# Comparing the Archival Rate of Arabic, English, Danish, and Korean Language Web Pages

Lulwah M. Alkwai, Michael L. Nelson, **Michele C. Weigle** (@weiglemc)

Web Sciences and Digital Libraries (WS-DL) Group (@WebSciDL)

Department of Computer Science
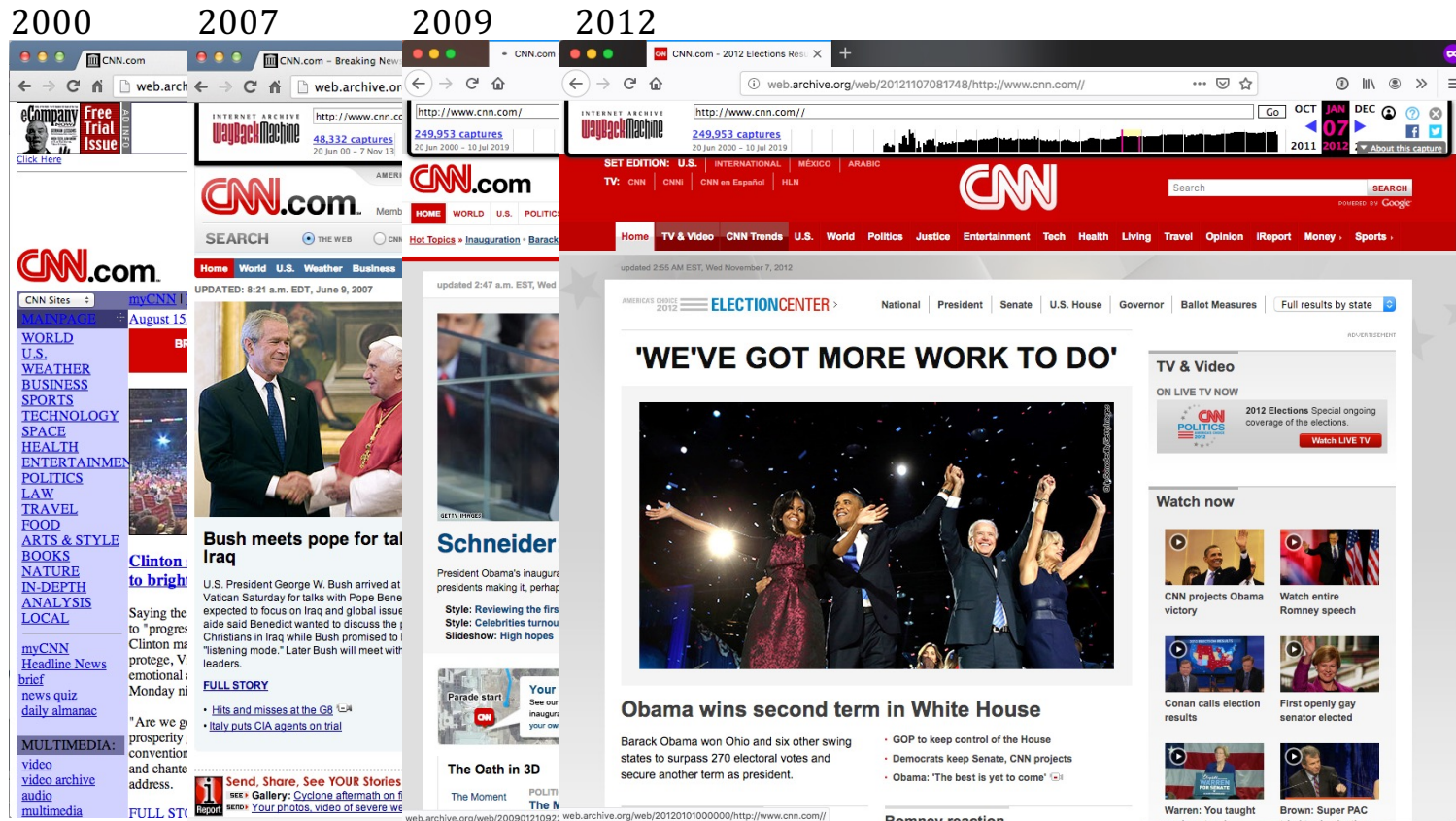
Old Dominion University

Norfolk, Virginia, USA

# Web archives are collections of web pages of the past

**2000**  **2007**  **2009**  **2012**

# Web archives are essential for studying recent history and culture



https://twitter.com/NetPreserve/status/1141321443373920256 *(photo cropped and enlarged)*



http://web.archive.org/web/19970222174751/http://www1.geocities.com/

# The Internet Archive holds the largest web archive



https://archive.org/web

# But it's not the only one



http://timetravel.mementoweb.org/list/19990518173206/http://geocities.com



https://en.wikipedia.org/wiki/List_of_Web_archiving_initiatives

# We've studied recent (2010s) events in the Middle East

- Iranian Elections and Protests - Jun 2009
  - SalahEldeen and Nelson, TPDL 2012

- Egyptian Revolution - Jan 20 - Mar 1, 2011
  - SalahEldeen and Nelson, TPDL 2012
  - AlNoamany, Weigle, Nelson, ACM WebSci 2017

- Syrian Uprising - Mar 2012
  - SalahEldeen and Nelson, TPDL 2012

- Egypt's Presidential Election - 2012
  - AlNoamany, Weigle, Nelson, TPDL 2015, IJDL 2016

# But, we can only study the past Web that already exists in the archives



11% of resources shared in social media disappear each year

Hany M. SalahEldeen, Michael L. Nelson, Losing My Revolution: How Many Resources Shared on Social Media Have Been Lost?, *Proceedings of TPDL* 2012.

# How well is the Arabic Web archived?

- Arabic is the 4th most popular language on the Internet

- Anecdotally known that web archives and search engines favor Western and English language pages

# We investigated the state of archival of Arabic language web pages in 2014-2015

- Gathered Arabic language web pages

- Analyzed domains, TLDs, GeoIP

- Analyzed presence in Google's index and web archives

- Compared this to the archival and indexing of English, Danish, and Korean language web pages

# *How to gather and detect Arabic language web pages?*

# Gathered URIs from Arabic website directories

| Directory | Registered Country | Year Estab. | Directory URI | URIs |
|-----------|-------------------|-------------|---------------|------|
| DMOZ | US | 1999 | Dmoz.org/world/arabic | 4,086 |
| Raddadi | Saudi Arabia | 2000 | Raddadi.com | 3,271 |
| Star28 | Lebanon | 2004 | Star28.com | 8,386 |
| **Total** | | | | **15,743** |

- 15,743 seed URIs
  - 15,092 were unique
  - **11,014** were live on the Web in March-May 2014

# We used four methods to determine the language of the 11,014 URIs

- HTTP Content-Language header

- HTML Title tag          https://github.com/kent37/guess-language

- Language Detection API         http://detectlanguage.com

- Trigrams         https://github.com/decultured/Python-Language-Detector

*Sample of test results evaluated by a native speaker (first author)*

# HTTP Content-Language classified 41% of the seed URIs as Arabic

Language Detection API    HTML Title Tag

Trigram                                    HTTP Content

386          1479

981    102    757

144    782    53    1615

872

60    56    926    141

97

HTTP response header
ex:
`Content-Language: `**`ar`**

HTML tag
ex:
```
<html dir="rtl"
xmlns="http://www.w3.org/19
99/xhtml" xml:lang="ar"
lang="ar">
```

# HTML Title tag classified 38% of the seed URIs as Arabic



Language Detection API  HTML Title Tag

Trigram  386  1479  HTTP Content

981  102  757

144  782  53  1615

872

60  141

56  926

97

Extract text from HTML title tag

ex:

<TITLE>دليل العرب الشامل</TITLE>

guess-language library
https://github.com/kent37/guess-language

# The Language Detection API test classified 39% of the seed URIs as Arabic

Language Detection API    HTML Title Tag

Trigram

HTTP Content

386
1479
981
102
757
144
782
53
1615
872
60
141
56
926
97

Extract title and text

Language Detection API
http://detectlanguage.com

# The Trigram test classified 36% of the seed URIs as Arabic



Extract title and text

Test sequences of three letters (trigrams)

Python-Language-Detector tool

https://github.com/decultured/Python-Language-Detector

# We took the union to obtain 7976 Arabic seed URIs



72.4% of the seed URIs were determined to be in Arabic

# We expanded the dataset by crawling the live Web and the past Web

- Crawled the 7976 live Arabic seed URIs, 2 levels deep
  – gathered all URIs linked from each seed URI
  – then, gathered all URIs linked from those URIs
  – 575,242 additional URIs

- Crawled the most recent memento of the Arabic seed URIs
  – 515,821 additional URIs

- Total of 663,443 unique crawled URIs gathered

# 482,905 were live



Live: 482,905
Unique: 663,443
Unavailable: 180,538

# 292,670 were live and Arabic



Live: 482,905
Unique: 663,443
Arabic: 292,670
Not Arabic: 190,235
Unavailable: 180,538

Total Arabic Dataset = **300,646 URIs**
  7,976 seed URIs
  + 292,670 crawled URIs

# *What are the characteristics of our Arabic language dataset?*

# Dataset has 17,536 unique domains

| Rank | Domain | URIs | GeoIP | Category |
|------|--------|------|-------|----------|
| 1 | Alarab.net | 284 | US | News |
| 2 | Aljarida.com | 248 | US | News |
| 3 | Arabic.cnn.com | 245 | US | News |
| 4 | Alarabiya.net | 231 | US | News |
| 5 | Ar.wikipedia.org | 230 | US | Encyclopedia |
| 6 | Aljazeera.net | 213 | US | News |
| 7 | Moheet.com | 142 | US | News |
| 8 | Facebook.com | 133 | US | Social |
| 9 | Al-sharq.com | 132 | US | Middle East Portal |
| 10 | Lakii.com | 123 | US | General Portal |
| 17 | Kuwaitclub.com.kw | 71 | Kuwait | Sport |

# First Arabic GeoIP is at rank 17

| Rank | Domain | URIs | GeoIP | Category |
|------|--------|------|-------|----------|
| 1 | Alarab.net | 284 | US | News |
| 2 | Aljarida.com | 248 | US | News |
| 3 | Arabic.cnn.com | 245 | US | News |
| 4 | Alarabiya.net | 231 | US | News |
| 5 | Ar.wikipedia.org | 230 | US | Encyclopedia |
| 6 | Aljazeera.net | 213 | US | News |
| 7 | Moheet.com | 142 | US | News |
| 8 | Facebook.com | 133 | US | Social |
| 9 | Al-sharq.com | 132 | US | Middle East Portal |
| 10 | Lakii.com | 123 | US | General Portal |
| 17 | Kuwaitclub.com.kw | 71 | Kuwait | Sport |

# 6 top domains are news websites

| Rank | Domain | URIs | GeoIP | Category |
|------|--------|------|-------|----------|
| 1 | Alarab.net | 284 | US | News |
| 2 | Aljarida.com | 248 | US | News |
| 3 | Arabic.cnn.com | 245 | US | News |
| 4 | Alarabiya.net | 231 | US | News |
| 5 | Ar.wikipedia.org | 230 | US | Encyclopedia |
| 6 | Aljazeera.net | 213 | US | News |
| 7 | Moheet.com | 142 | US | News |
| 8 | Facebook.com | 133 | US | Social |
| 9 | Al-sharq.com | 132 | US | Middle East Portal |
| 10 | Lakii.com | 123 | US | General Portal |
| 17 | Kuwaitclub.com.kw | 71 | Kuwait | Sport |

# Popular Western domains are in the Top 10

| Rank | Domain | URIs | GeoIP | Category |
|---|---|---|---|---|
| 1 | Alarab.net | 284 | US | News |
| 2 | Aljarida.com | 248 | US | News |
| 3 | Arabic.cnn.com | 245 | US | News |
| 4 | Alarabiya.net | 231 | US | News |
| 5 | Ar.wikipedia.org | 230 | US | Encyclopedia |
| 6 | Aljazeera.net | 213 | US | News |
| 7 | Moheet.com | 142 | US | News |
| 8 | Facebook.com | 133 | US | Social |
| 9 | Al-sharq.com | 132 | US | Middle East Portal |
| 10 | Lakii.com | 123 | US | General Portal |
| 17 | Kuwaitclub.com.kw | 71 | Kuwait | Sport |

# Over half have a .com TLD

| TLD | Percent |
|---|---|
| com | 57.97% |
| net | 15.07% |
| org | 6.40% |
| gov.sa | 1.94% |
| info | 1.68% |
| edu.sa | 1.27% |
| ws | 1.16% |
| org.sa | 0.97% |
| com.sa | 0.80% |
| gov.eg | 0.80% |
| Other | 11.94% |

# Only ~10% have an Arabic ccTLD

| TLD | Percent |
|---|---|
| com | 57.97% |
| net | 15.07% |
| org | 6.40% |
| gov.sa | 1.94% |
| info | 1.68% |
| edu.sa | 1.27% |
| ws | 1.16% |
| org.sa | 0.97% |
| com.sa | 0.80% |
| gov.eg | 0.80% |
| Other | 11.94% |

| ccTLD | Country | Percent |
|---|---|---|
| .sa | Saudi Arabia | 5.33% |
| .eg | Egypt | 2.00% |
| .jo | Jordan | 2.00% |
| .ae | United Arab Emirates | 1.06% |
| .kw | Kuwait | 0.82% |

# Most are geo-located in the US

| Geo-location | Percent |
|---|---|
| United States | 57.97% |
| Arabic Countries | 10.53% |
| Germany | 9.75% |
| Netherlands | 5.29% |
| France | 4.37% |
| Canada | 3.31% |
| United Kingdom | 3.07% |
| Other | 5.71% |

# Within Arabic countries, most are geo-located in Saudi Arabia

| Geo-location | Percent |
|---|---|
| United States | 57.97% |
| Arabic Countries | 10.53% |
| Germany | 9.75% |
| Netherlands | 5.29% |
| France | 4.37% |
| Canada | 3.31% |
| United Kingdom | 3.07% |
| Other | 5.71% |

| Geo-location | Percent |
|---|---|
| Saudi Arabia | 4.75% |
| Egypt | 1.97% |
| Jordan | 1.42% |
| Kuwait | 0.71% |
| United Arab Emirates | 0.67% |

# *How well are these Arabic web pages indexed and archived?*

# 53.77% of Arabic language web pages are archived

- Used Memento aggregator to determine if archived
  - checks multiple web archives

- 161,678 URIs were archived

- 97% of those were found in the Internet Archive

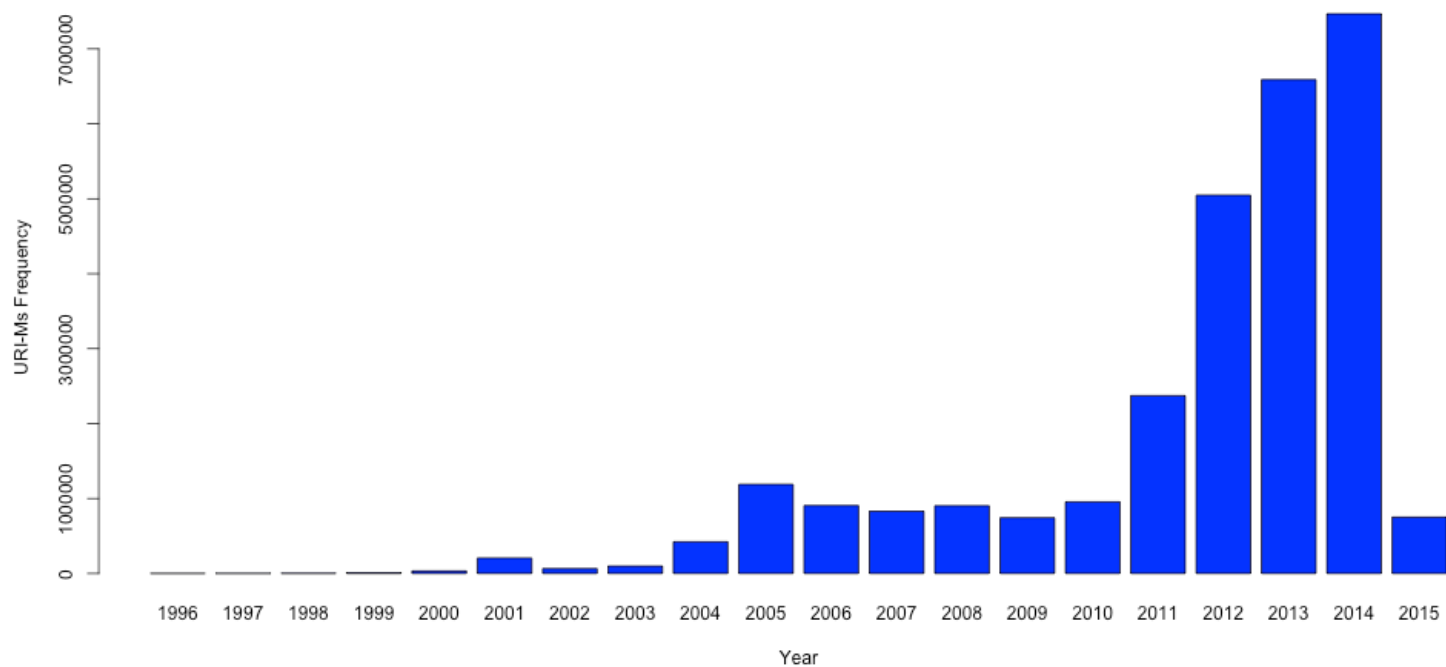- Mementos also found in 9 other archives

# 6 out of the 10 most archived are news websites

| URI-Rs | Mementos | Category |
|---|---|---|
| gulfup.com | 10,987 | File Sharing |
| masrawy.com | 9,144 | Egyptian portal |
| arabic.cnn.com | 9,022 | News |
| aljazeera.net | 8,906 | News |
| maktoob.yahoo.com | 8,478 | Search Engine |
| shorooknews.com | 7,548 | News |
| arabnews.com | 6,274 | News |
| bbc.co.uk/arabic | 6,268 | News |
| ahram.org.eg | 5,347 | News |
| google.com.sa | 4,968 | Search Engine |

# Many are Arabic versions of globally popular sites

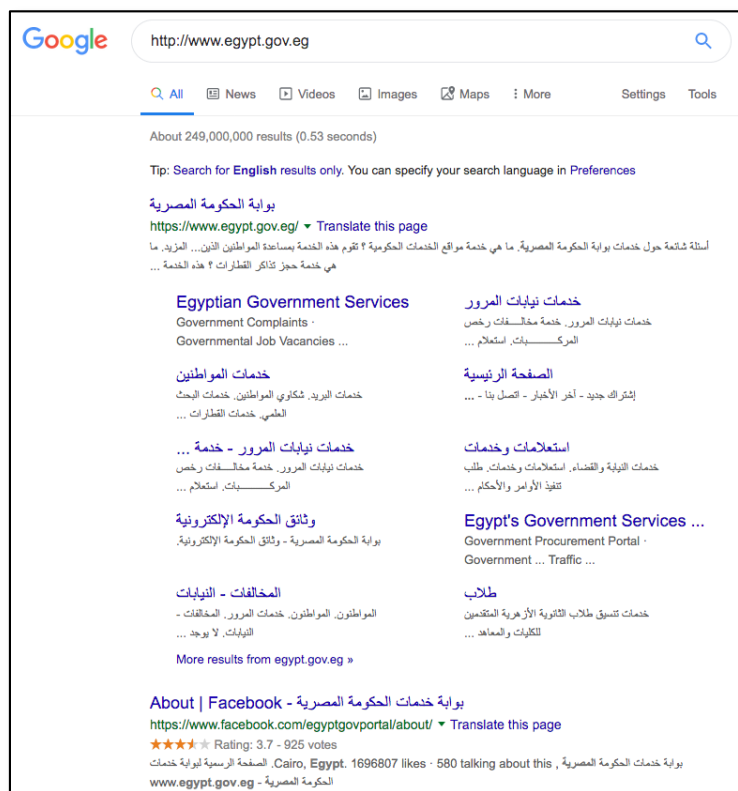| URI-Rs | Mementos | Category |
|---|---:|---|
| gulfup.com | 10,987 | File Sharing |
| masrawy.com | 9,144 | Egyptian portal |
| arabic.cnn.com | 9,022 | News |
| aljazeera.net | 8,906 | News |
| maktoob.yahoo.com | 8,478 | Search Engine |
| shorooknews.com | 7,548 | News |
| arabnews.com | 6,274 | News |
| bbc.co.uk/arabic | 6,268 | News |
| ahram.org.eg | 5,347 | News |
| google.com.sa | 4,968 | Search Engine |

# Most mementos are from recent years



*Analysis done in 2015*

# We looked at indexing of Arabic seed URIs

- Used Google Custom Search API to query
  - limited to 1000 queries/day
  - tested only seeds

# 69% of Arabic language seed URIs were indexed by Google

| Arabic Seed Dataset | |
|---|---|
| (Live, Indexed, Archived) | Percent |
| (1, 1, 1) | 43.34% |
| (1, 1, 0) | 25.59% |
| (1, 0, 1) | 15.27% |
| (1, 0, 0) | 15.76% |

*All seeds were live*

- 82% of those listed in DMOZ were indexed

- 74.6% of the top-level seeds (path depth = 0) were indexed

# 58% of seeds were archived

| Arabic Seed Dataset | |
|---|---|
| (Live, Indexed, Archived) | Percent |
| (1, 1, 1) | 43.34% |
| (1, 1, 0) | 25.59% |
| (1, 0, 1) | 15.27% |
| (1, 0, 0) | 15.76% |

*But, 42% were not archived*

# 43% were indexed and archived

Good!
(discovered
and saved)

| Arabic Seed Dataset | |
|---|---|
| (Live, Indexed, Archived) | Percent |
| (1, 1, 1) | 43.34% |
| (1, 1, 0) | 25.59% |
| (1, 0, 1) | 15.27% |
| (1, 0, 0) | 15.76% |

# 31% were not indexed by Google

| Arabic Seed Dataset | |
|---|---|
| (Live, Indexed, Archived) | Percent |
| (1, 1, 1) | 43.34% |
| (1, 1, 0) | 25.59% |
| (1, 0, 1) | 15.27% |
| (1, 0, 0) | 15.76% |

# 16% were neither indexed nor archived

| Arabic Seed Dataset | |
|---|---|
| (Live, Indexed, Archived) | Percent |
| (1, 1, 1) | 43.34% |
| (1, 1, 0) | 25.59% |
| (1, 0, 1) | 15.27% |
| (1, 0, 0) | 15.76% |

Bad!
(undiscovered and not saved)

# *How does this compare to other languages?*

# We chose English, Danish, and Korean

- English
  - most popular language on the Internet
  - over 65 countries have English as an official language

- Danish
  - European language
  - 96% of population in Denmark uses the Internet
  - Denmark has government initiative to archive Danish cultural heritage on the Web

- Korean
  - Asian language
  - 92% of population in South Korea uses the Internet (highest in Asian countries)

# We gathered seed URIs from DMOZ

- English
  - sample of size 10,000

- Danish
  - sample of size 10,000
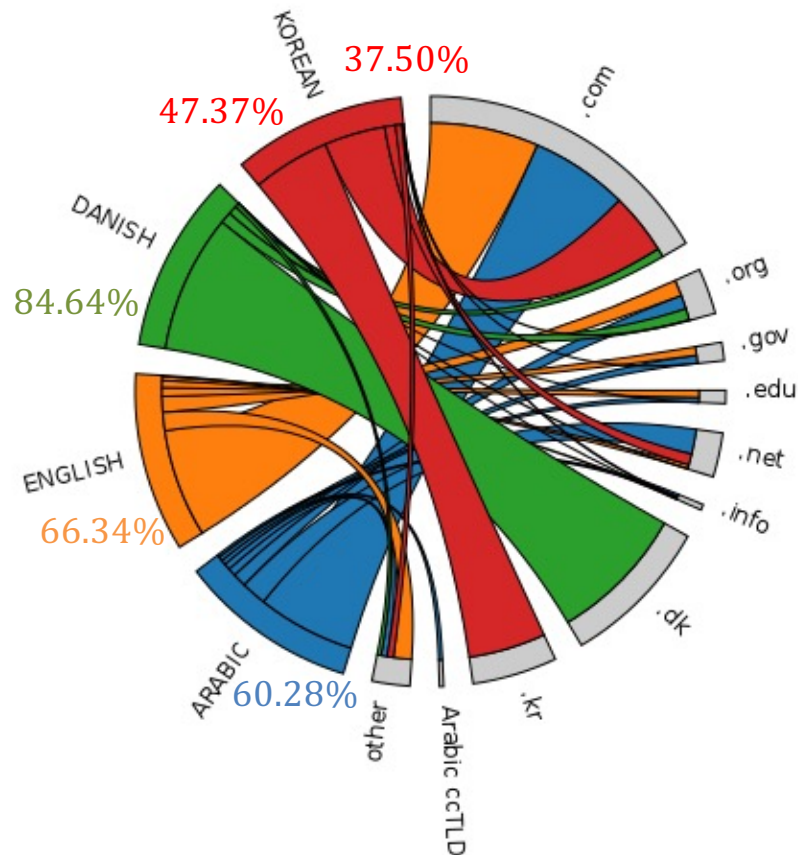
- Korean
  - all 2,347 URIs in DMOZ

# Crawled the seed URIs to expand the dataset

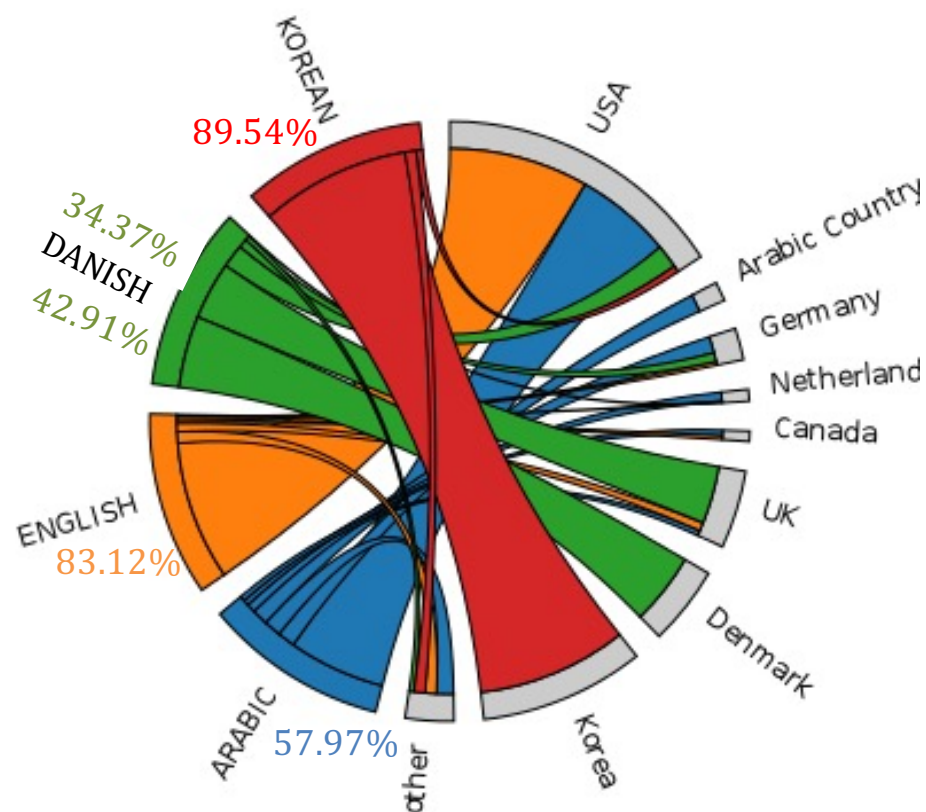|  | Arabic | English | Danish | Korean |
|---|---|---|---|---|
| Seeds | 15,742 | 10,000 | 10,000 | 2,347 |
| Live | 11,014 | 9,384 | 9,245 | 2,070 |
| Language | 7,976 | 8,576 | 6,331 | 1,157 |
|  |  |  |  |  |
| Crawled | 663,443 | 224,249 | 174,369 | 16,016 |
| Live | 482,905 | 176,261 | 131,484 | 11,099 |
| Language | 292,670 | 137,950 | 99,019 | 7,965 |
|  |  |  |  |  |
| **Total** | **300,646** | **146,526** | **105,350** | **9,482** |

**562,004 total**

December 2015 - March 2016

# English and Arabic web pages were .com, Danish and Korean were ccTLD

# English and Arabic web pages were located in the US, Danish and Korean were in their countries

# Arabic is archived less than English, but more than Danish and Korean

|  | Arabic | English | Danish | Korean |
|---|---|---|---|---|
| Seeds and Crawled | 300,646 | 146,526 | 105,350 | 9,482 |
| Archived | 161,678 | 107,398 | 41,703 | 3,972 |
| Percent | 53.77% | 73.30% | *39.59% | 41.89% |

\* Danish government archive is dark (not publicly available)

# >90% of pages listed in DMOZ were archived, regardless of language

|  | **Arabic** | **English** | **Danish** | **Korean** |
|---|---|---|---|---|
| DMOZ Seeds | 2,904 | 8,576 | 6,331 | 1,157 |
| Archived | 2,774 | 8,014 | 6,164 | 1,358 |
| Percent | 95.52% | 93.44% | 97.36% | 89.52% |

\* DMOZ closed in 2017, but curlie.org began in 2018 with DMOZ data

# There are a few caveats

- DMOZ is no longer operational

- Some web pages are multilingual

- Not easy to characterize "the Arabic Web"
  - 67% of Arabic dataset had neither Arabic ccTLD nor Arabic geo-location

- The language of a web page may shift over time

# *How well are Arabic language web pages archived?*

# Only about half of Arabic language web pages are archived

- Analyzed over 500,000 web pages (2014-2016)

- Arabic is archived less than English, but more than Danish and Korean

- Most Arabic and English pages are geo-located in the US, Danish in Denmark, and Korean in South Korea

- Web pages present in DMOZ are highly likely to be archived, regardless of language

Old Dominion University
Department of Computer Science

Web Science and Digital Libraries (WS-DL)
https://ws-dl.cs.odu.edu/        @WebSciDL