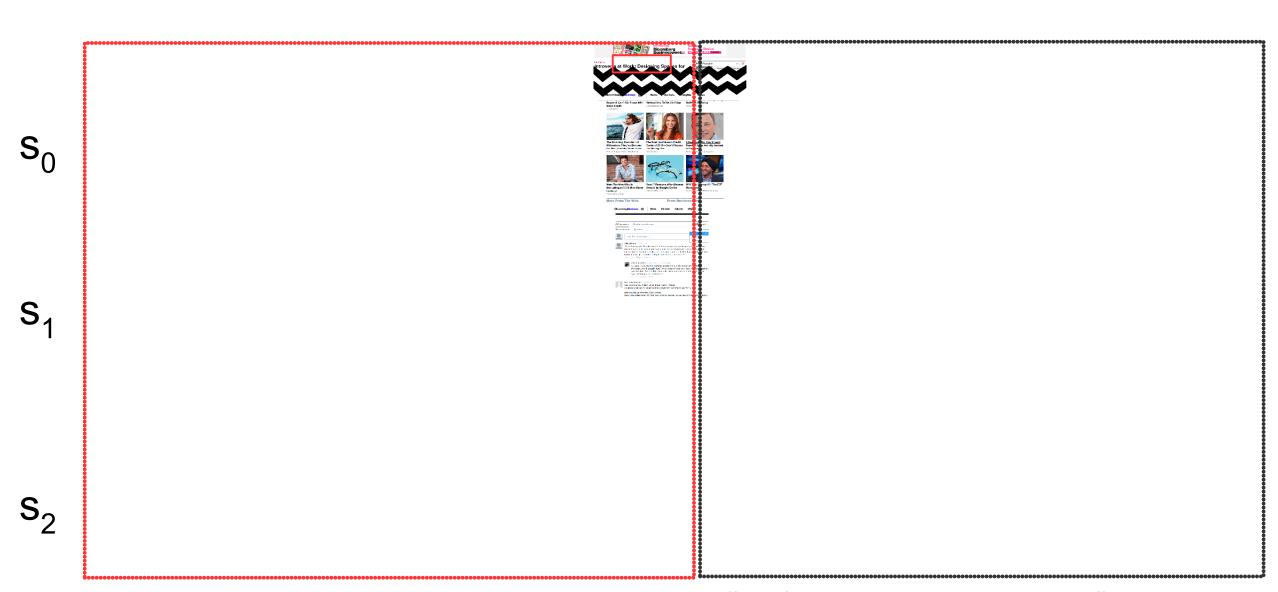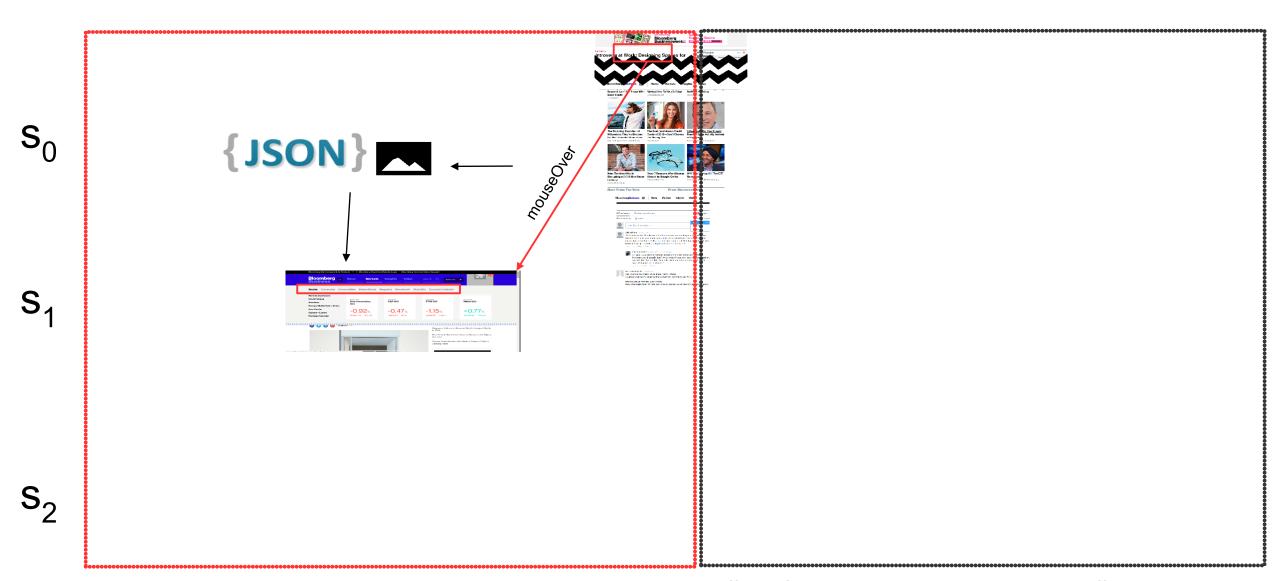


Click → Pan → Zoom
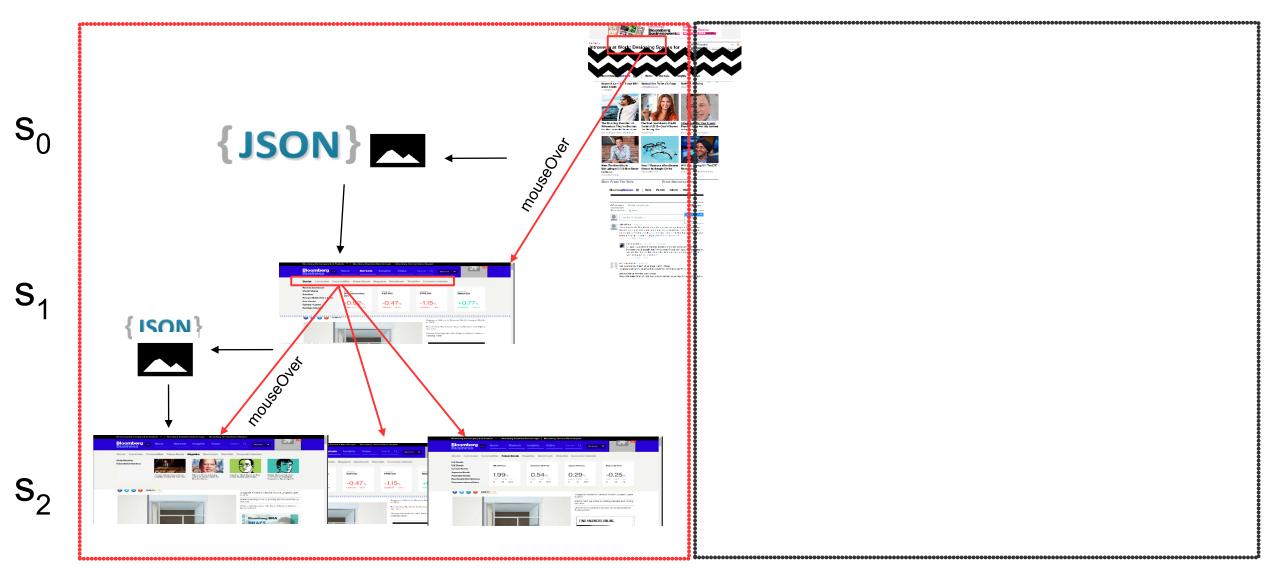
# Finding descendants

- Return to the same 440 URIs from before

-

- Use VisualEvent to identify interactive elements

- http://ws-dl.blogspot.com/2015/06/2015-06-26-phantomjsvisualevent-or.html

-

- Adapting work on state equivalency based on DOM equivalency, we define state equivalency as requiring the same embedded resources

- Report & code:

- http://arxiv.org/abs/1601.05142

- https://github.com/jbrunelle/clientSideState

-

-

# Interaction Trees are 2 Levels Deep

$s_0$

$s_1$

$s_2$

http://www.bloomberg.com/bw/articles/2014-06-16/open-plan-offices-for-people-who-hate-open-plan-offices

# Interaction Trees are 2 Levels Deep

$s_0$

$s_1$

$s_2$

mouseOver

# Interaction Trees are 2 Levels Deep



S₀

S₁

S₂

mouseOver

mouseOver

{JSON}

{JSON}

# Interaction Trees are 2 Levels Deep



$s_0$

$s_1$

$s_2$

mouseOver

mouseOver

# Interaction Trees are 2 Levels Deep



$s_0$

$s_1$

$s_2$

{JSON}

mouseOver

click

# Expanding the Crawl Frontier

Nondeferred

Deferred

3,640  $S_0$  6,983

$S_1$  53,706

$S_2$ 10,208

**Level $s_1$ provides the greatest benefit to the crawl frontier**

# Crawling Descendants



**New embedded resources at levels $s_1$ are largely unarchived**

# Expanding the crawl frontier

| Event Type | Percent of URI-Rs | | Contribution to $RP_{new}$ |
|---|---|---|---|
| | Deferred | Nondeferred | |
| click | 62.11% | 4.29% | 63.2% |
| mouseover | 25.26% | 3.00% | 4.7% |
| mousedown | 16.84% | 1.72% | 2.8% |
| blur | 14.74% | 0.86% | 9.8% |
| change | 11.58% | 2.14% | 0.0% |

**Click events lead to the most descendants**

# Future Work

- Modeling user interactions, tendencies, and simulation
  - form filling
  - click & navigation likelihood
  - Added frontier 92% unarchived
- Archival Halting Problem: How much is enough?

  - Mapping Applications – How many pans and zooms gets all the Norfolk, VA Google map tiles?
  - How many CNN.com pages get all the Google Ads?
  - Game walkthrough metaphor?  (insert url here)
- Playing back WARCs with IIPC metadata of deferred representations and descendants

# Contributions

- Defined:
  - *deferred representations*: representations that need client-side processing to load all required embedded resources
  - *descendants*: representation states reachable only via client-side events
- Two-tiered crawling of deferred representations
  - 10.5 times slower
  - 1.5 times larger frontier
  - 2 levels of descendants
- 2 levels are sufficient for descendants
  - Added frontier 92% unarchived
- More info:
  - http://arxiv.org/abs/1508.02315
  - http://arxiv.org/abs/1601.05142