# How Much of the Web is Archived?*

**Scott Ainsworth, Ahmed AlSum, Hany SalahEldeen,
Michele C. Weigle, Michael L. Nelson**

**Old Dominion University, USA**
{sainswor, aalsum, hany, mweigle, mln}@cs.odu.edu

**JCDL 2011, Ottawa, Canada**

# How Much of the Web is Archived?*

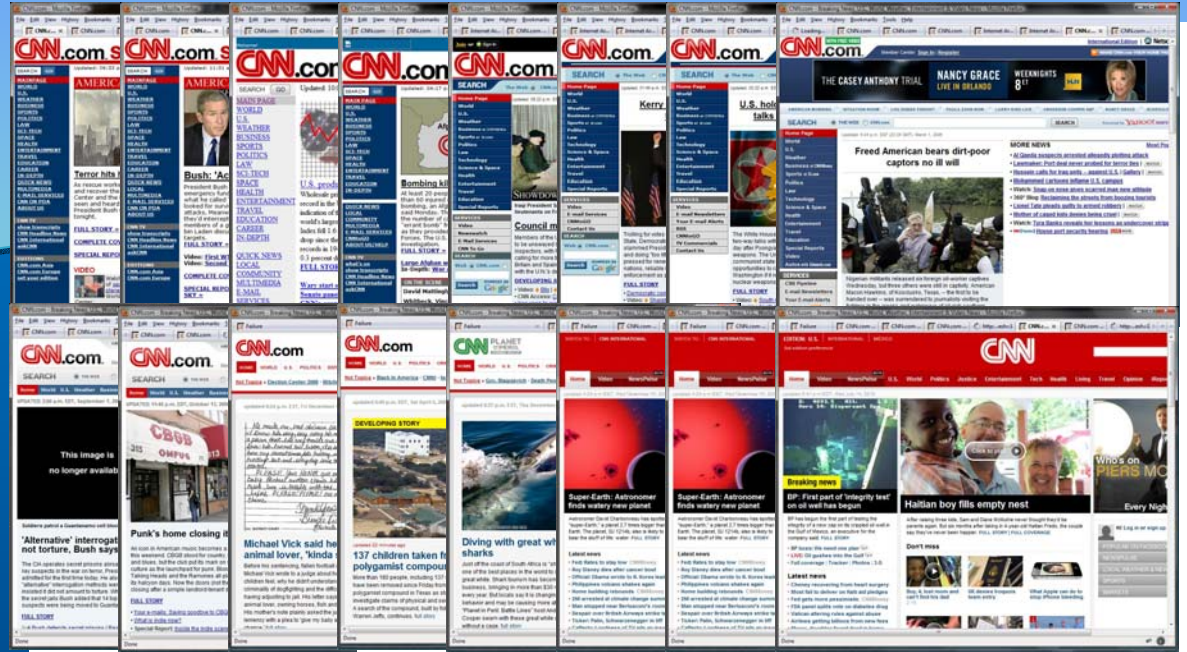**Scott Ainsworth, Ahmed AlSum, Hany SalahEldeen,
Michele C. Weigle, Michael L. Nelson**

**Old Dominion University, USA**
{sainswor, aalsum, hany, mweigle, mln}@cs.odu.edu

JCDL 2011, Ottawa, Canada

# How Many Archive Copies?

**http://www.cnn.com**

**http://cs.odu.edu/~aalsum**

# Research Question

## How Much of the Web is Archived?

# Experiment
## Sampling Techniques

- 4 Sample sets – 1000 URIs each

# Sampling Techniques
## DMOZ

**dmoz** open directory project

- **All the available DMOZ history:**

  **100 snapshots made from July 20, 2000 through October 3, 2010**

- **Remove duplicates, non-HTTP, invalid URIS**

- **Randomly selected 1000 URIs**

| HTTP Status | # |
|---|---|
| 200 | 507 |
| 3xx→200 | 192 |
| 3xx→ Others | 50 |
| 4xx | 135 |
| 5xx | 4 |
| Timeout | 112 |

| Indexed by: | bing | YAHOO! | Google API | Google API + Web |
|---|---|---|---|---|
| | 49% | 41% | 30% | 54% |

# Sampling Techniques
## Delicious

delicious

- **We used Delicious Recent random URI generator** (http://www.delicious.com/recent/?random=1) **on Nov. 22, 2010 to get 1000 URIs.**

| HTTP Status | # |
|---|---|
| 200 | 958 |
| 3xx→200 | 27 |
| 3xx→ Others | 1 |
| 4xx | 8 |
| 5xx | 3 |
| Timeout | 3 |

| Indexed by: | bing | YAHOO! | Google API | Google API + Web |
|---|---|---|---|---|
| | 95% | 86% | 88% | 95% |

# Sampling Techniques
**Bitly**

**bitly**

- **Bitly URI consists of a 1–6 character, alphanumeric hash value appended to http://bit.ly/**
- **Random hash values were created and dereferenced until 1,000 target URIs were discovered**

| HTTP Status | # |
|---|---|
| 200 | 488 |
| 3xx→200 | 243 |
| 3xx→ Others | 36 |
| 4xx | 197 |
| 5xx | 6 |
| Timeout | 30 |

| Indexed by: | bing | YAHOO! | Google API | Google API + Web |
|---|---|---|---|---|
| | 21% | 22% | 24% | 30% |

# Sampling Techniques
## Search Engine



- **Random Sample from Search Engine technique** [Bar-Yossef 2008]
- **The phrase pool was selected from the 5-grams in Google's N-gram data.**
- **A random sample of these queries was used to obtain URIs, of which 1,000 were selected at random.**

| HTTP Status | # |
|---|---|
| 200 | 943 |
| 3xx→200 | 17 |
| 3xx→ Others | 3 |
| 4xx | 16 |
| 5xx | 0 |
| Timeout | 21 |

| Indexed by: | bing | YAHOO! | Google API | Google API + Web |
|---|---|---|---|---|
| | 55% | 98% | 70% | 73% |

# Experiment

- **For each sample set, we used Memento Aggregator to get all the possible archived copies (Mementos).**
- **For each URI, Memento Aggregator responded with TimeMap for this URI.**

Example

  <http://memento.waybackmachine.org/memento/20010819194233/http://jcdl2002.org>;rel="first memento";datetime="Sun, 19 Aug 2001 19:42:33 GMT",

  <http://memento.waybackmachine.org/memento/20011216220248/http://jcdl2002.org>; rel="memento"; datetime="Sun, 16 Dec 2001 22:02:48 GMT",

# Results

| #Mementos | dmoz | delicious | bitly | SE YAHOO! |
|---|---|---|---|---|
| 0 (Not archived) | 10% | 3% | 65% | 22% |
| 1 | 4% | 8% | 10% | 34% |
| 2–5 | 14% | 48% | 17% | 32% |
| 6–10 | 9% | 4% | 2% | 4% |
| More than 10 | 63% | 37% | 6% | 8% |

# Archive categories

We have 3 categories of archives

- Internet Archive (classic interface)

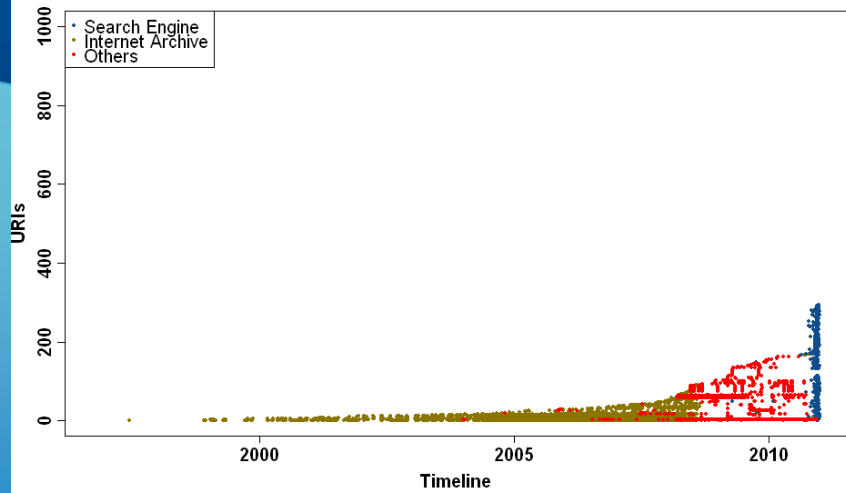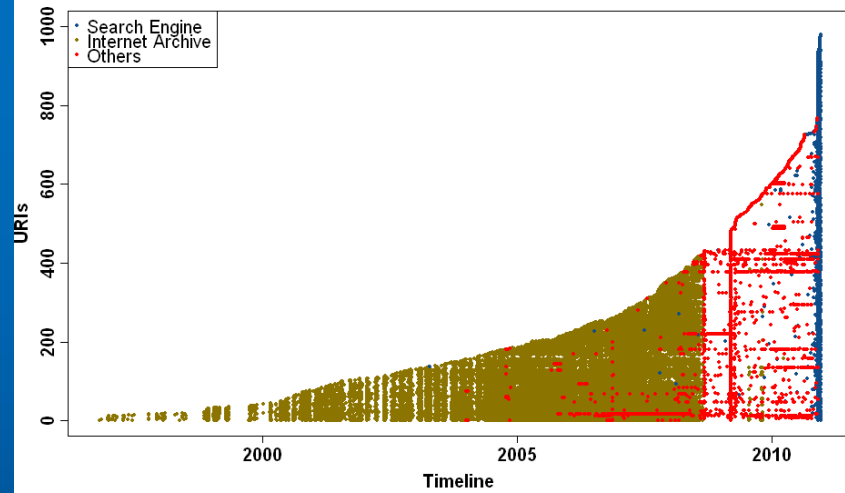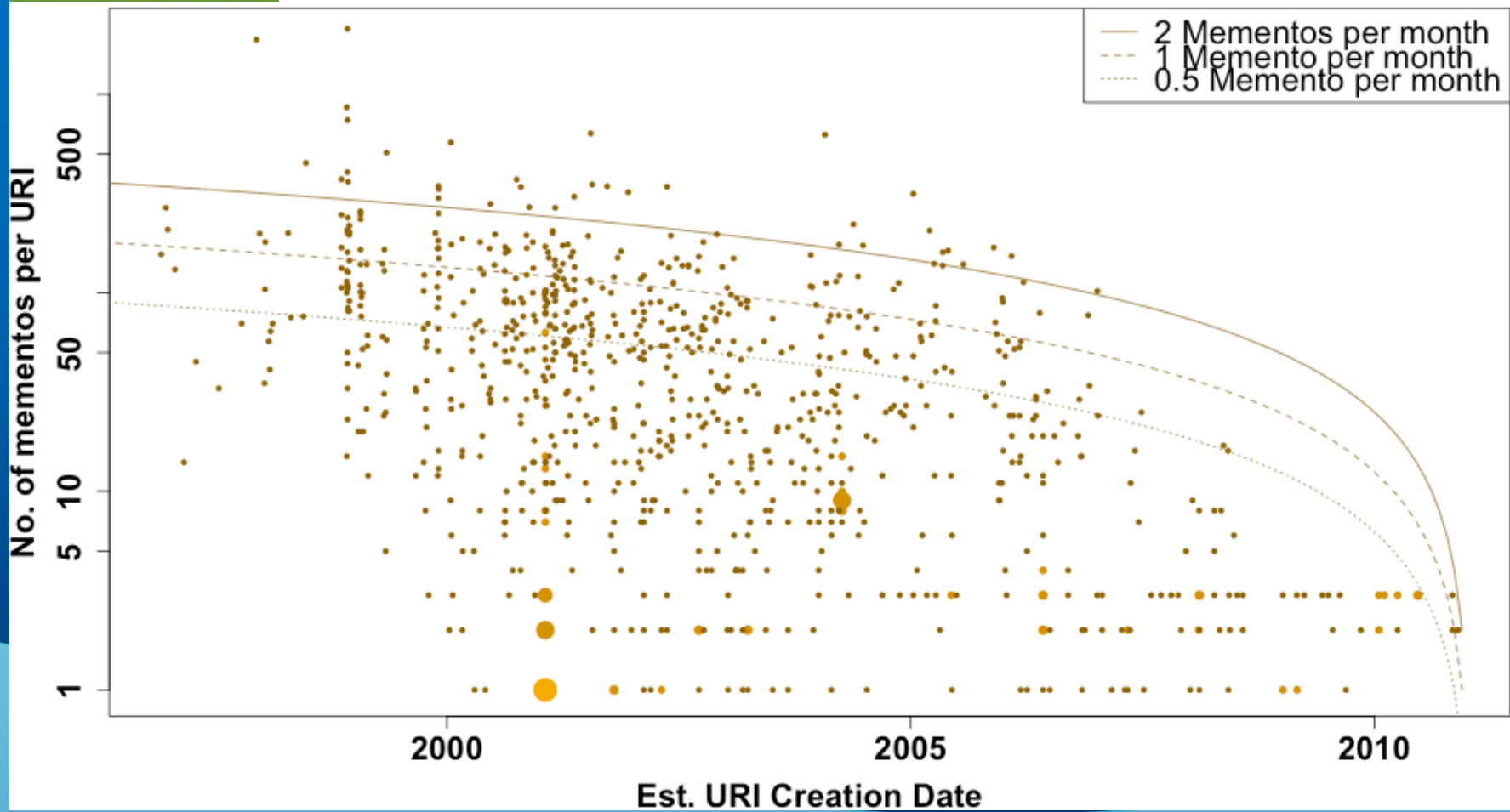- Search engine

- Other archives

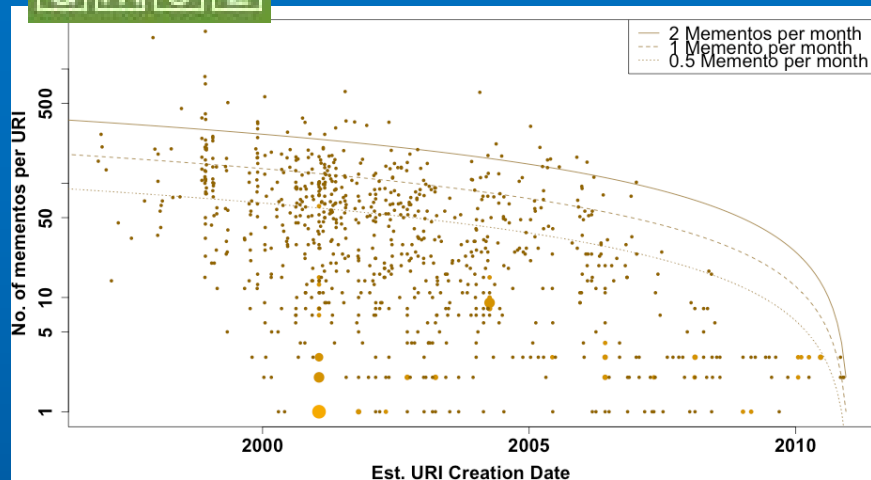**Memento Distribution, ordered by the first observation date.**

Memento Distribution, ordered by the first observation date.

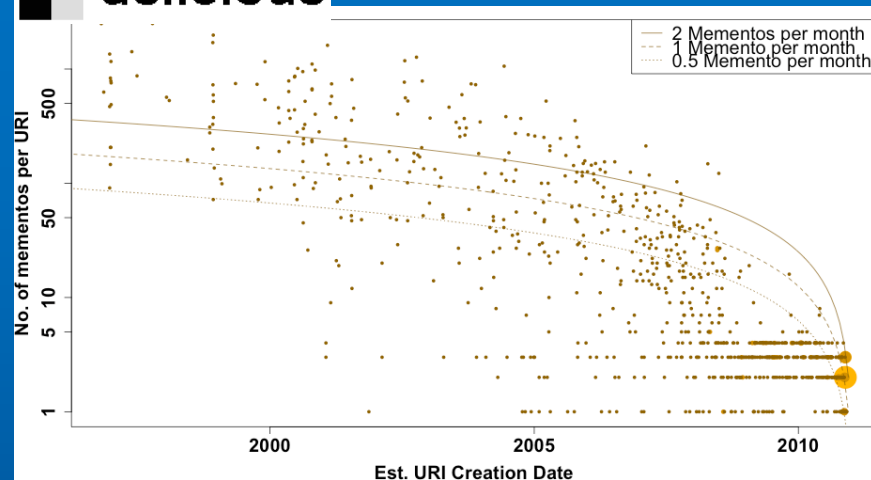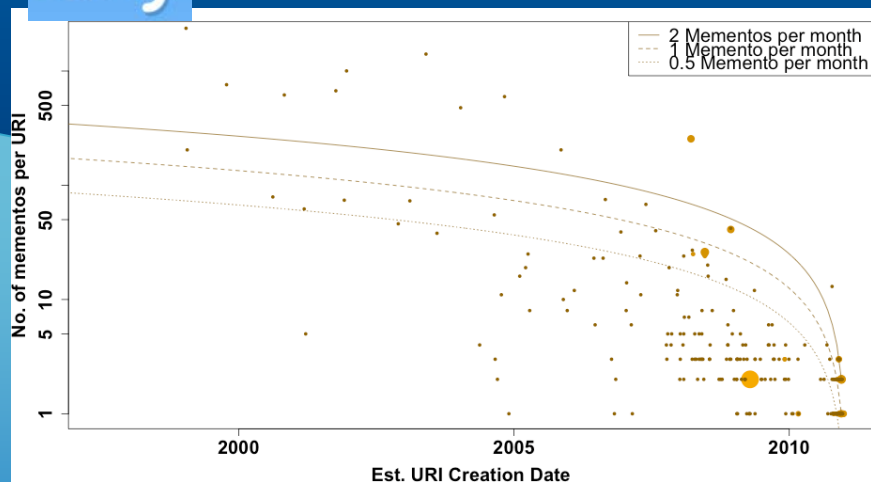(14.6% archived >= once per month)

How often are URIs archived?

(14.6% archived >= once per month)

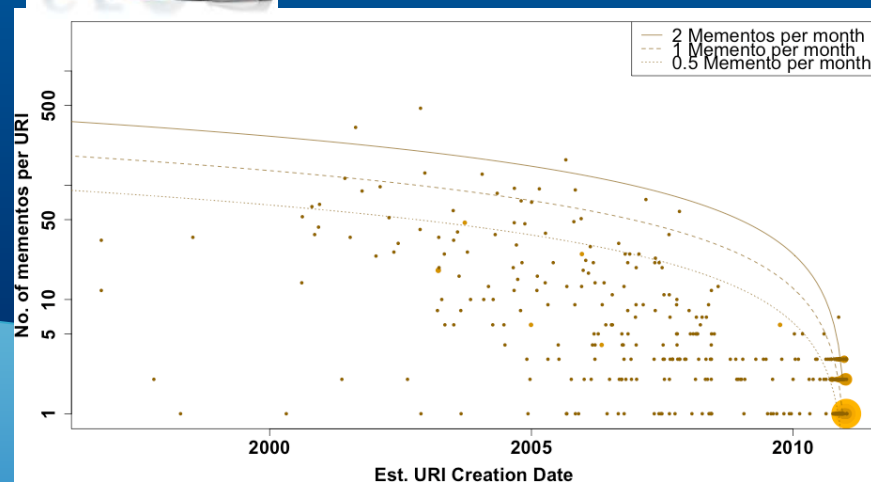(31% archived >= once per month)

(14.7% archived >= once per month)

(28.7% archived >= once per month)

**How often are URIs archived?**

# Analysis

- ## What affects the archival rate?
  - – URI source (More popular – More archival)
  - – BackLinks (Weak positive relationship)
- ## Archive Coverage
  - – Internet Archive (Best)
  - – Search Engine (Good), but only one copy.
  - – The other archives cover only a small fraction with recent copies.

# Conclusion
**How much of the Web is Archived?**

- Tell me what is your URI source!!

| | Including SE cache | Excluding SE Cache |
|---|---|---|
| dmoz | 90% | 79% |
| delicious | 97% | 68% |
| bitly | 35% | 16% |
| SE YAHOO! | 88% | 19% |

**Contact:** Ahmed AlSum (aalsum@cs.odu.edu)

# References

- Ziv Bar-Yossef and Maxim Gurevich. Random sampling from a search engine's index. *J. ACM*, 55(5), 2008.

- Frank McCown and Michael L. Nelson. Characterization of search engine caches. In *Proceedings of IS&T Archiving 2007*, May 2007.

- Herbert Van de Sompel, Michael Nelson, and Robert Sanderson. HTTP framework for time-based access to resource states — Memento, November 2010. http://datatracker.ietf.org/doc/draft-vandesompelmemento/.

# Thank You

**Contact:**

Ahmed AlSum
aalsum@cs.odu.edu