

ARCHIVESPARK

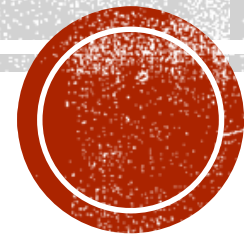
Andrej Galad

12/6/2016

CS-5974 - Independent Study

Virginia Polytechnic Institute and State University, Blacksburg, VA 24061

Professor Edward A. Fox



AGENDA

- ArchiveSpark Overview/Recap
- Benchmarking
- Demo



ARCHIVESPARK

- Framework - efficient data access, extraction and derivation on Web archive data
- [ArchiveSpark: Efficient Web Archive Access, Extraction and Derivation](#)
 - Helge Holzmann, Vinay Goel, Avishek Anand
 - Published in [JCDL 2016](#)
 - Nominated for the Best Paper Award
- OS project - <https://github.com/helgeho/ArchiveSpark>



TWO TECHNIQUES

1. Pre-generated CDX metadata index

- Smaller dataset
- Reduction based on Web archive metadata

2. Incremental filtering workflow

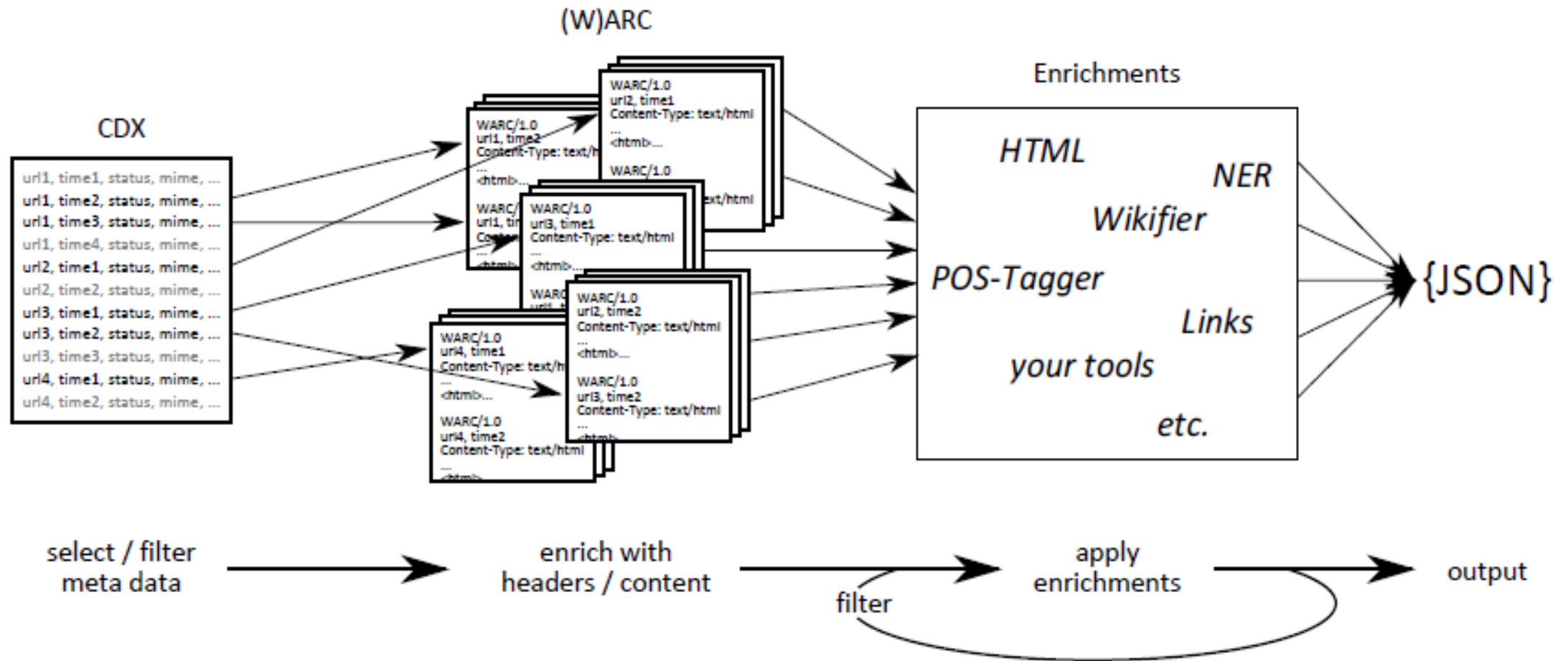
- “Extract only what you need”
- Augment -> Filter -> Repeat

▪ Concept – Enrichments

- ArchiveSpark Record extension
- Featured – StringContent.scala, Html.scala, Json.scala, Entities.scala, Prefix.scala, ...
- Custom – mapEnrich[Source, Target](sourceField, targetField) (f: Source => Target)



WORKFLOW



FLEXIBLE DEPLOYMENT

- Ultimately a Scala/Spark library
- Environments
 - Standalone solitary Spark instance
 - Local HDFS-backed Spark cluster
 - Large-scale YARN/Mesos-orchestrated cluster running Cloudera/Hortonworks
- Quickstart – Docker (latest version)
- ArchiveSpark version – 2.1.0
 - Spark 2.0.2
 - Scala 2.11.7 -> Java 8



WARC FILES

- Standard Web archiving format - [ISO 28500](#)
- Single capture of web resource at a particular time
 - Header section - Metadata (URL, timestamp, content length...)
 - Payload section - Actual response body (HTML, JSON, binary data)
 - HTTP Response - HTTP headers (origin, status code)

```
WARC/1.0
WARC-Type: response
WARC-Record-ID: <urn:uuid:9e6f625c-74f6-4f9b-bec0-cbebc1102f4f>
WARC-Date: 2015-06-13T18:07:56Z
WARC-Target-URI: http://netbootcamp.s3.amazonaws.com/wp-content/uploads/netbootcamp-logo-internet-training.png
WARC-IP-Address: 54.231.65.41
Content-Type: application/http;msgtype=response
WARC-Payload-Digest: sha1:IS672GHFV5GLGOWBVVEN6BLF3DTLS7NB
Content-Length: 18196
WARC-Block-Digest: sha1:LI2ZUQH5PH2UDHGOODXHSGBEHC3UXDG4

HTTP/1.1 200 OK
x-amz-id-2: Bw5mQr6nHrJs+CtzJ0QImShahazsS7hqEgNswjI1jOsZ/+tByVayRWUjL4ggMht
x-amz-request-id: D774CC57D37483F3
Date: Sat, 13 Jun 2015 18:08:35 GMT
Last-Modified: Wed, 24 Dec 2014 14:38:19 GMT
Etag: "1563196cca285cd1aa801a7a1ee91519"
Accept-Ranges: bytes
Content-Type: image/png
Content-Length: 17850
Server: AmazonS3
```



CDX INDEX

- “Reduced” WARC file
 - WARC metadata
 - Pointers to WARC records – offsets in WARC file
- CDX
 - Header – specifies metadata fields contained in the index
 - Body – typically 9 – 11 fields
 - Original URL, SURT, date, filename, MIME type, response code, checksum, redirect, meta tags, compressed offset

```
CDX N b a m s k r M S V g
com,example)/?example=1 20140103030321 http://example.com?example=1 text/html 200 B2LTWWPUOYAH7UIPQ7ZUPQ4VMBSVC36A - - 1043
333 example.warc.gz
com,example)/?example=1 20140103030341 http://example.com?example=1 warc/revisit - B2LTWWPUOYAH7UIPQ7ZUPQ4VMBSVC36A - - 553
1864 example.warc.gz
org,iana)/domains/example 20140128051539 http://www.iana.org/domains/example text/html 302 JZ622UA23G5ZU6Y3XAKH4LINONUEICEG
- - 577 2907 example.warc.gz
```



TOOLS

- CDX Writer
 - Python script for CDX extraction
 - Alternative to Internet Archive's Wayback Machine - <http://archive.org>
- Jupyter Notebook
 - Web application for code sharing/results visualization
- Warcbase (only benchmarking)
 - “State-of-the-art” platform for managing and analyzing Web archives
 - Hadoop/HBase ecosystem - CDH
 - Archive-specific Scala/Java objects for Apache Spark and HBase
 - HBase command-line utilities - IngestFiles



BENCHMARKING

- Evaluation of 3 approaches
 1. ArchiveSpark
 2. Pure Spark using Warcbase library
 3. HBase using Warcbase library
- Preprocessing
 - ArchiveSpark - CDX index files extraction
 - HBase - WARC ingestion
- ArchiveSpark Benchmark subproject
 - Requirements:
 - Built and included Warcbase
 - `sbt assemblyPackageDependency -> sbt assembly`



ENVIRONMENT

- Development
 - Cloudera Quickstart VM - CDH 5.8.2
- Benchmarking
 - Cloudera CDH 5.8.2 cluster hosted on AWS (courtesy of Dr. Zhiwu Xie)
 - 5-node cluster consisting of m4.xlarge AWS EC2 instances
 - 4 vCPUs
 - 16 GiB RAM
 - 30 GB EBS storage
 - 750 Mbps network

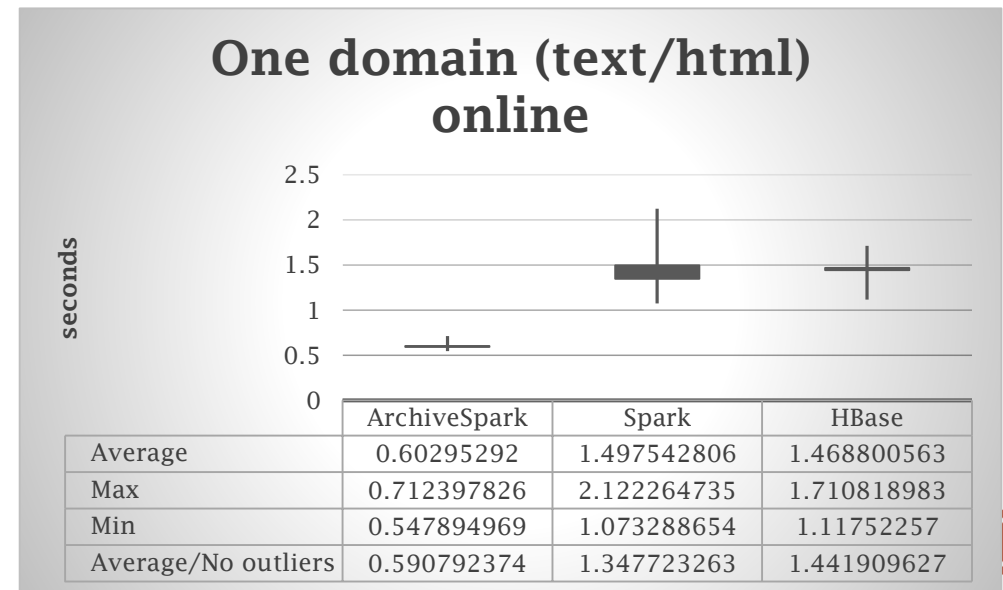
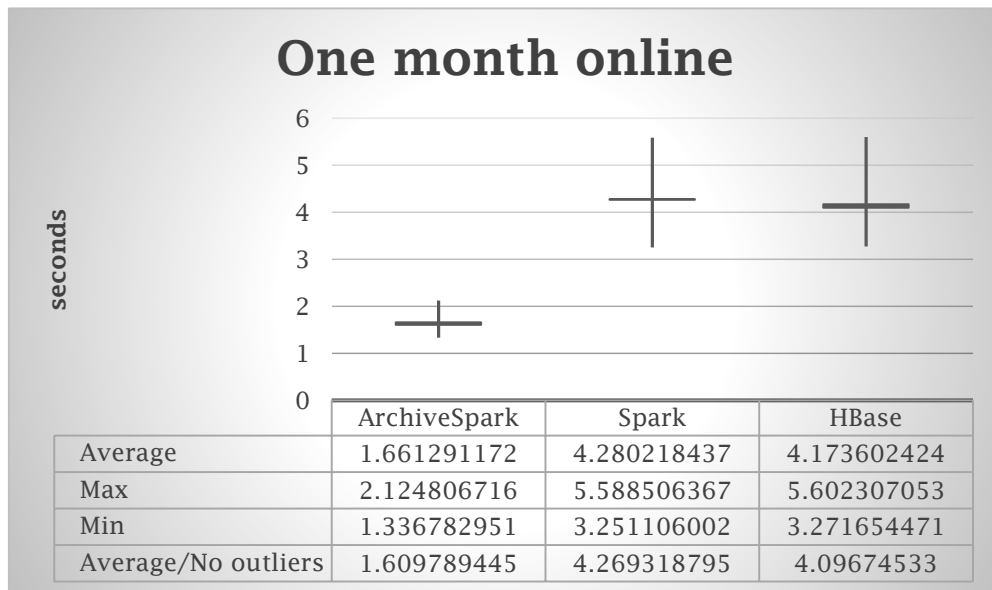
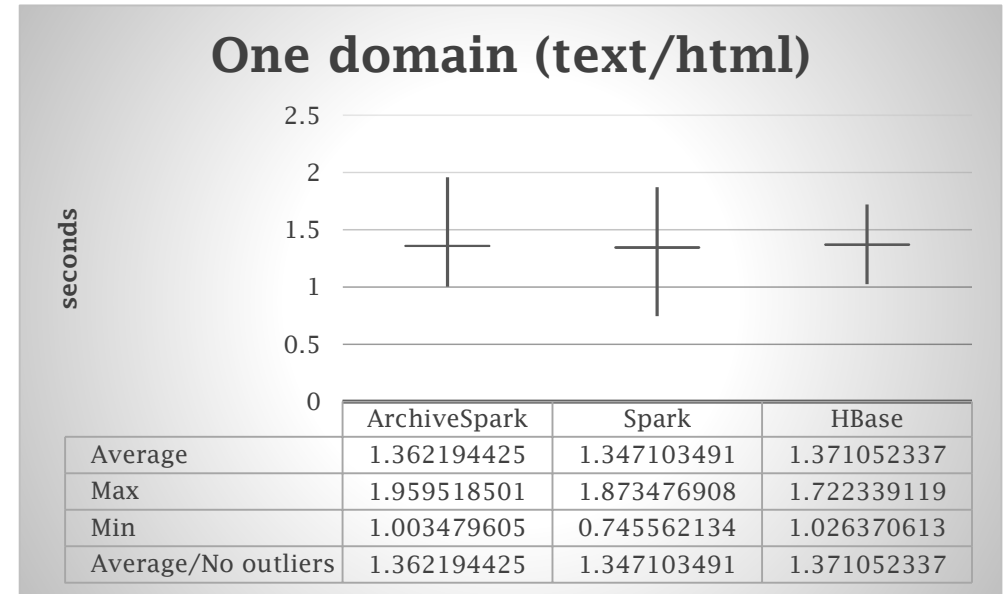
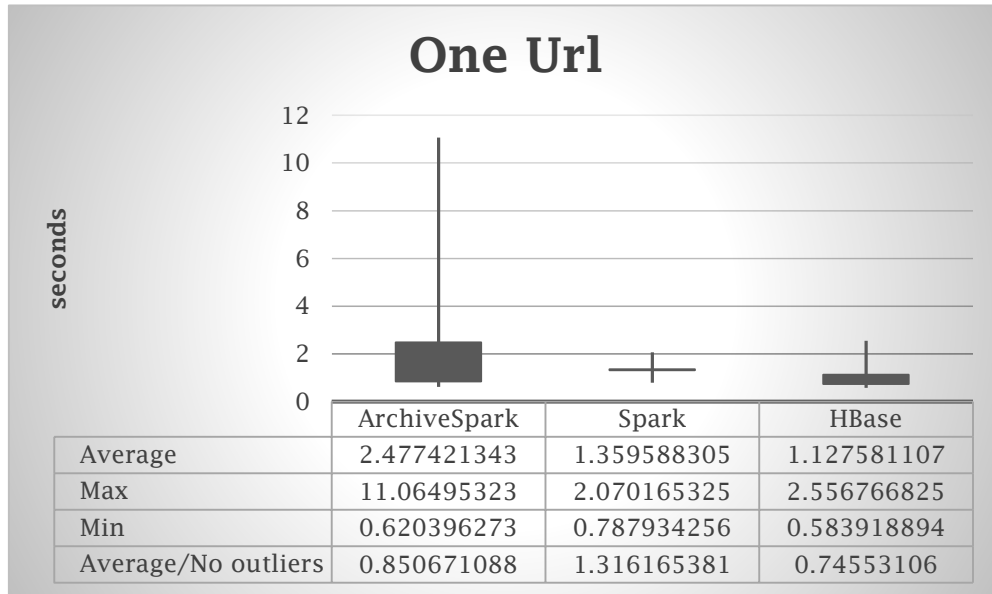


BENCHMARK 1 - SMALL SCALE

- Filtering & corpus extraction
- 4 scenarios
 - Filtering of the dataset for a specific URL (one URL benchmark)
 - Filtering of the dataset for a specific domain (one domain benchmark)
 - Filtering of the dataset for a date range of records (one month benchmark)
 - Filtering of the dataset for a specific active (200 OK) domain (one active domain benchmark)
- Dataset - `example.warc.gz`
 - One capture of archive.it domain
 - 261 records
 - 2.49 MB



BENCHMARK 1 - RESULTS



BENCHMARK 2 - MEDIUM SCALE

- Filtering & corpus extraction
- 4 scenarios
 - Filtering of the dataset for a specific url (one url benchmark)
 - Filtering of the dataset for a specific domain (one domain benchmark)
 - Filtering of the dataset for a specific active (200 OK) domain (one active domain benchmark)
 - Filtering of the dataset for pages containing scripts (pages with scripts benchmark)
- Dataset - WIDE collection
 - Internet Archive crawl data from Webwide Crawl (02/25/2011)
 - 214470 records
 - 9064 MB
 - 9 files - approx. 1 GB



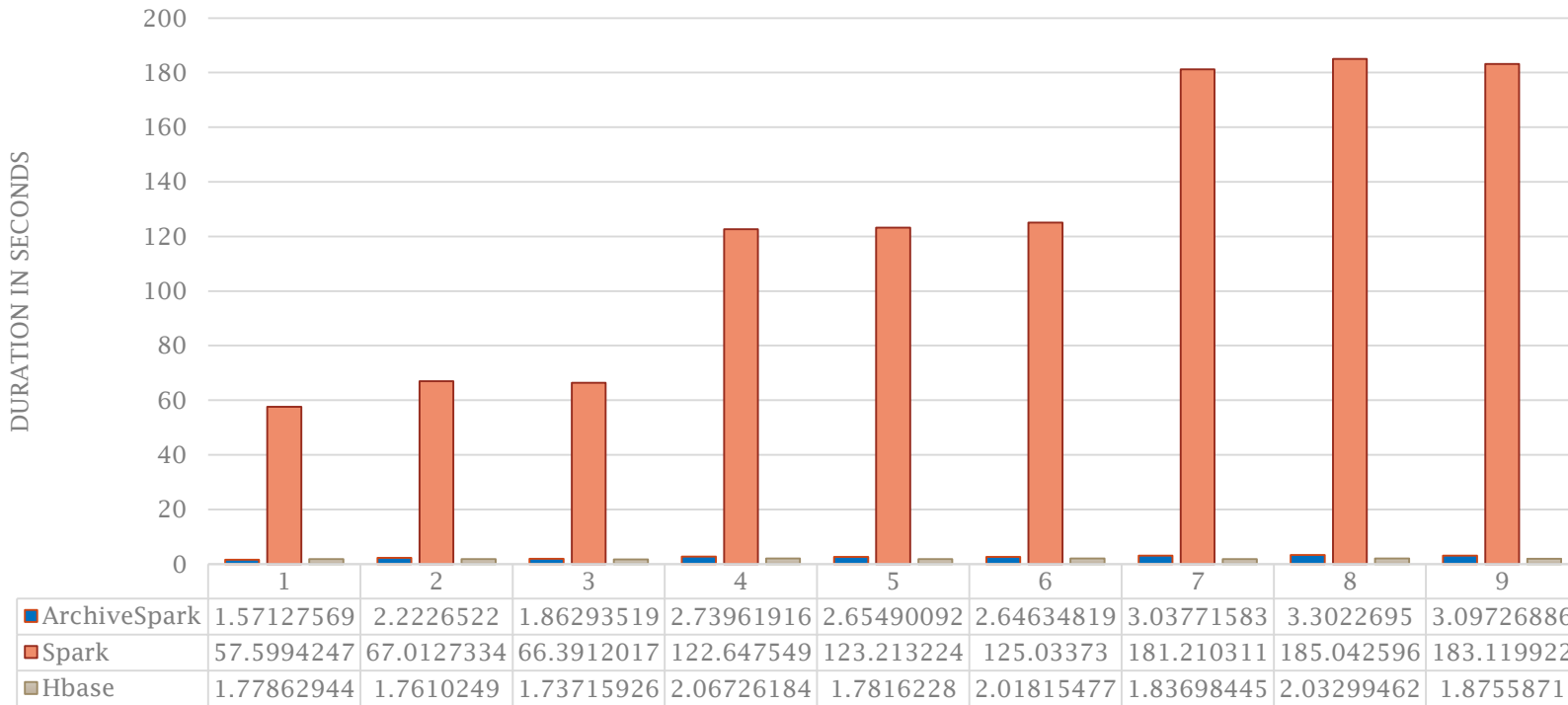
BENCHMARK 2 - PREPROCESSING

- CDX Extraction
 - 4 minutes 41 seconds
- HDFS Upload
 - 2 minutes 46 seconds
- HBase Ingestion x9
 - 1 file - (1 minute 10 seconds <-> 1 minute 32 seconds)
 - Sequential ingestion - approx. 13 minutes 54 seconds



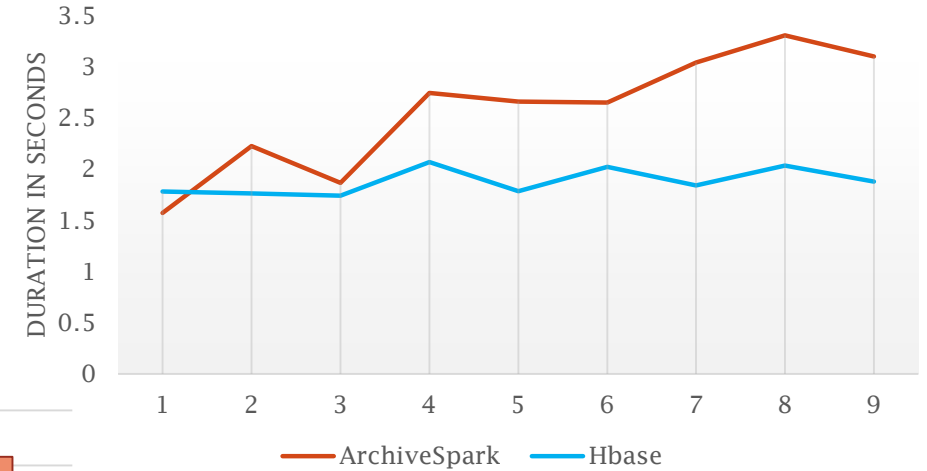
BENCHMARK 2 - RESULTS

One URL



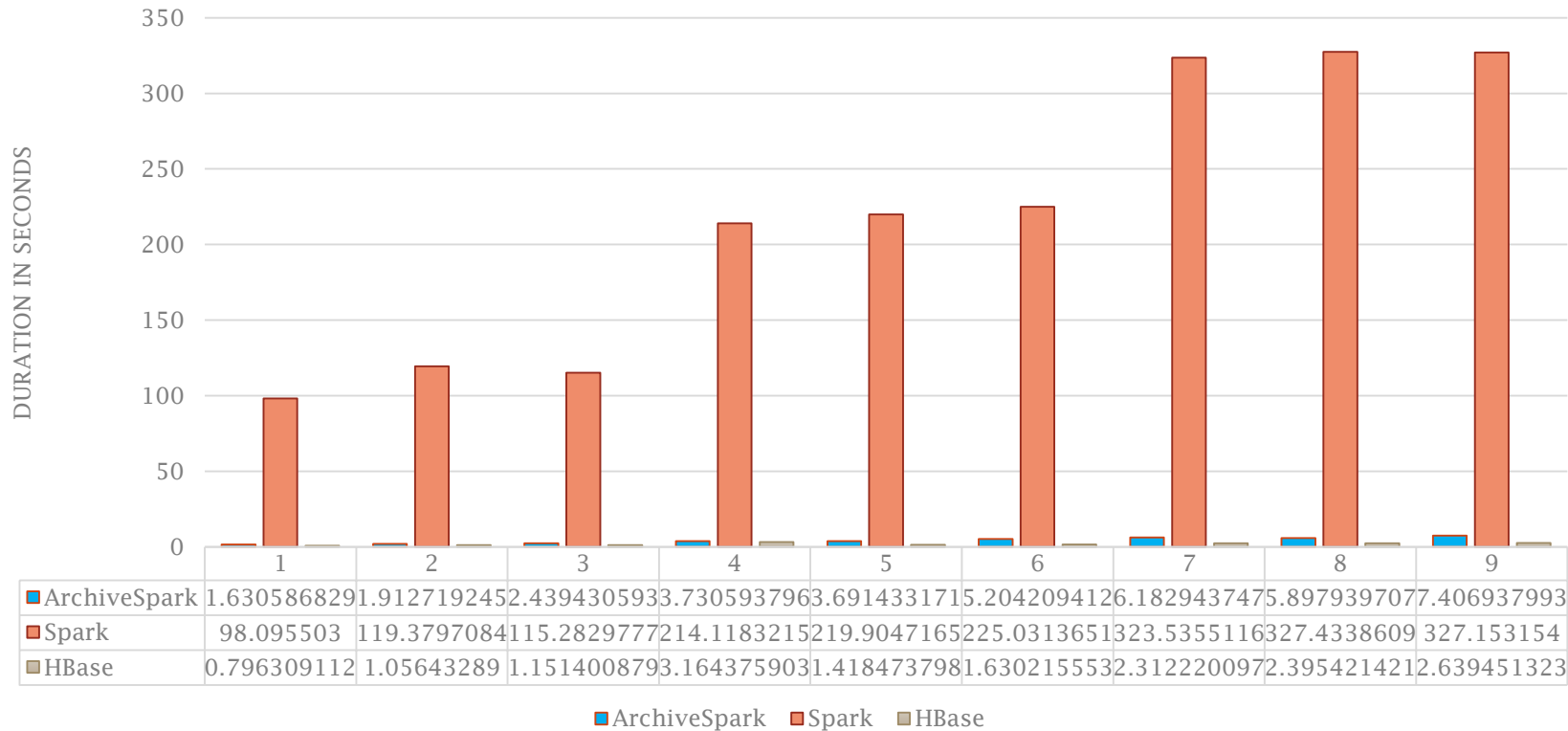
ArchiveSpark Spark Hbase

One URL

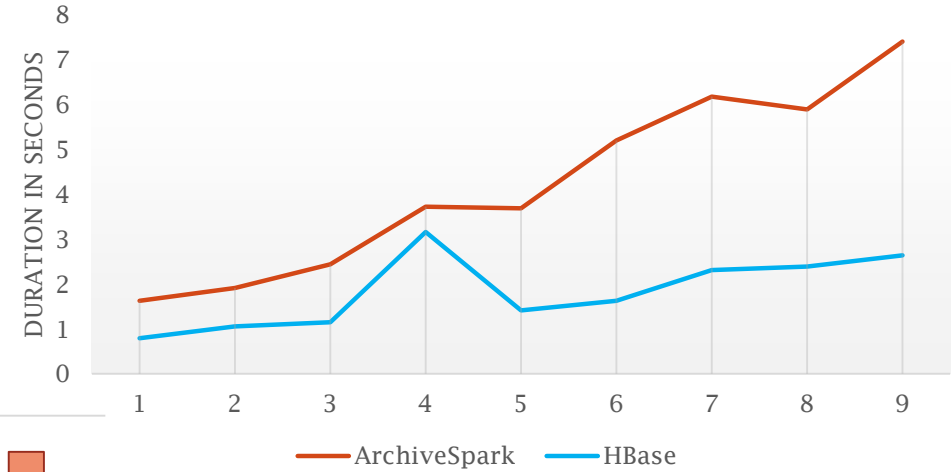


BENCHMARK 2 - RESULTS

One Domain (text/html)

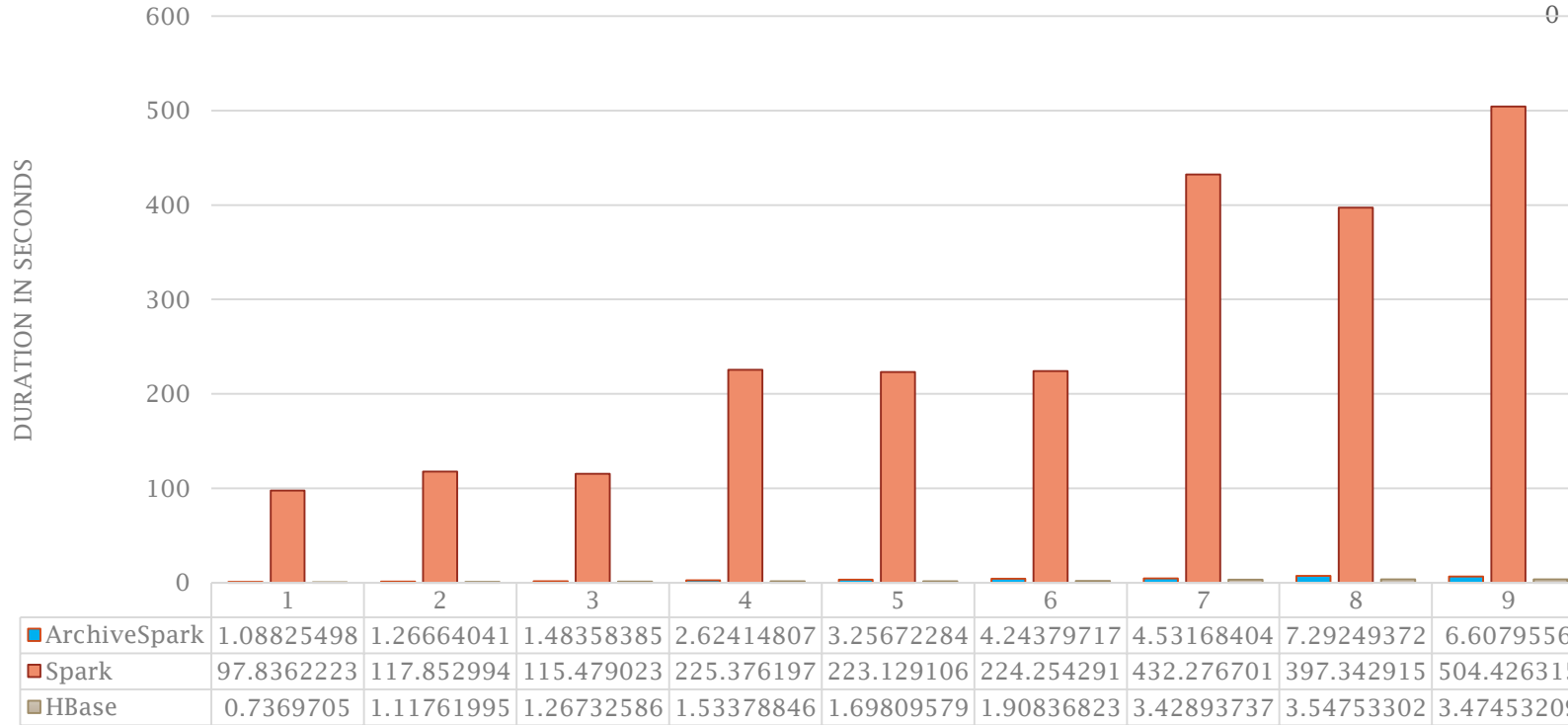


One Domain (text/html)



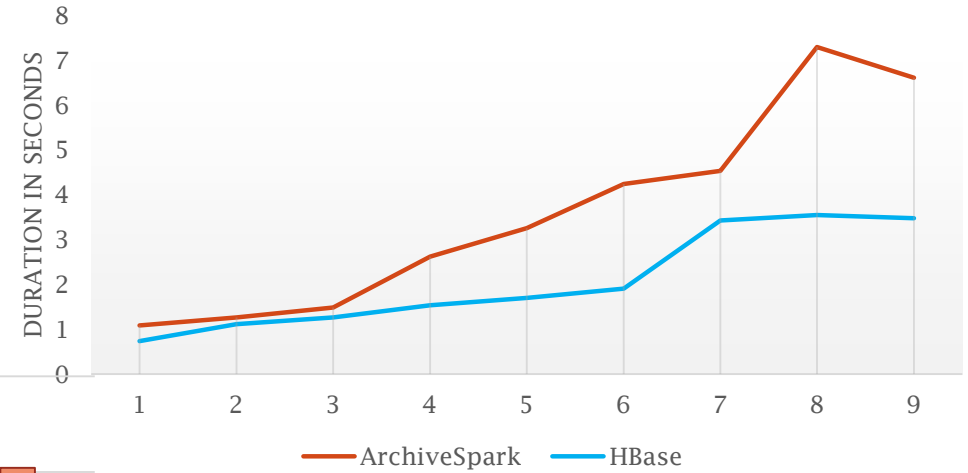
BENCHMARK 2 - RESULTS

One Domain (text/html) Online (Status Code - 200)



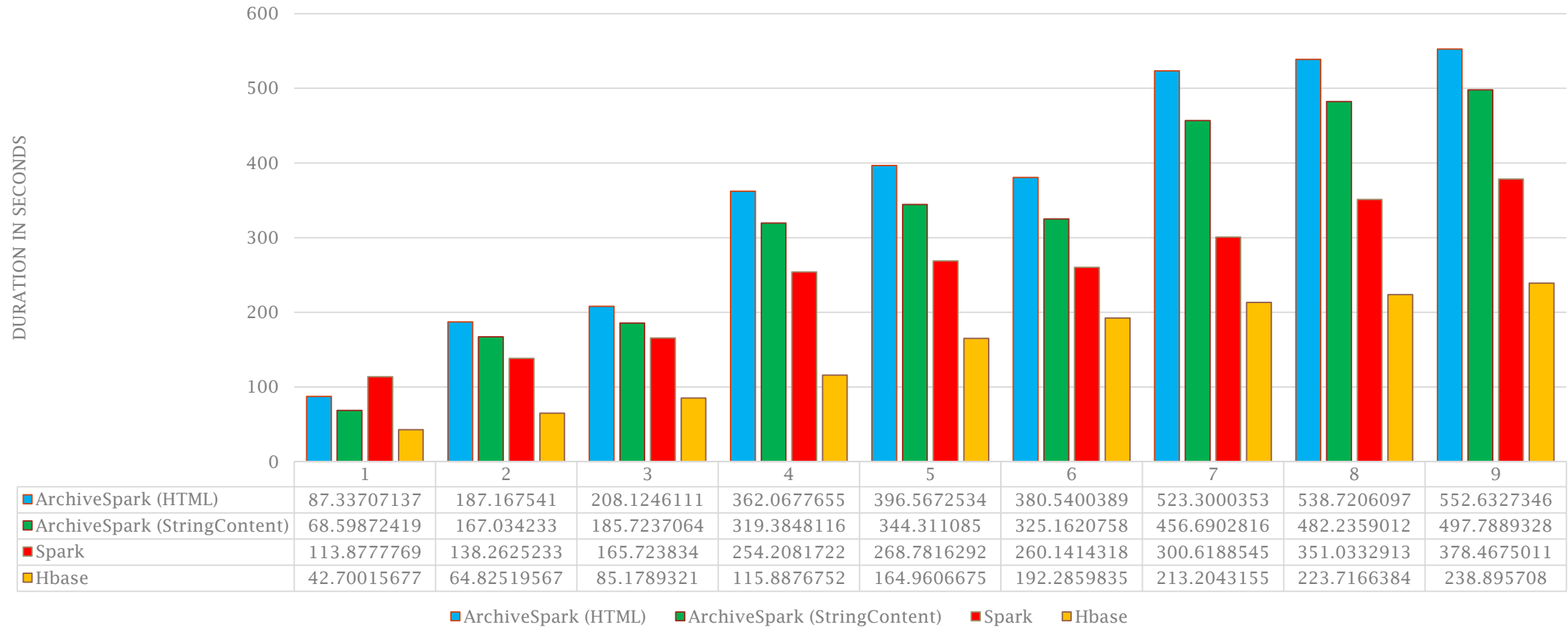
ArchiveSpark Spark HBase

One Domain (text/html) Online



BENCHMARK 2 - RESULTS

Web Pages (text/html) with Scripts



ACKNOWLEDGEMENT & DISCLAIMER

- This material is based upon work supported by following grants:
 - IMLS LG-71-16-0037-16: Developing Library Cyberinfrastructure Strategy for Big Data Sharing and Reuse
 - NSF IIS-1619028, III: Small: Collaborative Research: Global Event and Trend Archive Research (GETAR)
 - NSF IIS - 1319578: III: Small: Integrated Digital Event Archiving and Library (IDEAL)
- Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation.



THANK YOU



<http://tinyurl.com/zejgc9f>



TUTORIAL

