

9-23-2013

Profiling Web Archive Coverage for Top-Level Domain & Content Language

Ahmed AlSum

Old Dominion University

Michele C. Weigle

Old Dominion University, mweigle@odu.edu

Michael L. Nelson

Old Dominion University, mnelson@odu.edu

Herbert Van de Sompel

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



Part of the [Archival Science Commons](#)

Recommended Citation

AlSum, Ahmed; Weigle, Michele C.; Nelson, Michael L.; and de Sompel, Herbert Van, "Profiling Web Archive Coverage for Top-Level Domain & Content Language" (2013). *Computer Science Presentations*. 12.
https://digitalcommons.odu.edu/computerscience_presentations/12

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.

Profiling Web Archive Coverage for Top-Level Domain & Content Language

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

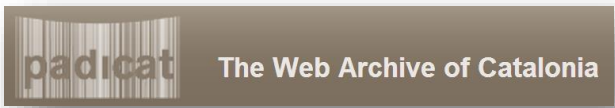
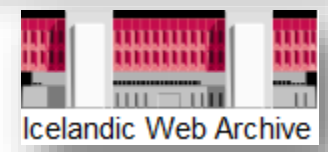


International Conference on Theory and Practice of Digital Libraries
September 22-26, 2013
Valletta, Malta





Aggregator




Where to find Mementos for ...



<http://www.japantimes.co.jp/>



 The National Archives

臺灣大學網站典藏庫
NTU Web Archiving System




Where to find Mementos for ...

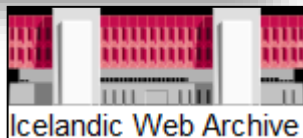


<http://www.japantimes.co.jp/>



 The National Archives

臺灣大學網站典藏庫
NTU Web Archiving System

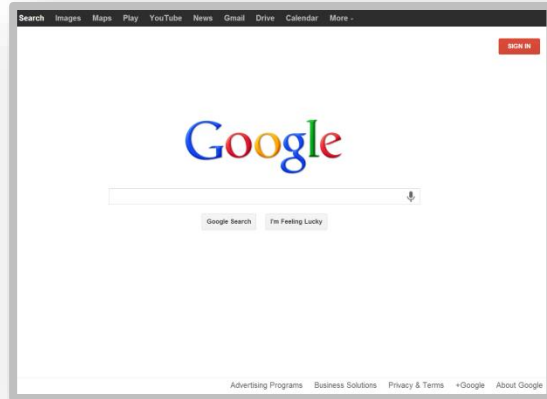


Icelandic Web Archive

WebCite



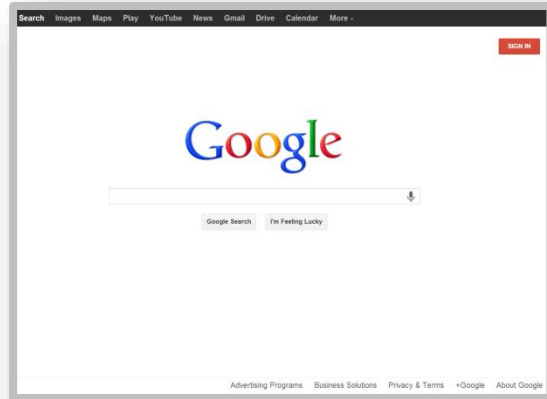
Where to find Mementos for ...



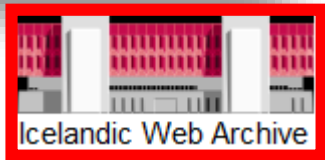
<http://www.google.com/>



Where to find Mementos for ...



<http://www.google.com/>



Research Question

Problem

- Profile public web archives according to the following dimensions:
 - Top-level domains
 - Languages
 - Growth rate
 - Archival date

Motivation

- To determine who is archiving what
- To optimize the query routing for a Memento Aggregator

Web Archives in this Experiment

	Full text	URI-lookup
Internet Archive		√
Library of Congress		√
Icelandic Web Archive		√
Library and Archives Canada	√	√
British Library	√	√
UK National Library	√	√
Portuguese Web Archive	√	√
Web Archive of Catalonia	√	√
Croatian Web Archive	√	√
Archive of the Czech Web	√	√
National Taiwan University	√	√
Archive It	√	√

Experiment Set Up

- Sample URIs from different sources
 - Details coming up
- Retrieve the TimeMap for each URI from all archives
 - A TimeMap lists all Mementos for a given URI
 - A Memento is an archived version of a resource
- Analyze
 - Details coming up

Sampling URIs

Web

1. DMOZ:Random
2. DMOZ:TLD - 2% of each TLD from DMOZ (.com, .org, .jp, etc 52 TLD)
3. DMOZ:Languages - 100 URIs for each Languages (24 lang.)

Web Archives Full Text

4. Top 1-Gram from Bing
5. Top 1000 queries term by Yahoo in 9 languages



User requests

6. IA Wayback Machine Log files
7. Memento aggregator log files

Sampling URIs - DMOZ

1. DMOZ:Random
 - 10,000 URIs randomly sampled from DMOZ directory (~5M URIs).
2. DMOZ:TLD - 2% for each TLD from DMOZ or 100 URIs whichever is greater
 - 52 TLDs (**com** 23,470) (**de** 6,332), (**org** 4,025), (**uk** 3,309), (**net** 2,073), (**it** 1,775), (**jp** 1379), (**ru** 1244), (**fr** 1154), (**pl** 1062), (**au** 764), (**ca** 642), (**at** 438), (**edu** 390), (**cz** 385), (**tr** 334), (**info** 319), (**cn** 278), (**us** 266), (**nz** 265), (**es** 238), (**ar** 213), (**no** 150), (**br** 149), (**tw** 141), (**za** 118), (**fi** 113), (100 URIs for [**ae**, **cat**, **cl**, **cu**, **eg**, **gov**, **id**, **in**, **ir**, **is**, **ke**, **kr**, **ma**, **mt**, **mx**, **my**, **na**, **pe**, **pk**, **pt**, **sa**, **to**, **uy**, **zw**])
3. DMOZ:Languages - 100 URIs for each language
 - 24 languages: Icelandic, Portuguese, Catalan, Afrikaans, Arabic, Indonesian, Chinese (Simplified), Chinese (Traditional), Dutch, Spanish, French, Greek, Hindi, Italian, Japanese, Korean, Norwegian, Persian, Polish , Russian, Turkish, Ukrainian

Sampling URIs – Web Archives Full Text

- Query the fulltext search interface of select web archives with two sets of query terms.
- 4. Top 1-Gram from Bing
 - Most are English
- 5. Top 1000 query terms from Yahoo in 9 languages
 - Excluding general keywords such as: Obama, Facebook.

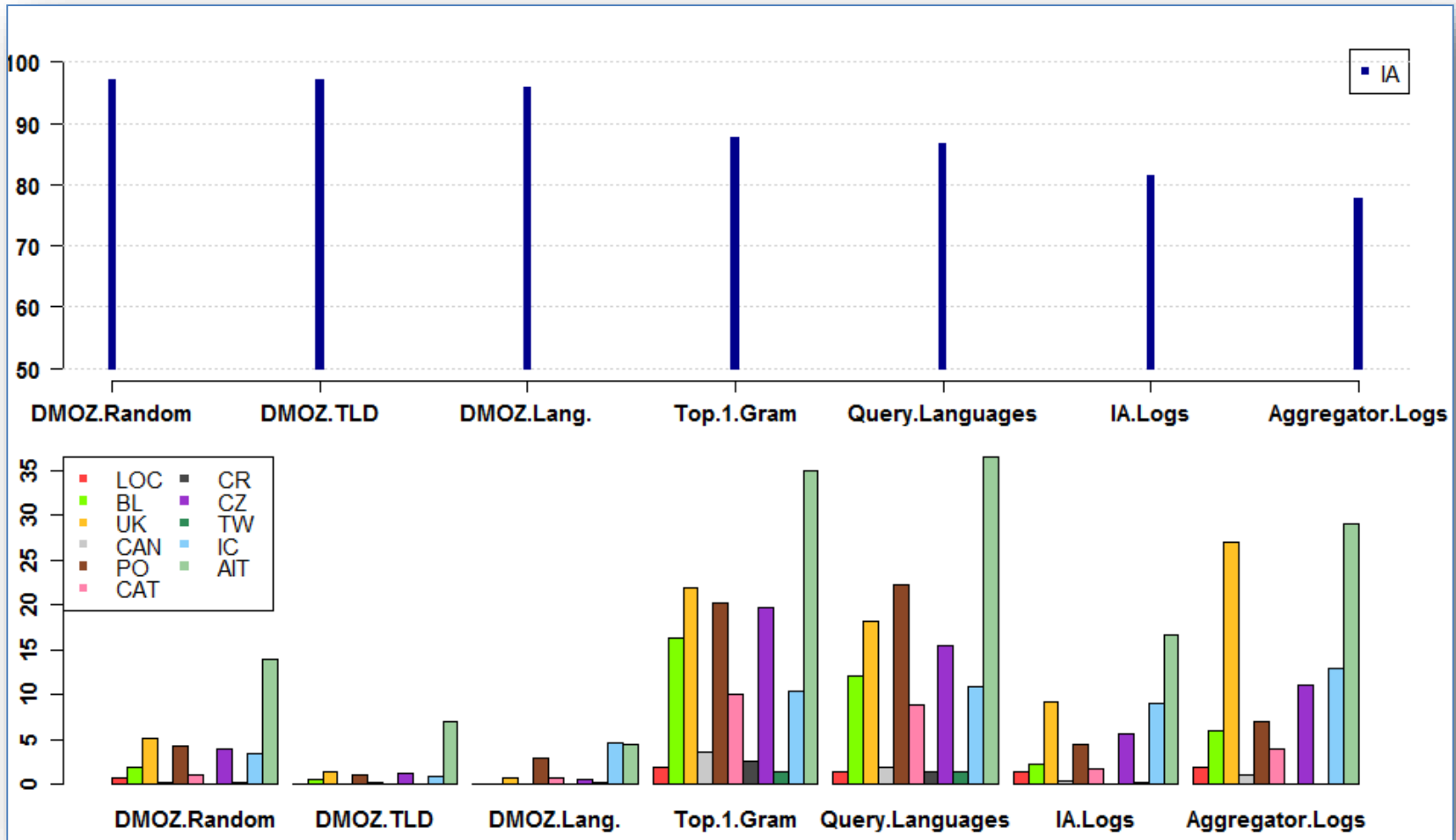
Sampling URIs – User Requests

- Sampling from user requests for archived web resources
6. Sample from IA Wayback Machine Log files
 - 1,000 URIs randomly sampled from Feb 22, 2012 to Feb 26, 2012.
 7. Sample from Memento Aggregator log files
 - 100 URIs randomly sampled from LANL Memento Aggregator between 2011 to 2013.

Archive Coverage per Sample

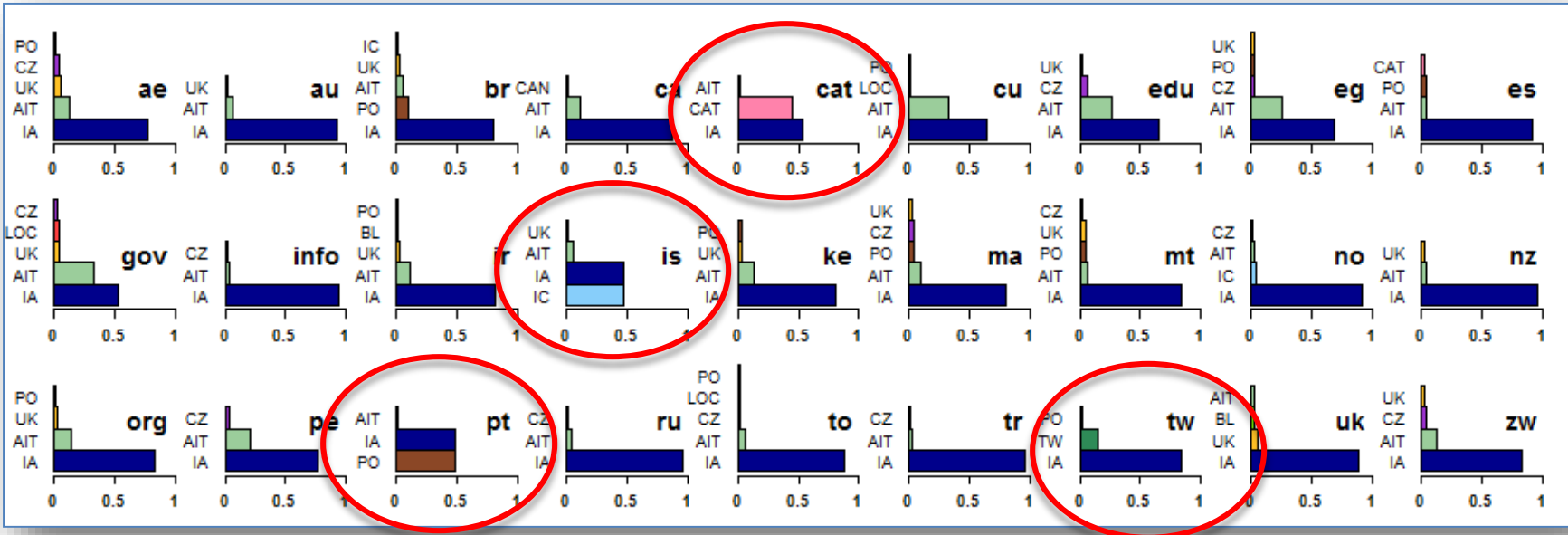
100%

35%

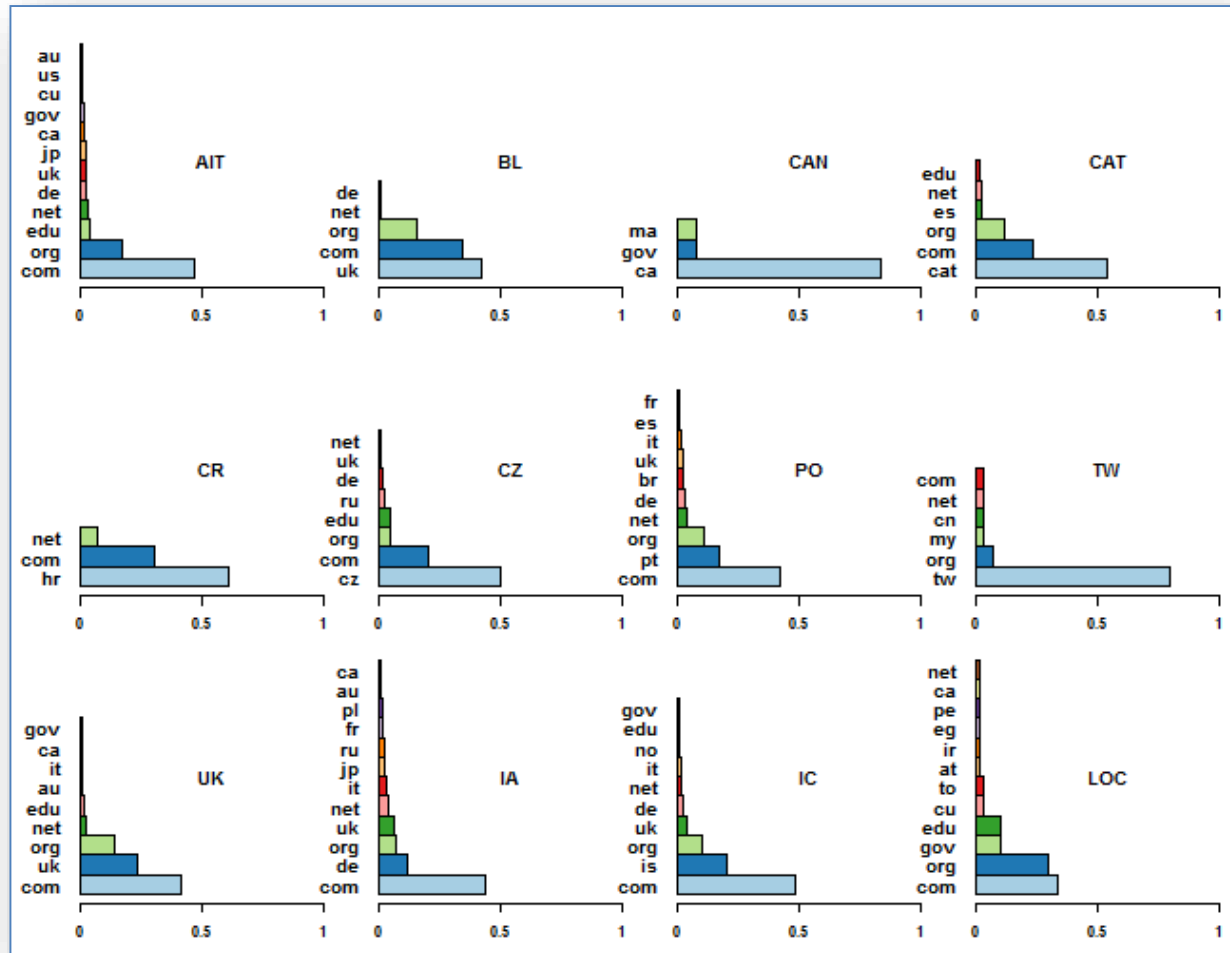


Entire Sample

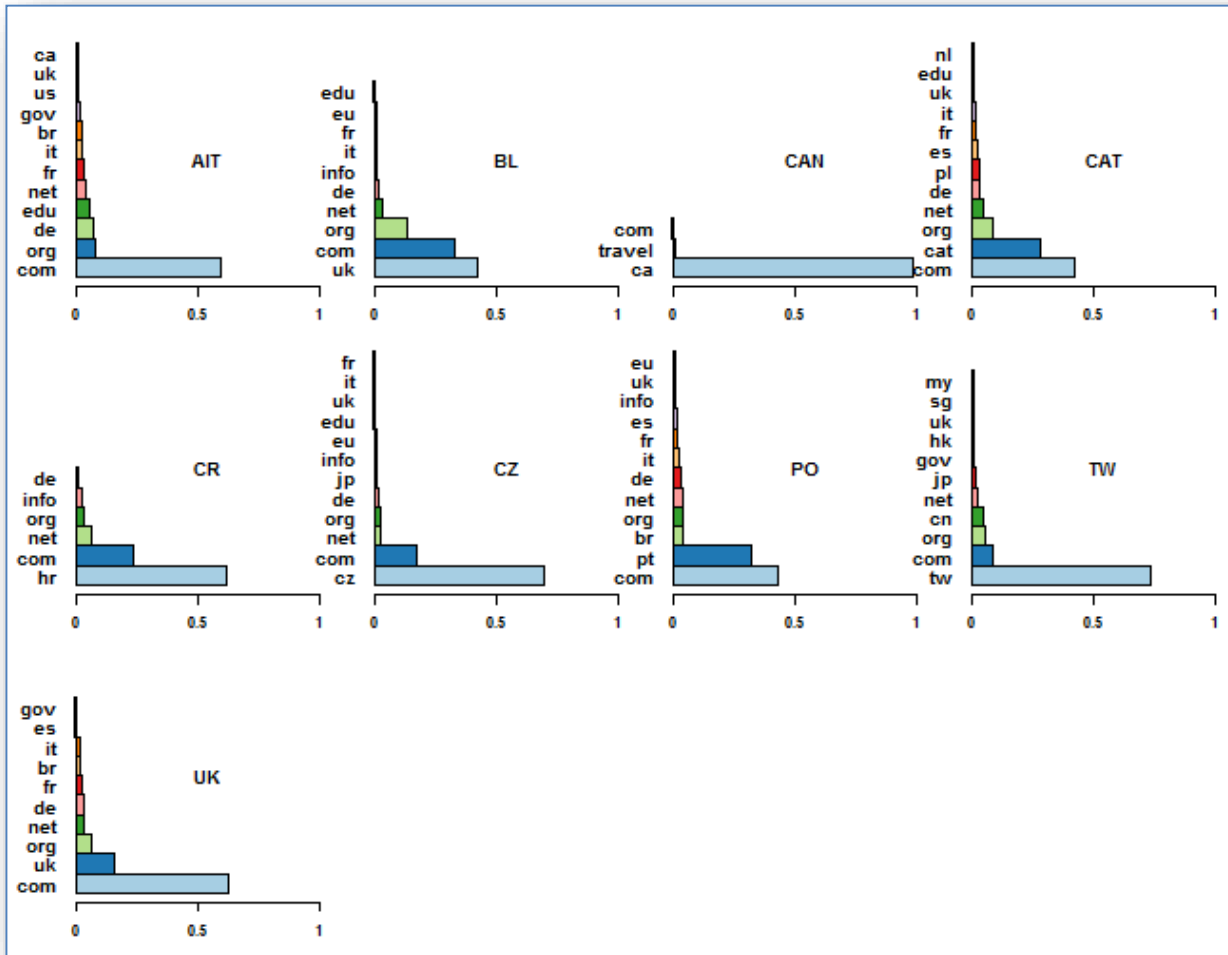
TLD Coverage across Archives (1)



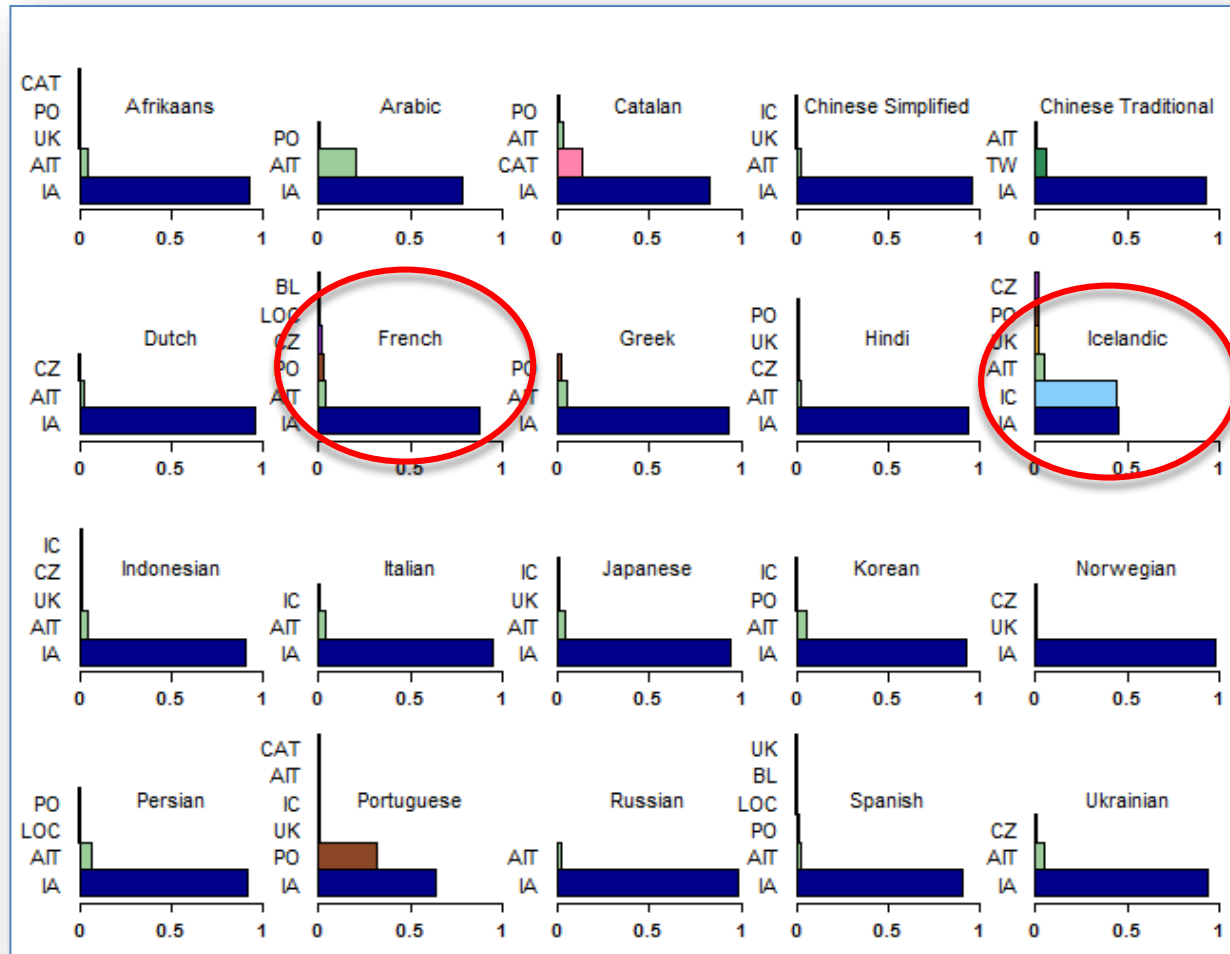
TLD Distribution per Archive



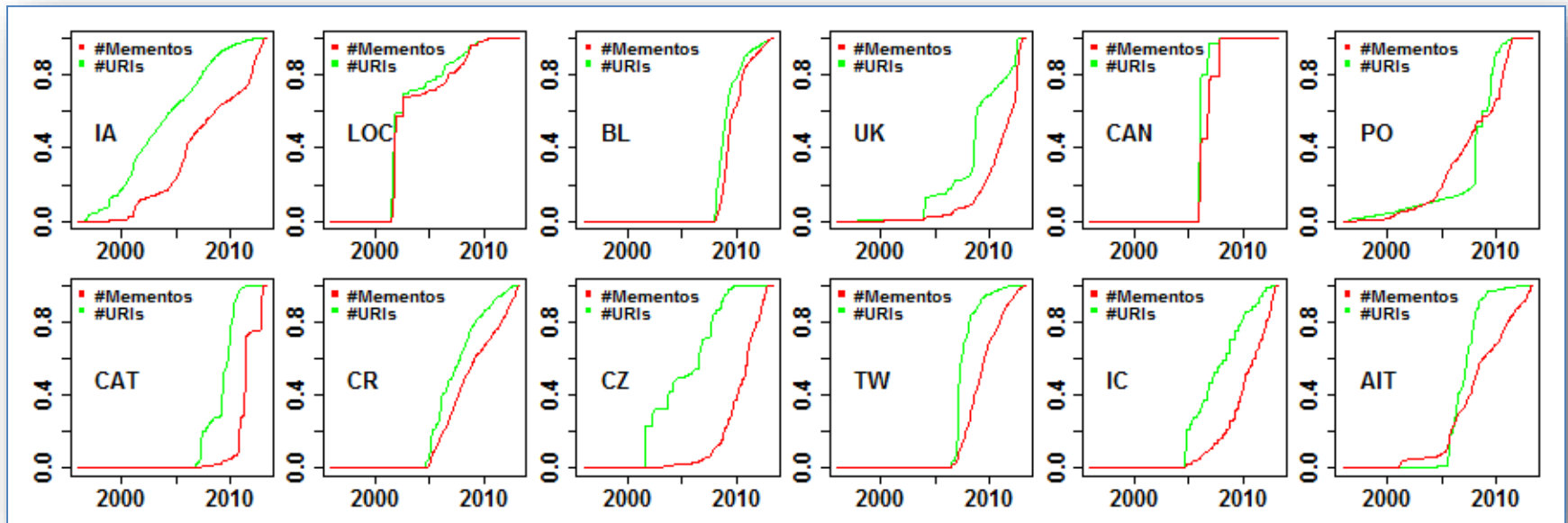
TLD Distribution per Archive



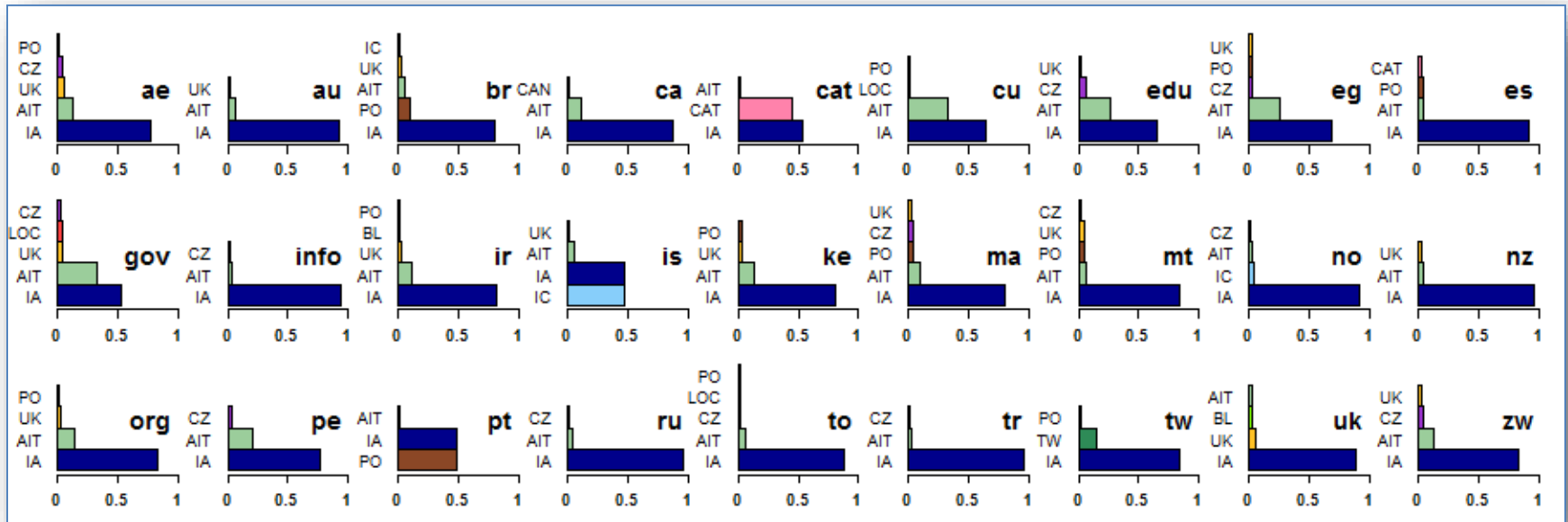
Language Coverage per Archive



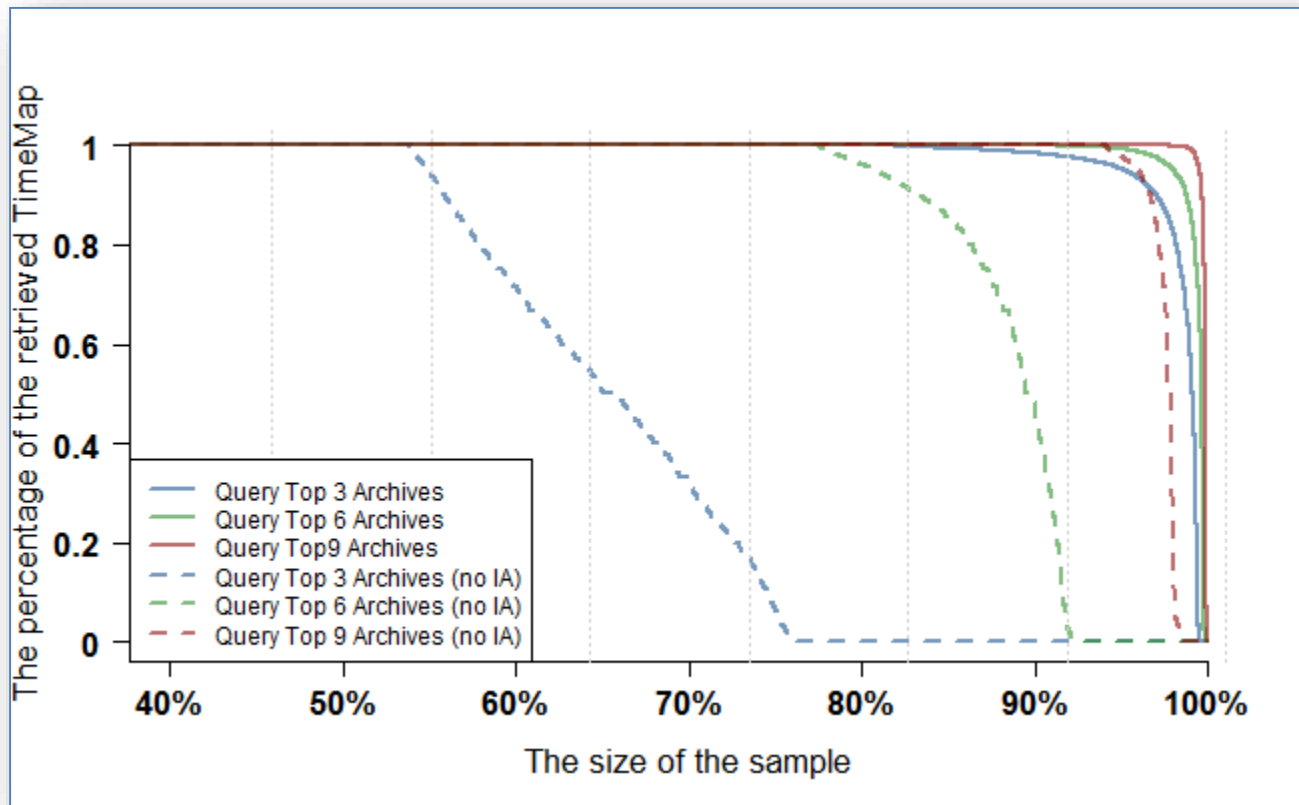
Archive Growth Rate



TLD Coverage across Archives



Query Routing Evaluation



Conclusions

- Introduced automated methodology to profile web archives using available infrastructure, no privileged access
- Coverage:
 - Internet Archive provides broad coverage
 - National archives have good coverage for their domains
 - Surprising coverage by certain archives
- Query Routing:
 - In 84% of the cases, all existing Mementos for a TLD can be found by using IA and two additional top archives for a TLD
 - In 55% of the cases, all existing Mementos for a TLD can be found by using the top 3 archives for a TLD, excluding IA