

4-23-2013

## More Archives, More Better

Michael L. Nelson

*Old Dominion University*, [mnelson@odu.edu](mailto: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

Nelson, Michael L., "More Archives, More Better" (2013). *Computer Science Presentations*. 16.  
[https://digitalcommons.odu.edu/computerscience\\_presentations/16](https://digitalcommons.odu.edu/computerscience_presentations/16)

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).

# More Archives, More Better

Michael L. Nelson  
Old Dominion University  
[ws-dl.blogspot.com](http://ws-dl.blogspot.com)

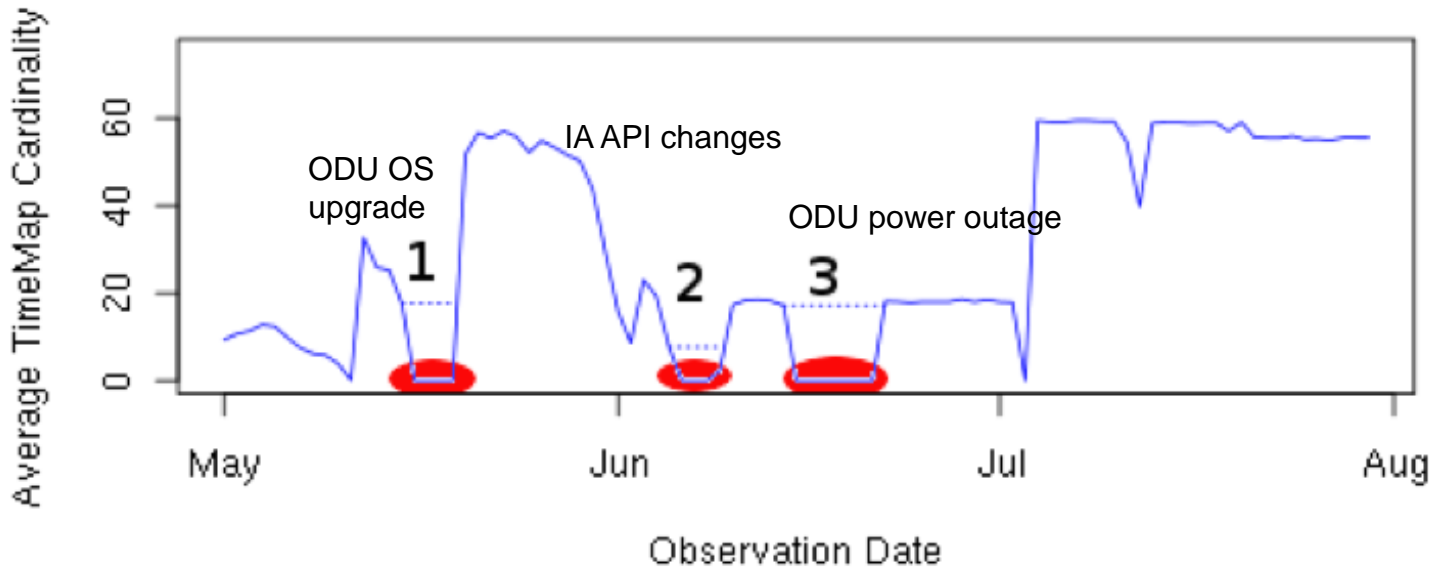
IIPC General Assembly  
Ljubljana, Slovenia  
April 23, 2013

# Three Easy Pieces

- "An Evaluation of Caching Policies for Memento TimeMaps"
  - 4000 aggregated TimeMaps downloaded daily for 3 months
  - 20% of the time the TimeMaps *shrink*
- "How Much of The Web Is Archived?"
  - 4000 URIs, 9 archives, 3 search engines
  - 16% -- 79% of the web archived
- "Profiling Web Archive Coverage for Top-Level Domain and Content Language"
  - 153329 URIs, 12 archives
  - querying only top 3 archives gives a complete TimeMap 84% of the time (52% of the time even if you exclude the IA)

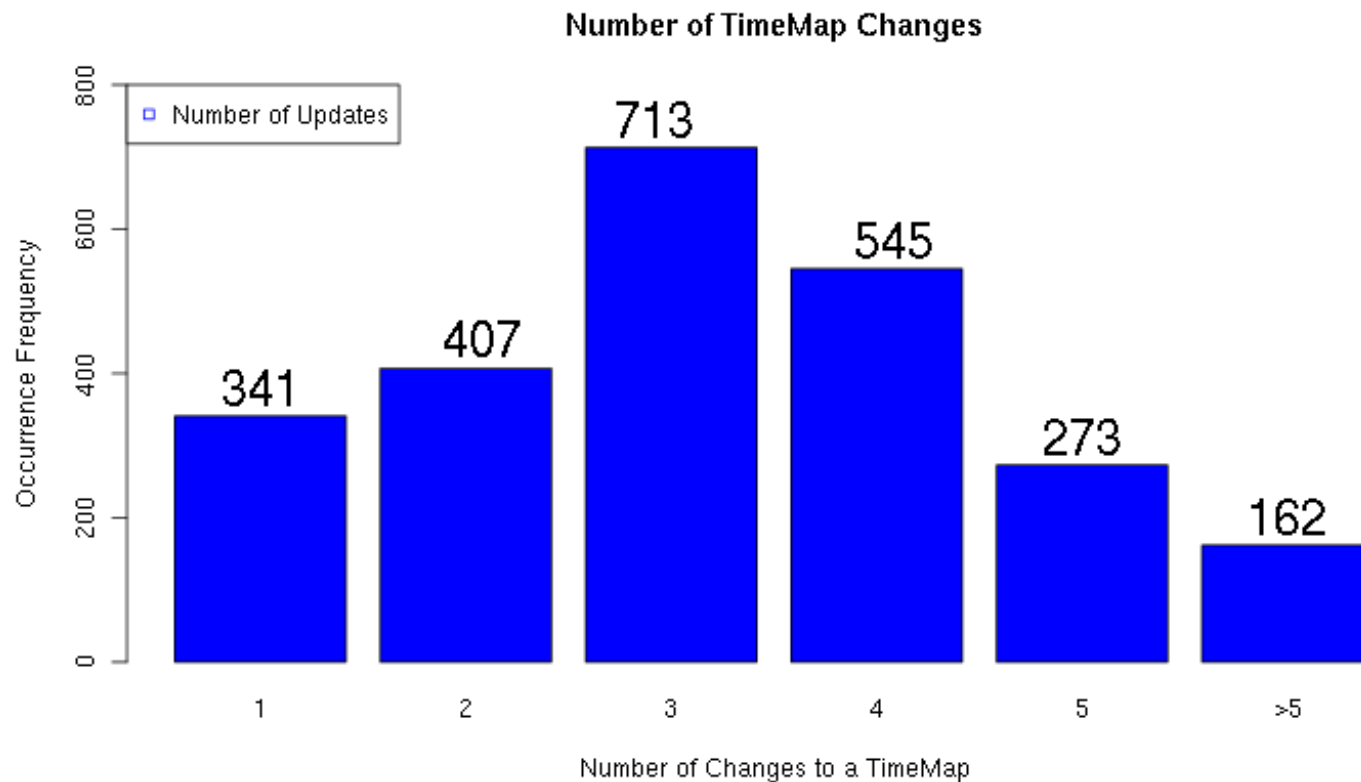
An Evaluation of Caching Policies  
for Memento TimeMaps  
JCDL 2013  
Justin Brunelle, Michael L. Nelson

# Mean # Mementos per TimeMap per Day

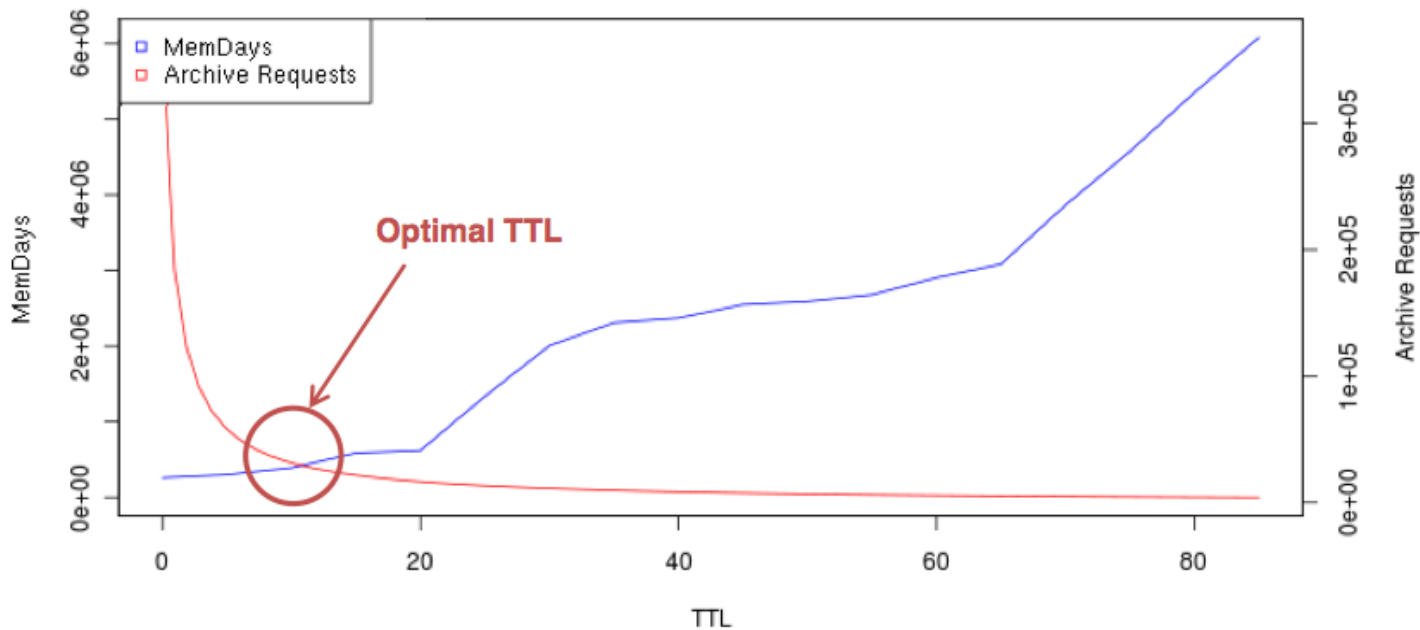


download the same 4000 TimeMaps everyday

# Frequency of TimeMap changes over 92 days



# Optimal TimeMap Cache TTL=15 days



minimizes queries to archives, minimizes "lost" mementos\*days,  
will only cache new TimeMap if it is "bigger"  
question: can we do this adaptively?

How Much of The Web Is Archived?

JCDL 2011

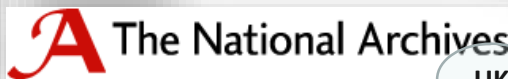
Scott Ainsworth, Ahmed AlSum, Hany SalahEldeen,  
Michele C. Weigle, Michael L. Nelson



# Public Archives, ca. late 2010 / early 2011

## Three categories of archives

- **Internet Archive** (classic interface)
- **Search engine**
- **Other archives**



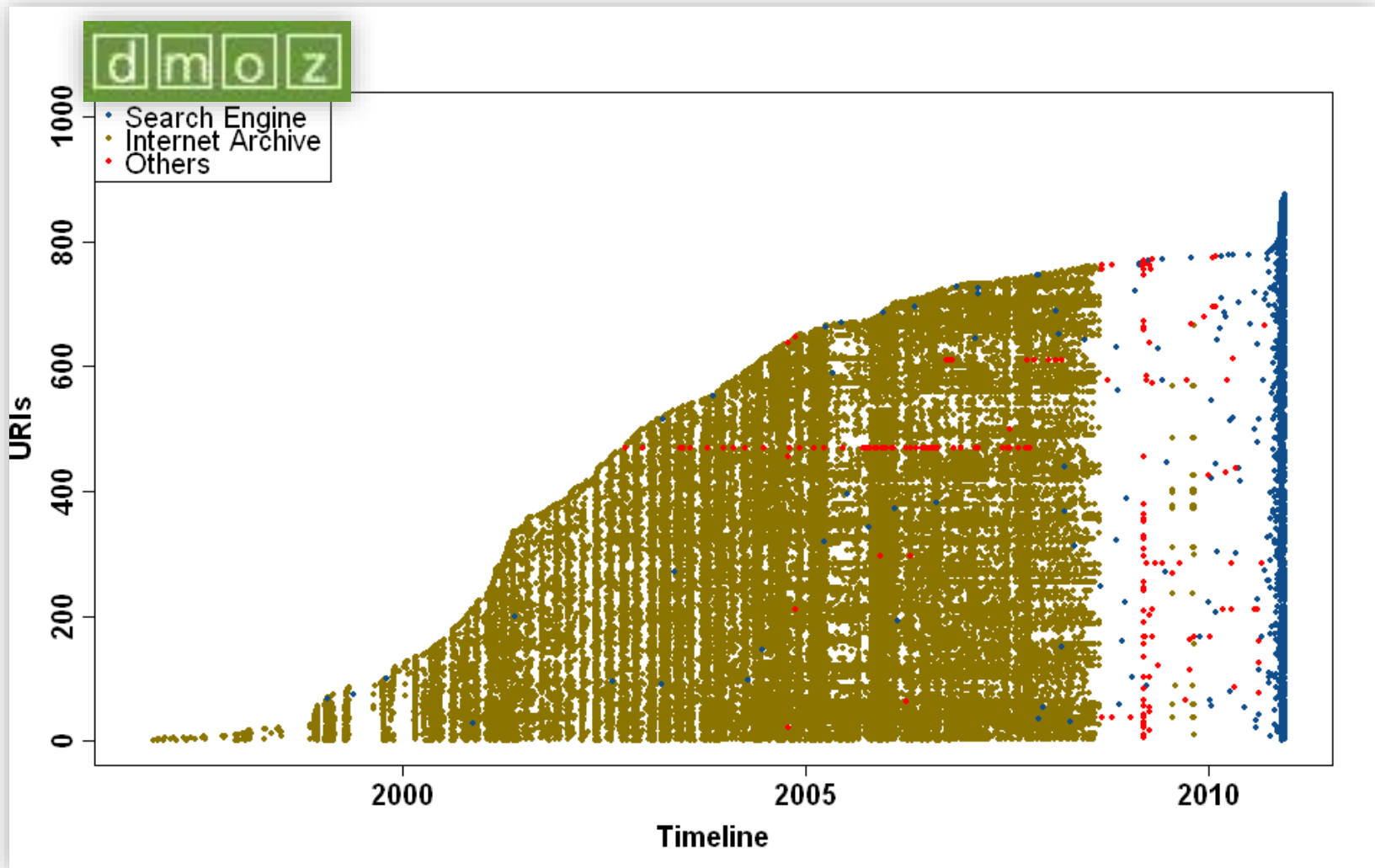
UK



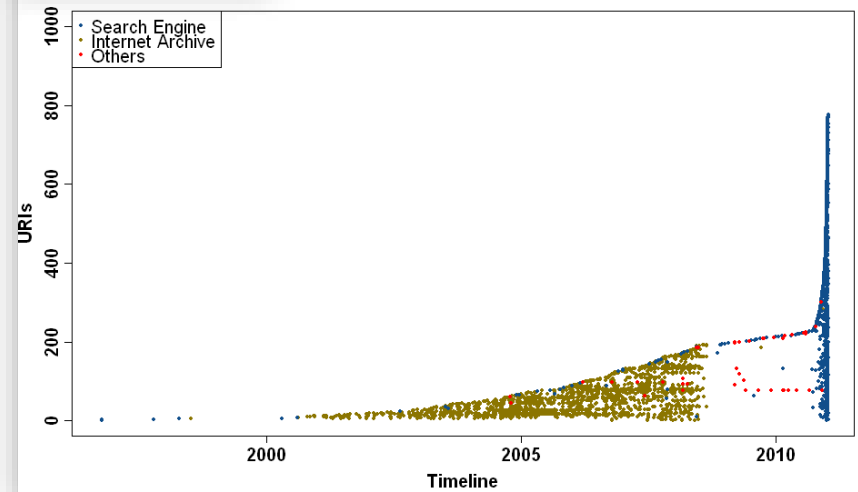
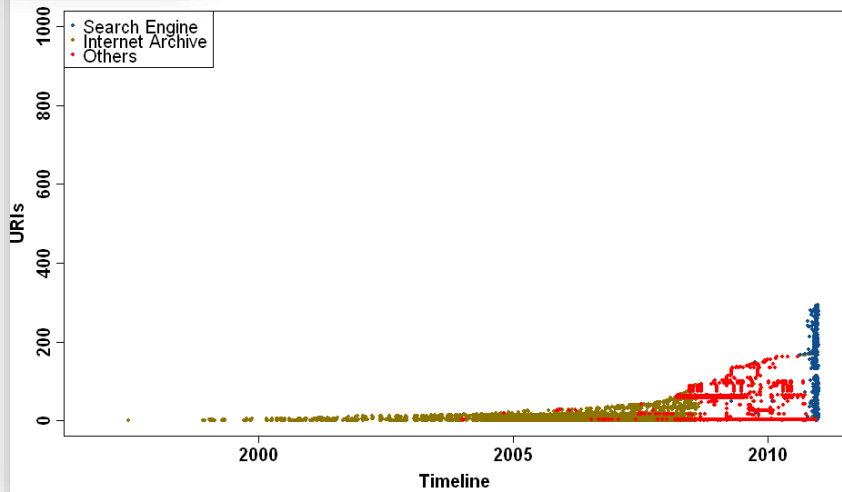
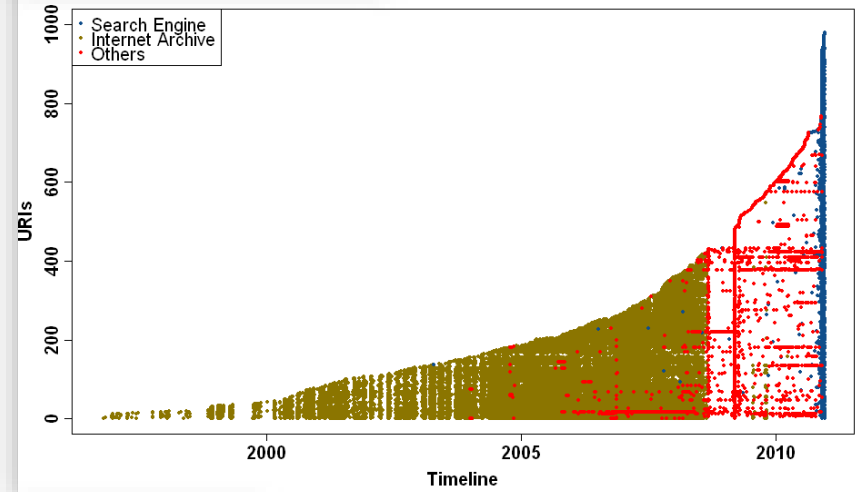
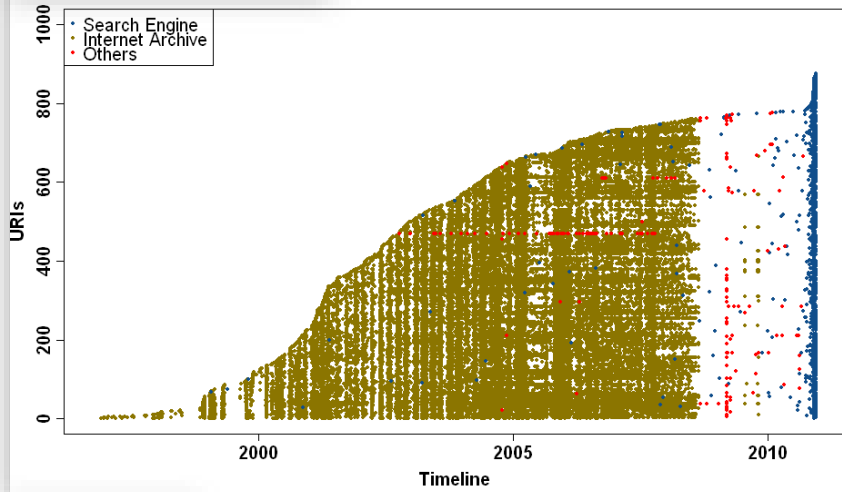
US



# 1000 URIs, ordered by first observation date







See also: <http://ws-dl.blogspot.com/2011/06/2011-06-23-how-much-of-web-is-archived.html>



see also: <http://ws-dl.blogspot.com/2013/04/2013-04-19-carbon-dating-web.html>

# How Much of the Web is Archived?

It depends on which web...

	Including SE cache	Excluding SE Cache
	90%	79%
	97%	68%
	35%	16%
	88%	19%

Changes since 2011: no more free SE APIs; greatly reduced IA quarantine period

Profiling Web Archive Coverage for  
Top-Level Domain and Content Language  
(submitted for publication)

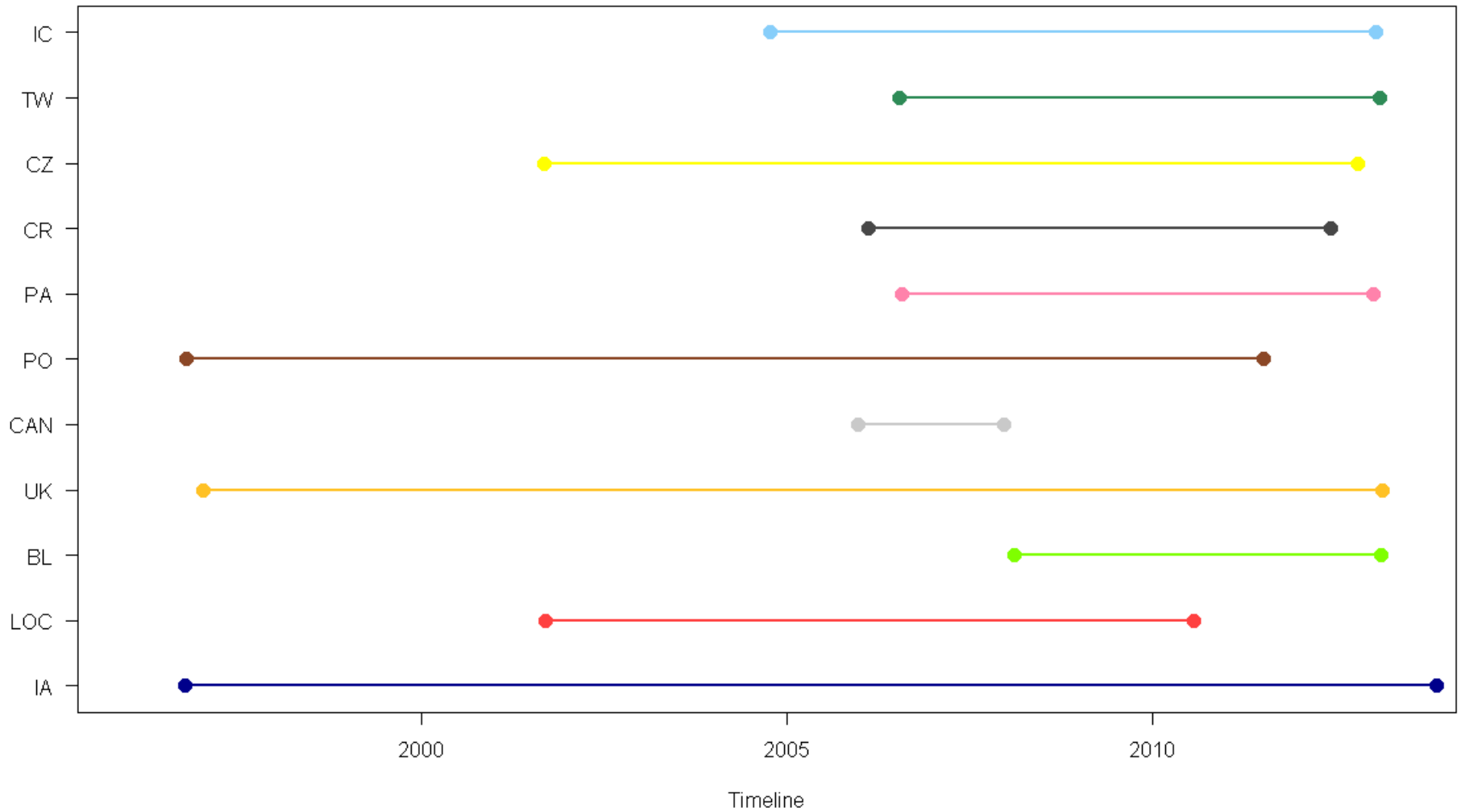
Ahmed AlSum, Michele C. Weigle, Michael L. Nelson,  
Herbert Van de Sompel

## 12 (IIPC) Archives

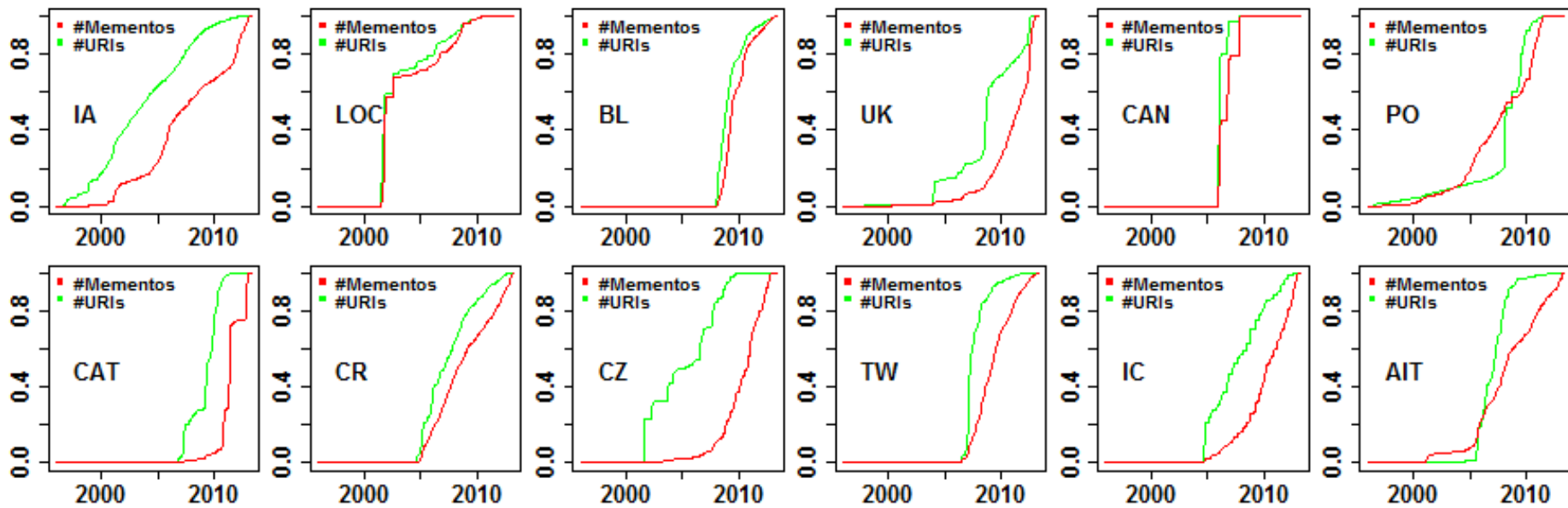
<b>Archive Name</b>	<b>FullText search</b>	<b>Website</b>
IA Internet Archive		<a href="http://web.archive.org">web.archive.org</a>
LoC Library of Congress		<a href="http://www.loc.gov/lcwa">www.loc.gov/lcwa</a>
IC Icelandic Web Archive		<a href="http://vefsafn.is">vefsafn.is</a>
CAN Library & Archives Canada	x	<a href="http://www.collectionscanada.gc.ca">www.collectionscanada.gc.ca</a>
BL British Library	x	<a href="http://www.webarchive.org.uk/ukwa">www.webarchive.org.uk/ukwa</a>
UK UK Gov. Web Archive	x	<a href="http://webarchive.nationalarchives.gov.uk">webarchive.nationalarchives.gov.uk</a>
PO Portuguese Web Archive	x	<a href="http://arquivo.pt">arquivo.pt</a>
CAT Web Archive of Catalonia	x	<a href="http://www.padi.cat">www.padi.cat</a>
CR Croatian Web Archive	x	<a href="http://haw.nsk.hr">haw.nsk.hr</a>
CZ Archive of the Czech Web	x	<a href="http://webarchiv.cz">webarchiv.cz</a>
TW National Taiwan University	x	<a href="http://webarchive.lib.ntu.edu.tw">webarchive.lib.ntu.edu.tw</a>
AIT Archive-It	x	<a href="http://www.archive-it.org">www.archive-it.org</a>

153329 URIs from DMOZ, archive fulltext search, IA logs, Memento aggregator logs

# Temporal Spread

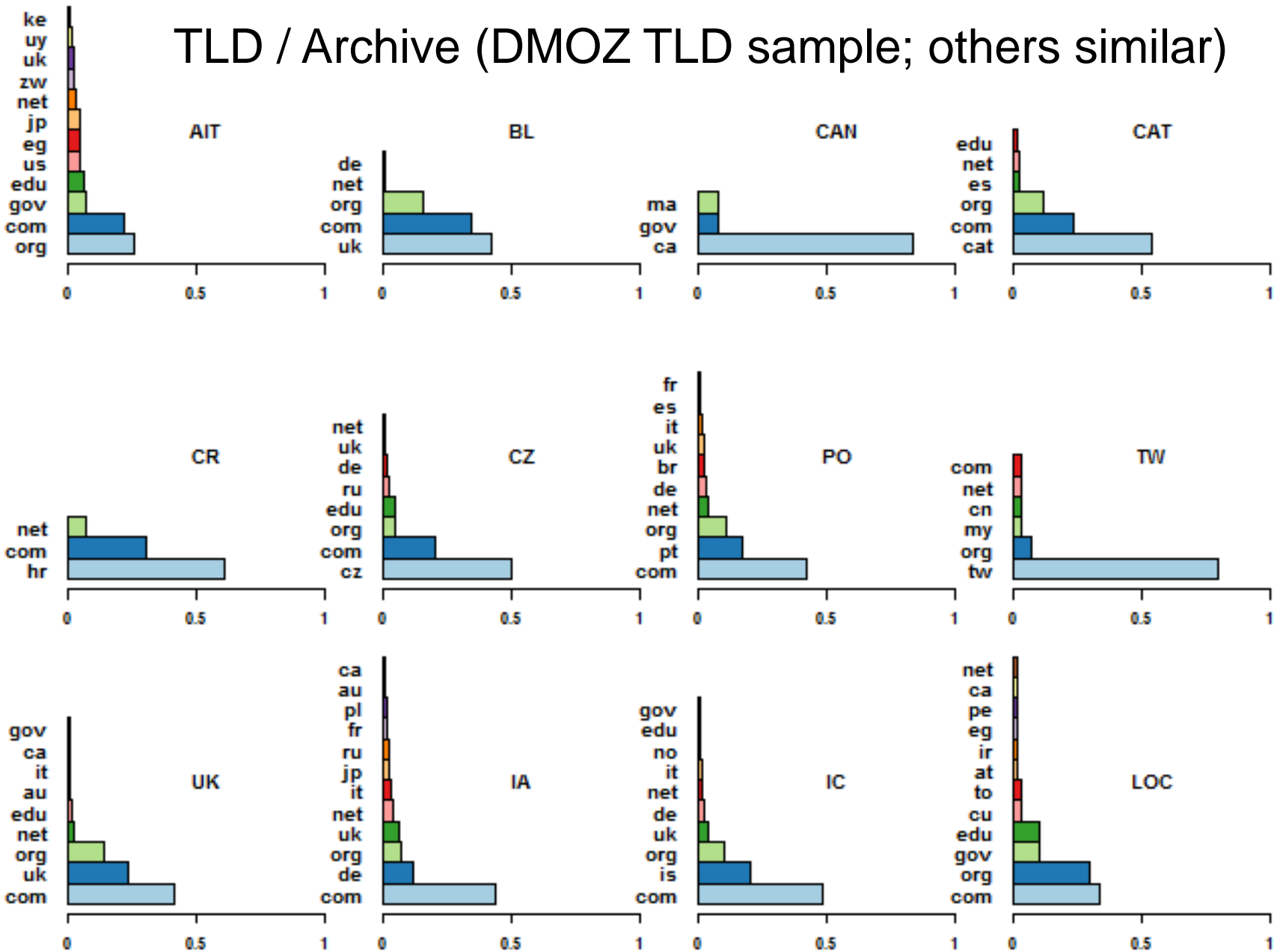


# Rate of Acquiring URI-Rs, URI-Ms



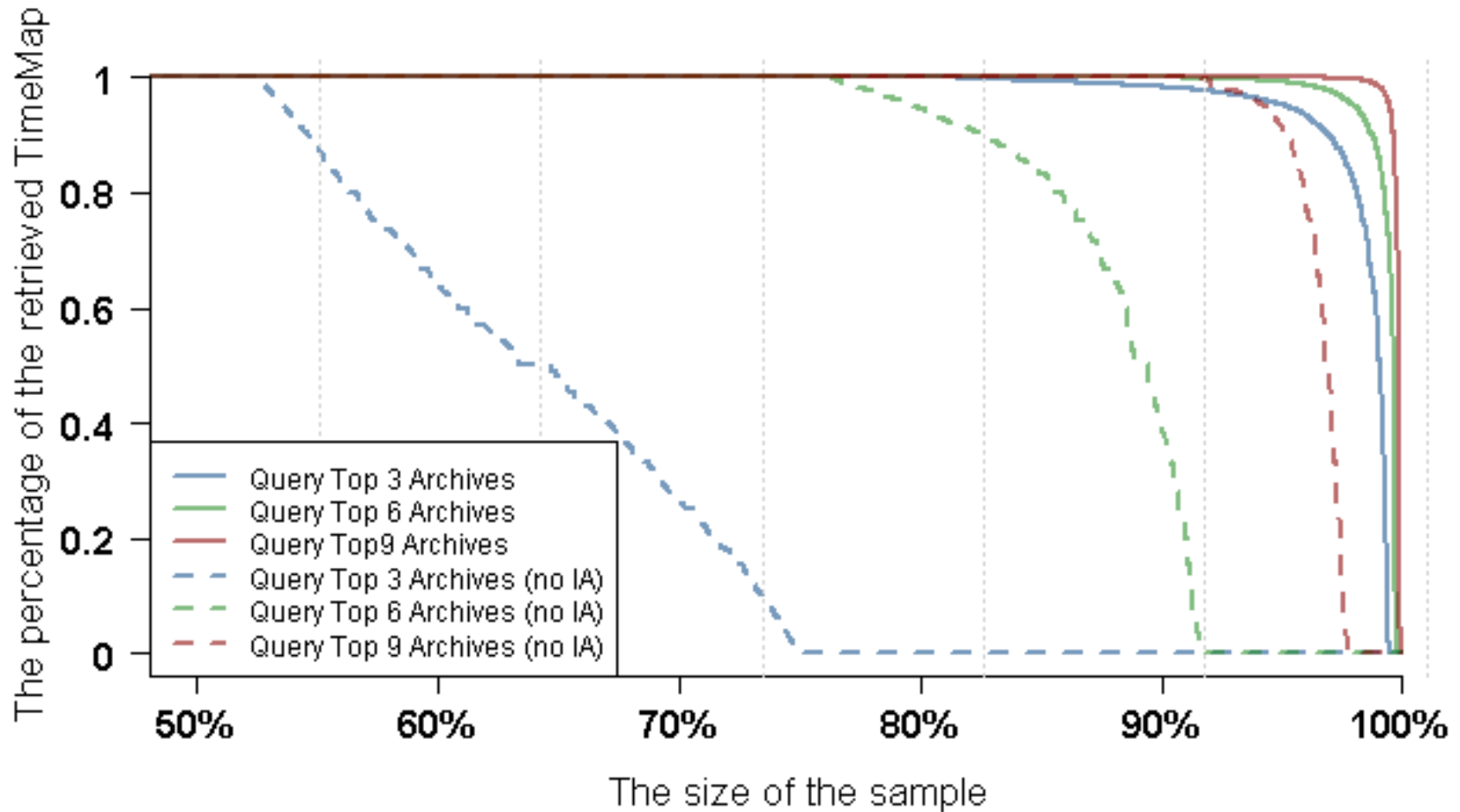


# TLD / Archive (DMOZ TLD sample; others similar)





# Using Only Top-k Archives for URI Lookup Yields Good Results



Even when there are 100s of archives, we only need to talk to a few.