LA-UR-21-32308

Approved for public release; distribution is unlimited.

Title: Developing Bloom Filters for Web Archives' Holdings

Author(s): Klein, Martin

Balakireva, Lyudmila Leonidovna

Hulob, Karolina Rudomino, Ingeborg Celjak, Drazenko

Intended for: Report

Issued: 2021-12-17









Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by Triad National Security, LLC for the National Nuclear Security Administration of U.S. Department of Energy under contract 89233218CNA000001. By approving this article, the publisher recognizes that the U.S. Government retains nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher dientify this article as work performed under the auspices of the U.S. Department of Energy. Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

Developing Bloom Filters for Web Archives' Holdings

Martin Klein, Lyudmila Balakireva

Los Alamos National Laboratory (LANL) {@mart1nkle1n, @milbala}

Karolina Holub, Inge Rudomino

National and University Library in Zagreb, Croatian Web Archive (HAW) {@KarolHolu, @ingeRudomino}

Drazenko Celjak

SRCE - University of Zagreb University Computing Centre @dceljak



Thanks to the IIPC for generously funding this project!





Recall, the Problem

- Archival holdings largely unknown
 - To the public
 - To web archives
- Searching through a CDX file
 - Not available to the public
 - Requires insight knowledge and time
- CDX file sharing w/ 3rd parties
 - Not practical
 - Potential legal implications unclear



https://twitter.com/anjacks0n/status/466690812269846528





Investigating a Potential Solution

- Bloom Filters (see: https://llimllib.github.io/bloomfilter-tutorial/)
 - Data structure to reveal whether an element is present in a set
 - Database of hashed URIs that are part of a WA's holdings
- Pros:
 - Fast lookup (hash URI, check for match)
 - No active publication of plain URIs in archive
 - Benefit for dark/hybrid archives
- Unknowns:
 - Cost to create a BF
 - Performance
 - Scale





At the previous IIPC webinar...

- Sample CDX file from HAW
 - Contained 180 million URIs
- BF library:
 - Max of 200mio URLs
 - Desired FPR: 1%
 - Number of hash functions: 5
 - https://github.com/Bagend/Orestes-Bloomfilter





At the previous IIPC webinar...

Mode	Redis DB	
Partner	LANL	
Ingested entries	180,379,433	
Time to ingest	7,273.6s (121.2min)	
Size of the BF	246MB	
Query response time	0.033ms	
FPR	0.66%	





Since then...

BF implementation in two possible modes:

- Redis-based
 - Persistent store
 - Possibly distributed across systems
- Memory-based
 - Not persistent
 - "local"





Redis vs. Memory Mode

Mode	Redis DB		Memory	
Partner	LANL	HAW	LANL	HAW
Ingested entries	180,379,433	180,379,433	180,379,433	180,379,433
Time to ingest	7,273.6s (121.2min)	10,759.2s (179.3min)	1,382.7s (23m)	1505.5s (25.1min)
Size of the BF	246MB	246MB	246MB	246MB
Query response time	0.033ms	0.051ms	0.0008ms	0.0008ms
FPR	0.66%	0.66%	0.66%	0.66%
Time to serialize to CSV	16.8s	17.8s	2.6s	3.4s

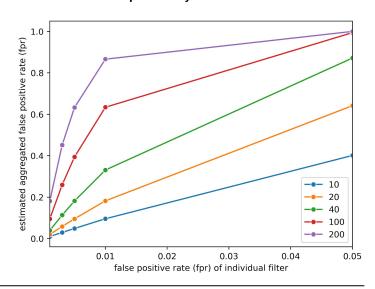




How does this scale?

- Redis mode:
 - Key limit of 512MB, so BF can not exceed 512MB in size
 - HAW, for example, would currently require ~ 1GB
 - Dynamic BFs previously introduced
 - FPR increases significantly with each added BF
- Memory mode:
 - Got unlimited memory?
 - BF library:
 - uses integer values
 - BF can not grow beyond 268.44MB

Estimated fpr of Array of identical Bloomfilters







Our Proposal

Modified DBFs:

- HAW: 800mio URLs & 200mio added per year
- Creation of 16 BFs, 200mio URLs each
- Name each BF after a hex character ('0', '1', '2', '3', '4', '5', '6', '7', '8', '9', 'a', 'b', 'c', 'd', 'e', and 'f')
- special routing rule for insertion, based on first hex character of the URL's MD5 hash value
- each URL will be inserted in its corresponding filter, matched by name
- Distribution done by hash values → filters have equal population of URLs (same probability for insertion)





Our Proposal

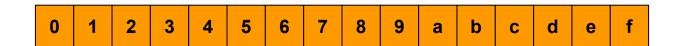
Modified DBFs:

- Lookup done by routing for appropriate BF based on URL's first hex character
- Avoids parallel querying, main bottleneck in DBFs
- For HAW, 16 BFs:
 - Desired FPR: 1% at 3.8 GB
 - Desired FPR: 5% at 1.6 GB
 - Good for next 10+ years
- Can be extended to first 2 hex characters, or 32 BFs, or





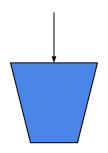


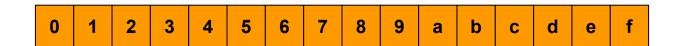






https://haw.nsk.hr/en

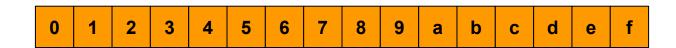






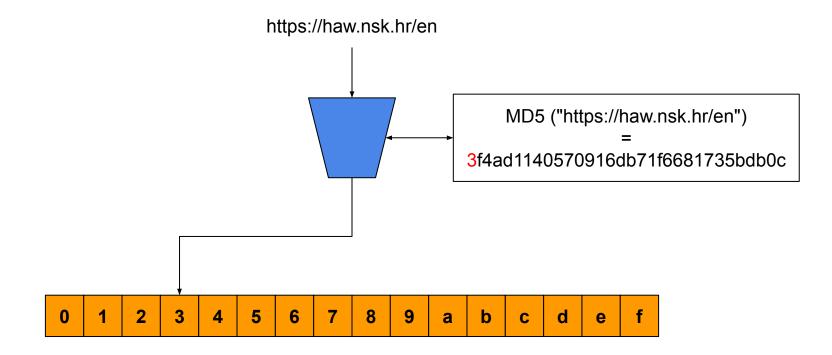






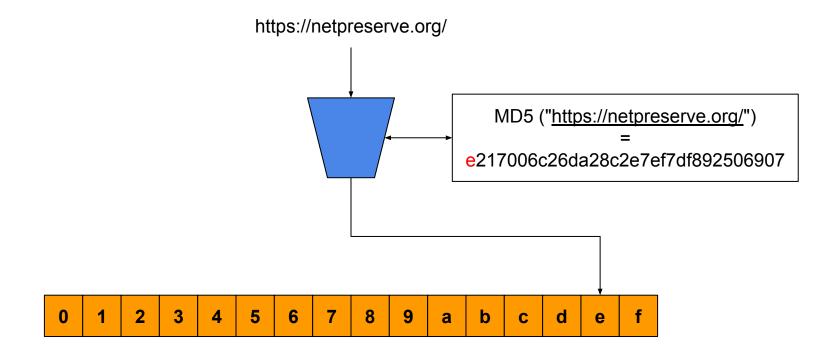
















Our Proposal - Implemented

```
♠ martin — -bash — 73×18
pn1935049:~ martin$ curl -I "http://bloom.mementodepot.org/haw/https://ww
w.zzzr.hr/Y/m/Europe/Europe/Zagreb"
HTTP/1.1 200 OK
Server: nginx/1.18.0
Date: Sun, 12 Dec 2021 23:22:30 GMT
Content-Type: application/octet-stream
Content-Length: 0
Connection: keep-alive
pn1935049:~ martin$
```

http://bloom.mementodepot.org/haw/https://www.zzzr.hr/Y/m/Europe/Europe/Zagreb



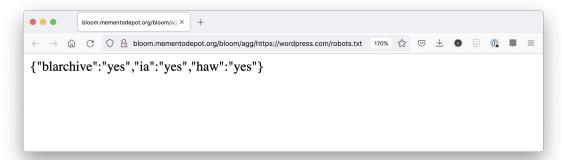


Pilot of Federated Search via BFs

- Implemented a simple lookup service, based on CDX files from:
 - HAW
 - Internet Archive
 - British Library
 - http://bloom.mementodepot.org/bloom/agg/{uri}
 - Takes URI as parameter







http://bloom.mementodepot.org/bloom/agg/https://wordpress.com/robots.txt



http://bloom.mementodepot.org/bloom/agg/http://www.econ.yale.edu/smith/econ116a/keynes1.pdf





Concept to Share BFs

- Croatian Web Archive (HAW)
 - Testing environment vs. production environment
 - Goal? To share archive's holdings
- The process of sharing an archive's BF has the following steps:
 - BF creation in memory
 - Serialization of BF into plain text file
 - Creation/update of sitemap, possibly with added metadata.





HAW's Sitemap

```
https://haw.nsk.hr/bloom-filter x +
                                                          Q 🖞 🛊 🐧 🕢 🕖 🤌 🛊 🗑 :
← → C n haw.nsk.hr/bloom-filter/sitemap.xml
This XML file does not appear to have any style information associated with it. The
document tree is shown below.
▼<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9"
 xmlns:dc="http://purl.org/dc/elements/1.1/">
 ▼<url>
     <loc>https://haw.nsk.hr/bloom-filter/haw-2020-hr-domain-
     harvest.csv</loc>
     <lastmod>2021-12-15</lastmod>
     <dc:title>Croatian TLD (hr) harvest bloom filter</dc:title>
     <dc:publisher>National and University Library in Zagreb (NSK)
     </dc:publisher>
     <dc:date>2021</dc:date>
     <dc:description>Bloom filter of URLs harvested during 2020 domanin
     crawl</dc:description>
     <dc:format>text/csv</dc:format>
     <dc:format>Number of URLs 180,379,433</dc:format>
     <dc:format>Size 319502017B</dc:format>
    </url>
 </urlset>
```

https://haw.nsk.hr/bloom-filter/sitemap.xml





Concept to Share (Use) BFs

- The process of using BF exposed by archives:
 - sitemap discovery
 - parse the sitemap XML file
 - extract the URL of new or updated BF files
 - get the serialized plain text file
 - ingest the file into the local BF implementation.





Final Thoughts

- Code currently undergoing LANL-internal review
- Open Source release as soon as approved
- This really works!!! (with pros and cons)
- Likely most suitable for:
 - Smaller archives
 - Individual collections
 - Live lookup of URLs during distributed crawl of (topic) collection





Developing Bloom Filters for Web Archives' Holdings

Martin Klein, Lyudmila Balakireva

Los Alamos National Laboratory (LANL) {@mart1nkle1n, @milbala}

Karolina Holub, Inge Rudomino

National and University Library in Zagreb, Croatian Web Archive (HAW) {@KarolHolu, @ingeRudomino}

Drazenko Celjak

SRCE - University of Zagreb @dceljak



Thanks to the IIPC for generously funding this project!





HAW's CDXs

File name	recordcount	loadTimeInSec	IoadTimeInMin
index-2011.cdx	65,525,756	2731.744635	42.5
index-2012.cdx	69,197,158	2884.026416	48.1
index-2013.cdx	80,006,504	3339.867583	55.7
index-2014.cdx	94,445,269	3946.300632	65.8
index-2015.cdx	88,303,363	3683.1083	61.4
index-2016.cdx	93,152,112	3891.03869	64.9
index-2017.cdx	85,003,467	3801.229491	63.4
index-2018.cdx	156,308,196	6476.537223	107.9
index-2019.cdx	164,433,243	6831.800261	113.9
index-2020.cdx	180,379,433	7469.478346	124.5



