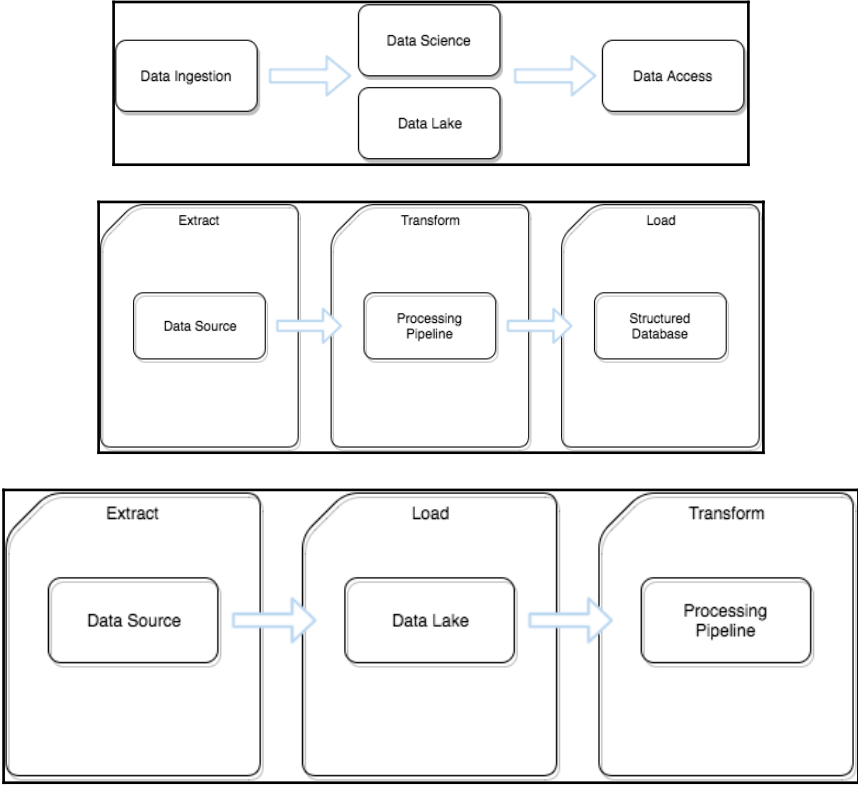
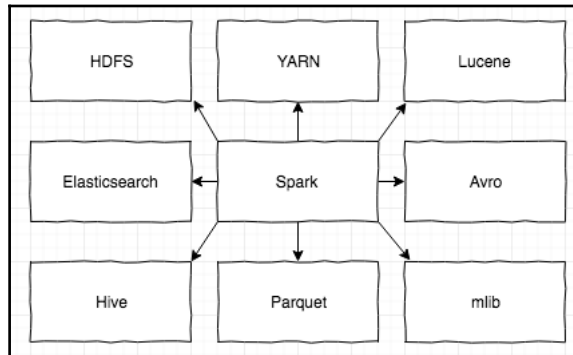
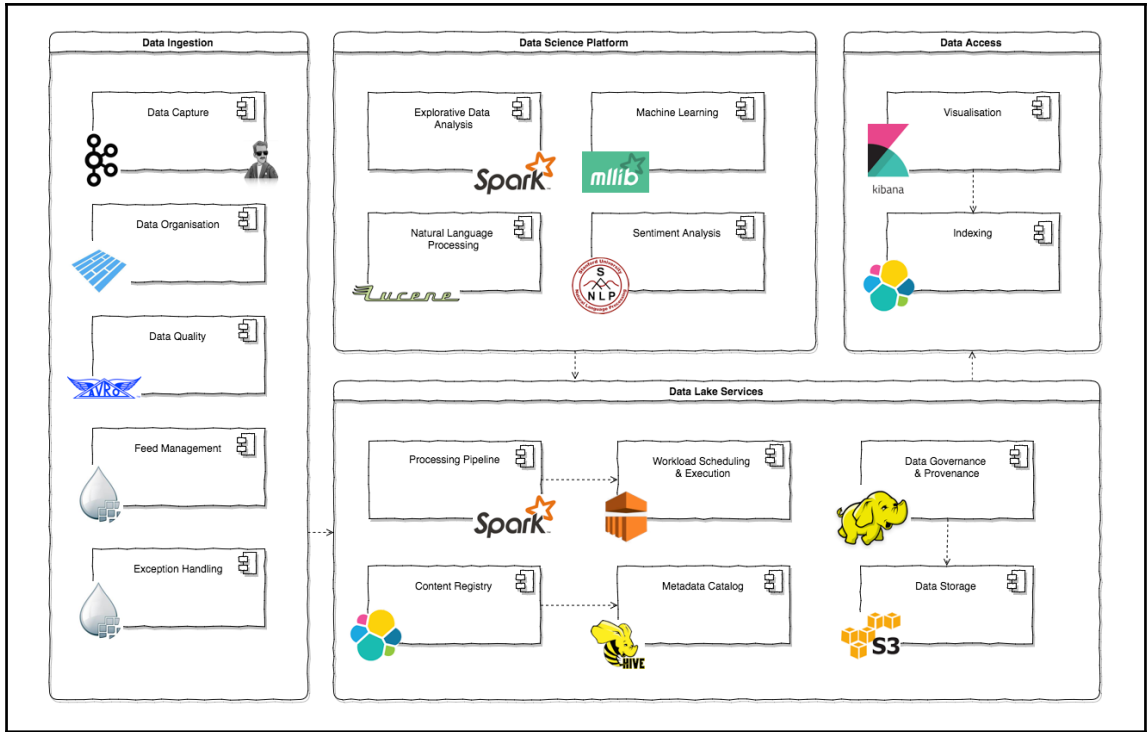
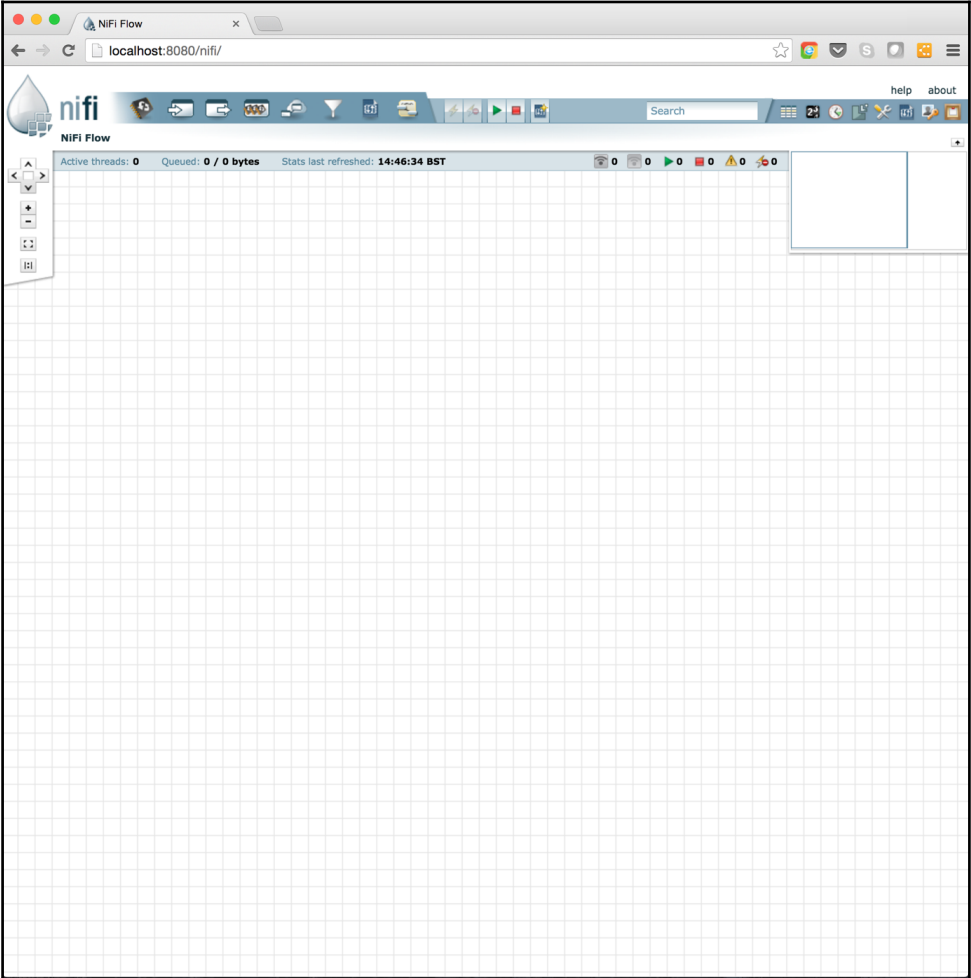


Chapter 01: The Big Data Science Ecosystem





Chapter 02: Data Acquisition






nifi


NiFi Flow

Active threads: 0 Queued: 0 / 0 bytes Stats last refreshed: 15:25:19 BST



Fetch GDELT Latest File List 
GetHTTP

In	0 / 0 bytes	(5 min)
Read/Write	0 bytes / 0 bytes	(5 min)
Out	0 / 0 bytes	(5 min)
Tasks/Time	0 / 00:00:00.000	(5 min)


Search

NiFi Flow
 Active threads: 0 Queued: 0 / 0 bytes Stats last refreshed: 15:25:50 BST

⚠
Fetch GDELT Latest File List
 GetHTTP

In: 0 / 0 bytes (5 min)
 Read/Write: 0 by
 Out: 0 / 0
 Tasks/Time: 0 / 0

Configure Processor

Settings | Scheduling | **Properties** | Comments

Required field + New property

Property	Value
URL	http://data.gdeltproject.org/gdeltv2/lastupdate.txt
Filename	update\${UUID()}.csv
SSL Context Service	No value set
Username	No value set
Password	No value set
Connection Timeout	30 sec
Data Timeout	30 sec
User Agent	No value set
Accept Content-Type	No value set
Follow Redirects	false
Proxy Host	No value set

```

data.gdeltproject.org/gdeltv2/lastupdate.txt

110399 9fbd87b04d55ee2eac90908de6369ba9 http://data.gdeltproject.org/gdeltv2/20160730143000.export.CSV.zip
212198 7fd3068395b0c819a21887acd0d39ffe http://data.gdeltproject.org/gdeltv2/20160730143000.mentions.CSV.zip
8254196 75810140bac401461de7a5baecc519a3 http://data.gdeltproject.org/gdeltv2/20160730143000.gkg.csv.zip
  
```

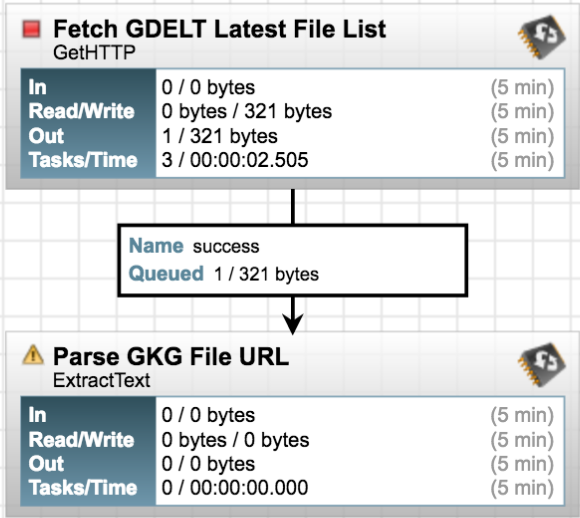


nifi

NiFi Flow



Active threads: 0 Queued: 1 / 321 bytes Stats last refreshed: 04:22



nifi

NIFI Flow

Active threads: 0 Queued: 0 / 0 bytes Stats last refreshed: 16:48:15 BST

Fetch GDELT Latest File List
GetHTTP

In: 0 / 0 bytes (5 min)
Read/Write: 0 bytes / 0 bytes (5 min)
Out: 0 / 0 bytes (5 min)
Taska/Time: 0 / 00:00:00.000 (5 min)

Name success
Queued: 0 / 0 bytes

Parse GKG File URL
ExtractText

In: 0 / 0 bytes (5 min)
Read/Write: 0 bytes / 0 bytes (5 min)
Out: 0 / 0 bytes (5 min)
Taska/Time: 0 / 00:00:00.000 (5 min)

Configure Processor

Settings Scheduling **Properties** Comments

Required field + New property

Property	Value
Enable Canonical Equivalence	<input checked="" type="checkbox"/> true
Enable Case-insensitive Matching	<input type="checkbox"/> false
Permit Whitespace and Comments in Pattern	<input type="checkbox"/> false
Enable DOTALL Mode	<input type="checkbox"/> false
Enable Literal Parsing of the Pattern	<input type="checkbox"/> false
Enable Multiline Mode	<input type="checkbox"/> false
Enable Unicode-aware Case Folding	<input type="checkbox"/> false
Enable Unicode Predefined Character Classes	<input type="checkbox"/> false
Enable Unix Lines Mode	<input type="checkbox"/> false
Include Capture Group 0	<input type="checkbox"/> true
url	<input type="text" value="([^\s]*gkg.csv.*)"/> ✕

nifi

NIFI Flow

Active threads: 0 Queued: 0 / 0 bytes Stats last refreshed: 16:49:15 BST

```

    graph TD
      A[Fetch GDELT Latest File List] -- "Name success" --> B[Parse GKG File URL]
      B -- "Name matched" --> C[Fetch GKG File From URL]
  
```

Fetch GDELT Latest File List
GetHTTP
In: 0 / 0 bytes (5 min)
Read/Write: 0 bytes / 0 bytes (5 min)
Out: 0 / 0 bytes (5 min)
Task(s)/Time: 0 / 00:00:00.000 (5 min)

Parse GKG File URL
ExtractText
In: 0 / 0 bytes (5 min)
Read/Write: 0 bytes / 0 bytes (5 min)
Out: 0 / 0 bytes (5 min)
Task(s)/Time: 0 / 00:00:00.000 (5 min)

Fetch GKG File From URL
InvokeHTTP
In: 0 / 0 bytes (5 min)
Read/Write: 0 bytes / 0 bytes (5 min)
Out: 0 / 0 bytes (5 min)
Task(s)/Time: 0 / 00:00:00.000 (5 min)

Configure Processor

Settings Scheduling **Properties** Comments

Required field + New property

Property	Value
HTTP Method	GET
Remote URL	\${url}
SSL Context Service	No value set
Connection Timeout	5 secs
Read Timeout	15 secs
Include Date Header	True
Follow Redirects	True
Attributes to Send	No value set
Basic Authentication Username	No value set
Basic Authentication Password	No value set
Proxy Host	No value set

Cancel Apply

NIIFI Flow

localhost:8080/nifi/

NIIFI Flow

Active threads: 1 Queued: 0 / 0 bytes Stats last refreshed: 18:55:11 BST

```
graph TD; A[Fetch GDELT Latest File List] -- "Name success" --> B[Parse GKG File URL]; B -- "Name matched" --> C[Fetch GKG File From URL]; C -- "Name Response, Retry" --> D[Unzip GKG Content]; D -- "Name success" --> E[Save To HDFS];
```

Fetch GDELT Latest File List
GetHTTP
In: 0 / 0 bytes (5 min)
Read/Write: 0 bytes / 320 bytes (5 min)
Out: 1 / 320 bytes (5 min)
Tasks/Time: 66 / 00:00:12.664 (5 min)

Name success
Queued: 0 / 0 bytes

Parse GKG File URL
ExtractText
In: 1 / 320 bytes (5 min)
Read/Write: 320 bytes / 0 bytes (5 min)
Out: 1 / 320 bytes (5 min)
Tasks/Time: 1 / 00:00:00.001 (5 min)

Name matched
Queued: 0 / 0 bytes

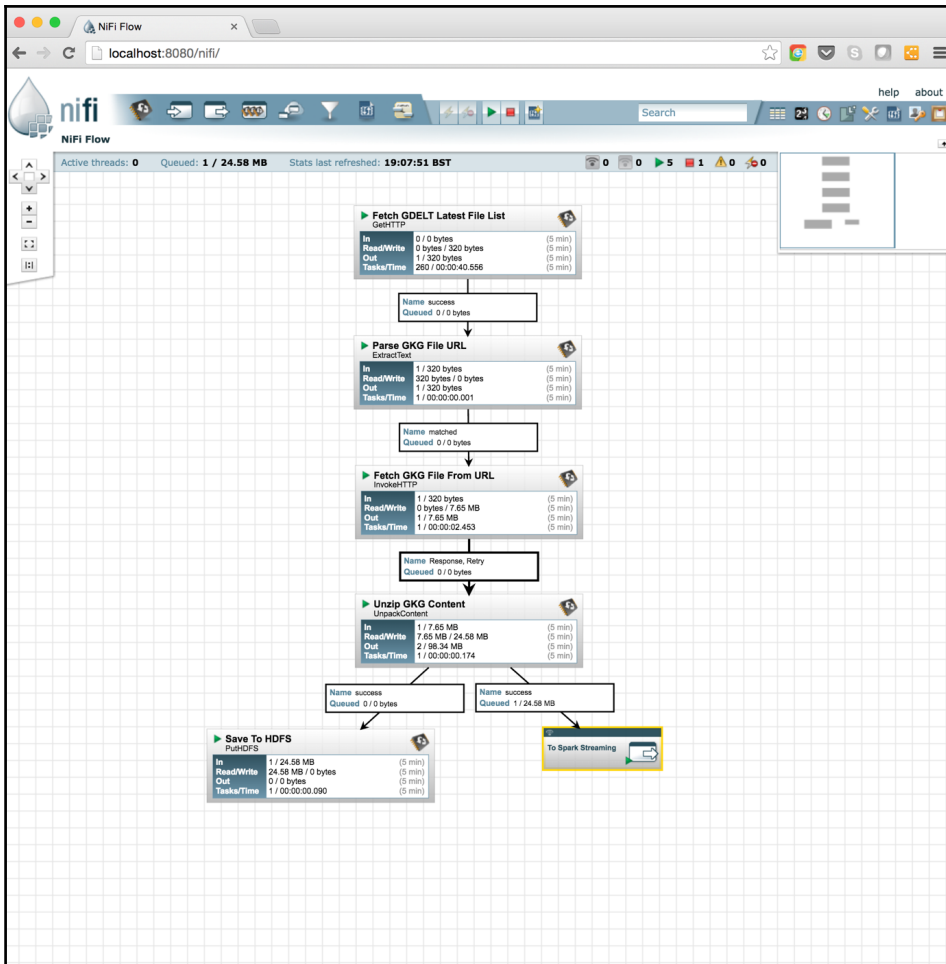
Fetch GKG File From URL
InvokeHTTP
In: 1 / 320 bytes (5 min)
Read/Write: 0 bytes / 8.38 MB (5 min)
Out: 1 / 8.38 MB (5 min)
Tasks/Time: 1 / 00:00:01.597 (5 min)

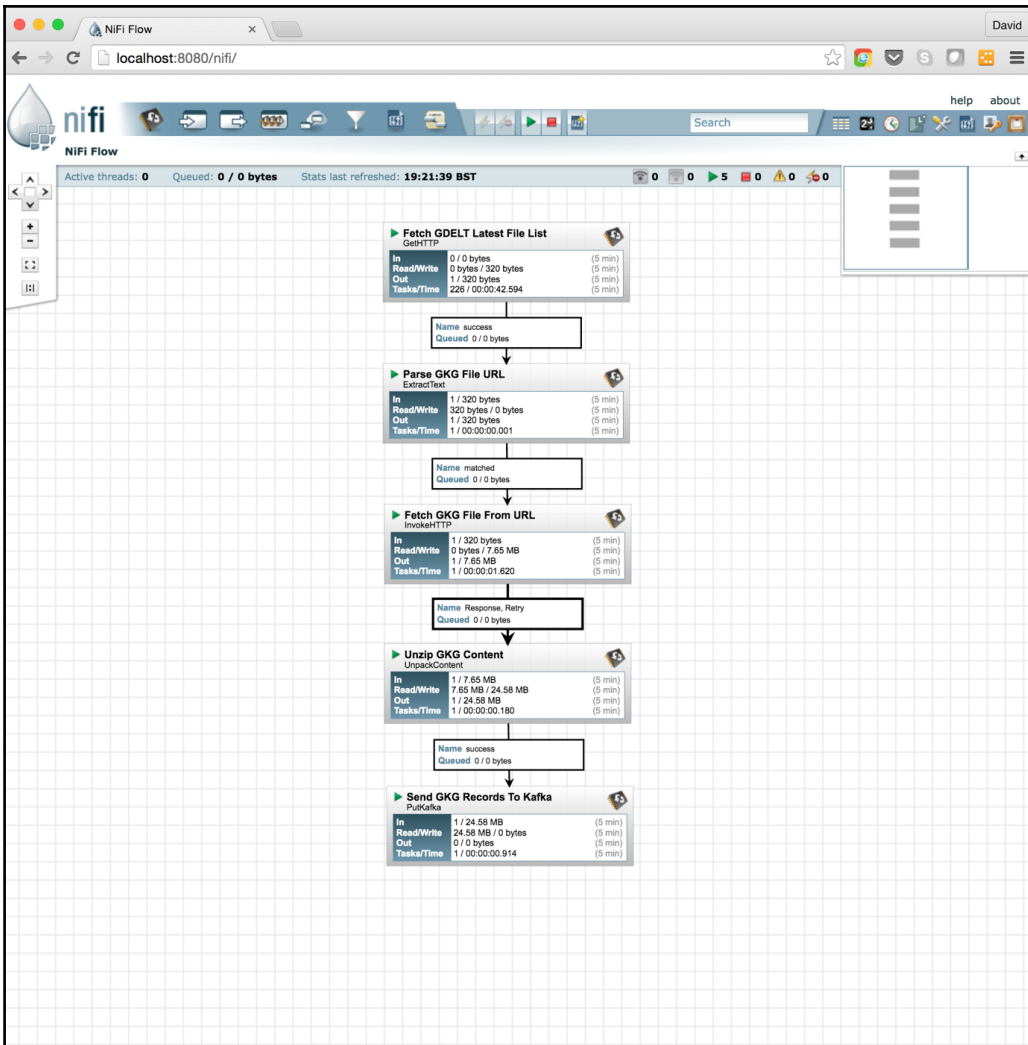
Name Response, Retry
Queued: 0 / 0 bytes

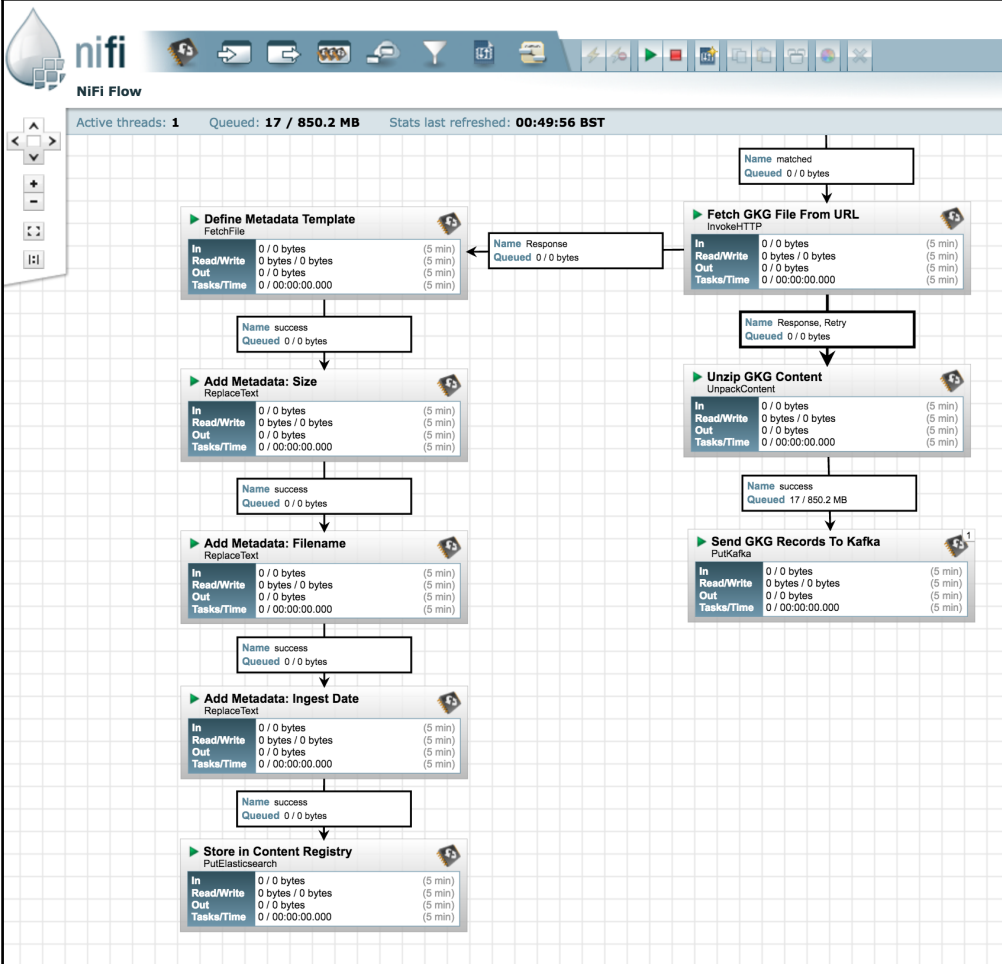
Unzip GKG Content
UnpackContent
In: 1 / 8.38 MB (5 min)
Read/Write: 8.38 MB / 26.77 MB (5 min)
Out: 1 / 26.77 MB (5 min)
Tasks/Time: 1 / 00:00:00.205 (5 min)

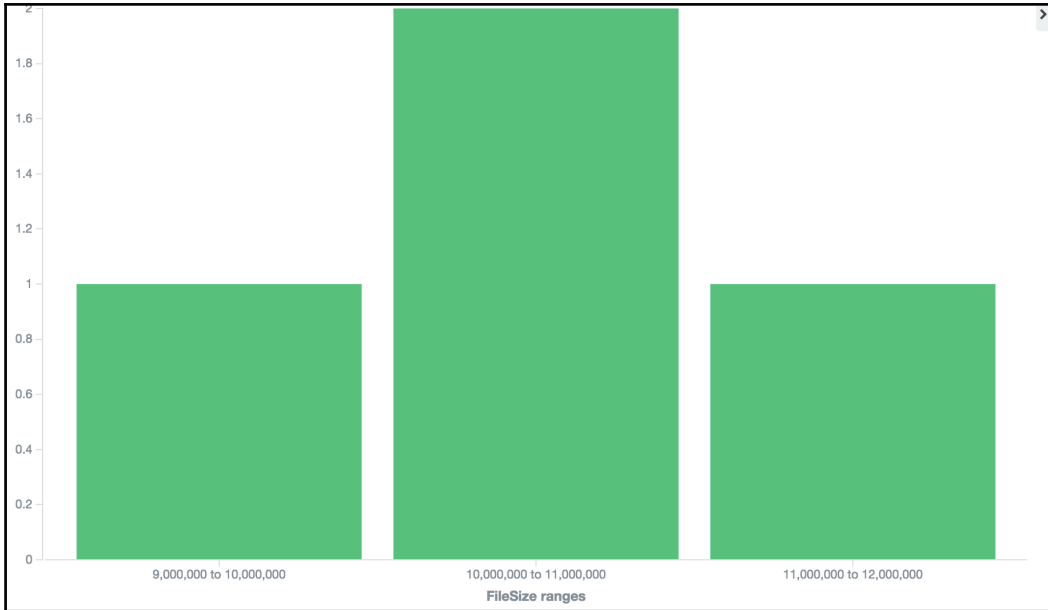
Name success
Queued: 0 / 0 bytes

Save To HDFS
PutHDFS
In: 1 / 26.77 MB (5 min)
Read/Write: 26.77 MB / 0 bytes (5 min)
Out: 0 / 0 bytes (5 min)
Tasks/Time: 1 / 00:00:00.277 (5 min)





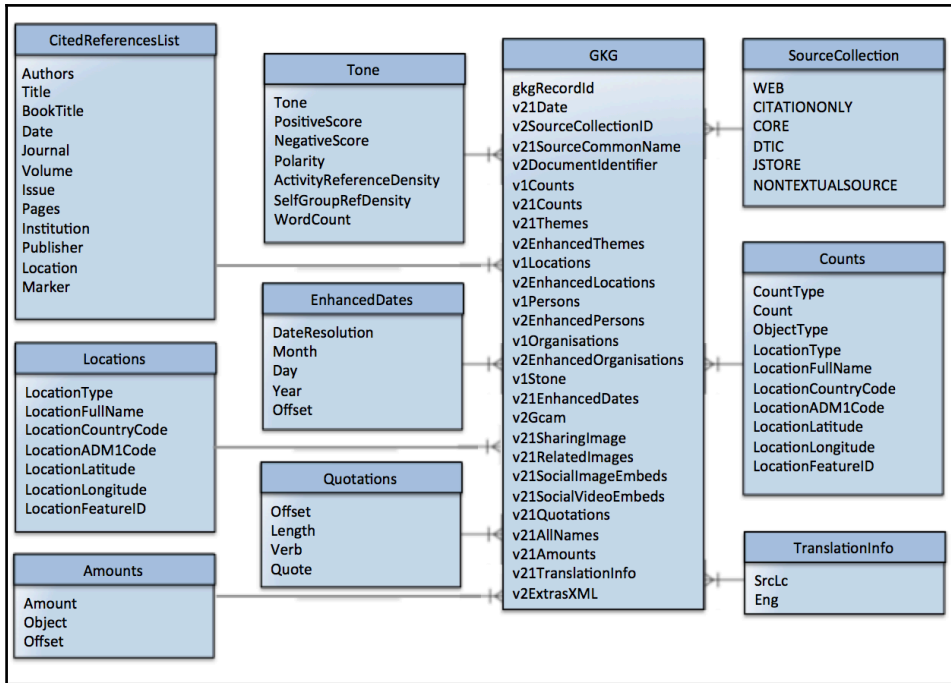


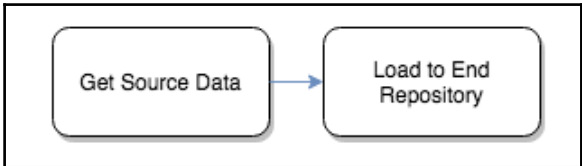
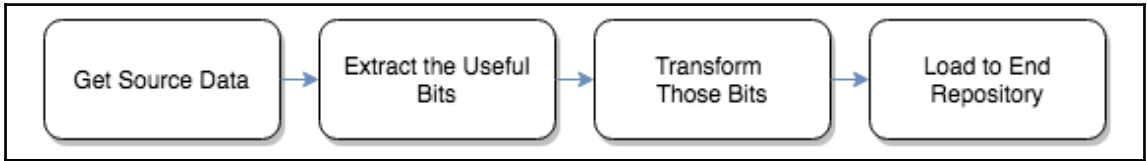
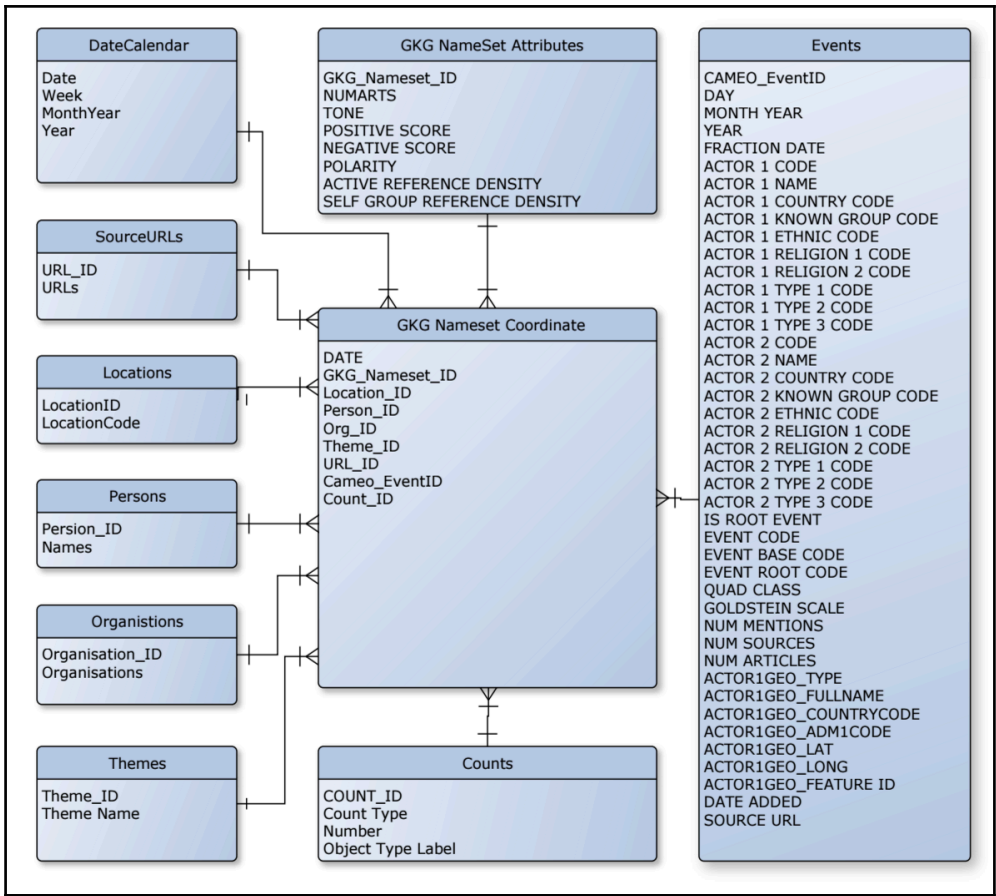


Chapter 03: Input Formats and Schema

GKG
gkgRecordId
v21Date
v2SourceCollectionIdentifier
v21SourceCommonName
v2DocumentIdentifier
v1Counts
v21Counts
v21Themes
v2EnhancedThemes
v1Locations
v2EnhancedLocations
v1Persons
v2EnhancedPersons
v1Organisations
v2EnhancedOrganisations
v1Stone
v21EnhancedDates
v2Gcam
v21SharingImage
v21RelatedImages
v21SocialImageEmbeds
v21SocialVideoEmbeds
v21Quotations
v21AllNames
v21Amounts
v21TranslationInfo
v2ExtrasXML

Type	Count
<u>WordCount</u>	125
General Inquirer <u>Bodypt</u>	4
<u>SentiWordNet</u>	40
<u>SentiWordNet</u> average	3.21111111





Chapter 04: Exploratory Data Analysis

Configure the Profiler's Mask FINISHED

YourMask: STRING HIGH

YourMask: String = ASCIICLASS_HIGHGRAIN

Took 34 seconds

Configure Delimiter FINISHED

YourDelimiter: \t

YourDelimiter: String = \t

Took 28 seconds

Configure the CSV file to profile FINISHED

YourFilePath: /user/feeds/gdelt/evven

YourFilePath: String = /user/feeds/gdelt/events/*.export.CSV

Took 1 seconds

Headers? FINISHED

YourHeaders: Has Header No Header

YourHeader: String = true

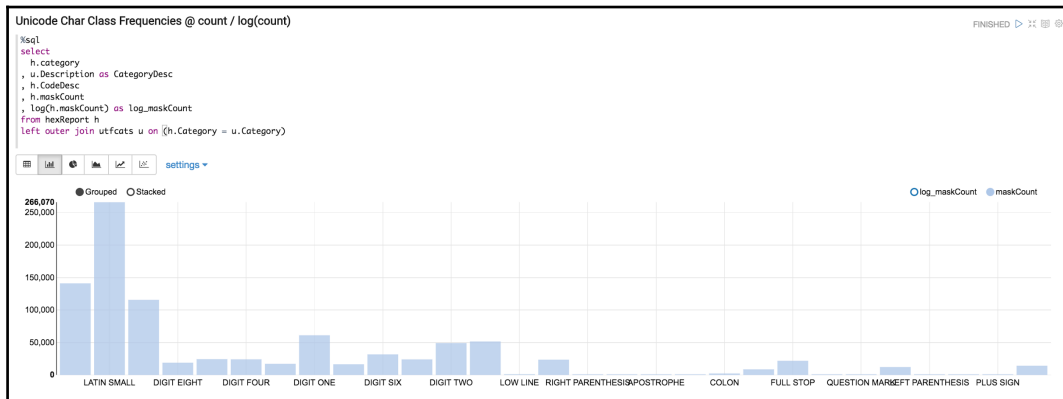
Took 0 seconds

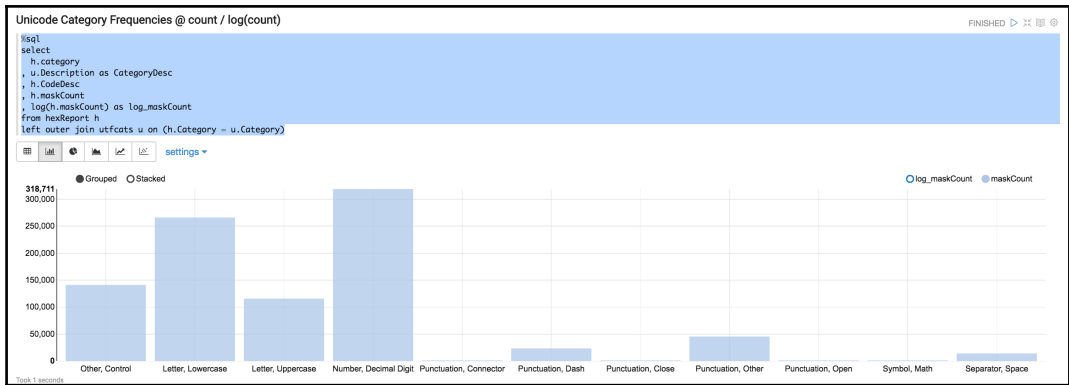
Inspect File as Tabled Data FINISHED

```
%sql
select * from RowData
limit 10
```

ime	ActionGeo_CountryCode	ActionGeo_ADM1Code	ActionGeo_ADM2Code	ActionGeo_Lat	ActionGeo_Long	ActionGeo_FeatureID	DateAdded	SourceURL
dia	IN	IN07	17,911	28.6	77.2	-2,106,102	20,161,220,124,500	https://www.zigwheels.com/news-features/news/toyota-fortuner-sales-cross-one-lakh-milestone-in-india/27107/
	JA	JA		36	138	JA	20,161,220,124,500	http://www.asahi.com/ajw/articles/AJ201612200061.html
	CA	CA02	12,552	49.1	-122.65	-567,692	20,161,220,124,500	http://www.scienceworldreport.com/articles/55092/20161220/controversial-files-secret-document-reveals-alien-ufo-existence-details-biology.htm
an,	PK	PK03	40,341	32.8109	70.7154	84,485	20,161,220,124,500	http://www.whio.com/news/world/women-vienna-get-pocket-alarms-for-new-year-eve?JaoPwK6ZG3tVXsJeltTV
r,	GM	GM		51	9	GM	20,161,220,124,500	http://www.pravdareport.com/society/stories/20-12-2016/136464-germany_migrants-0/
n,	GH	GH03	190,560	7.75	-1.5	-2,071,502	20,161,220,124,500	http://www.ghanaweb.com/GhanaHomePage/NewsArchive/Mahama-gun-...-to-be-released-by-Yeas-Fa...-49E99
..								

Took 0 seconds

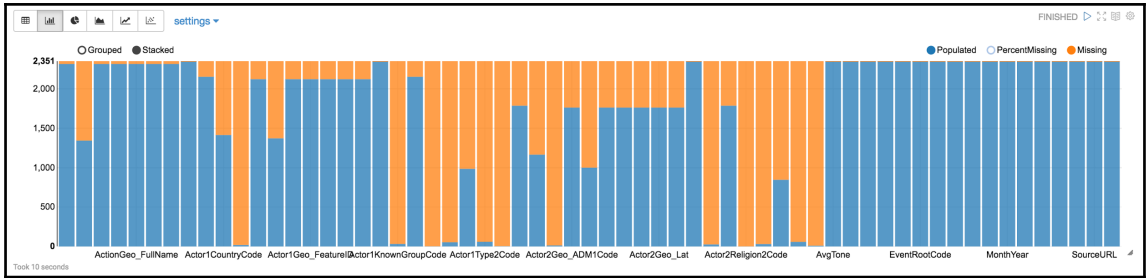




```
Metrics_POPCHECKS.toDF.show(1000, false)
```

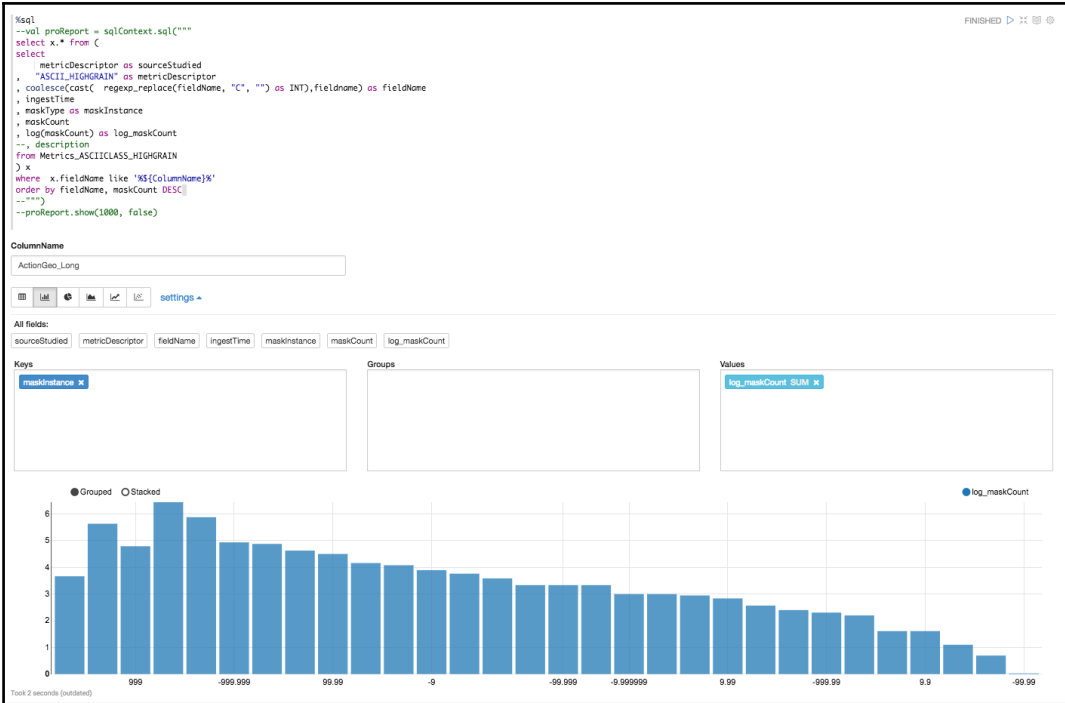
metricDescriptor	sourceStudied	ingestTime	maskType	fieldName	occurrenceCount	keyCount	maskCount	description
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		Actor1Type2Code	1100525	161	159	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		Actor1Geo_Lat	1100525	161	12122	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		Actor2Type2Code	1100525	161	158	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		NumSources	1100525	161	12351	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		ActionGeo_Long	1100525	161	12315	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		ActionGeo_FullName	1100525	161	12122	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		AvgTone	1100525	161	12351	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		Actor1Geo_Type	1100525	161	12351	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		ActionGeo_ADM2Code	1100525	161	11343	<function>
/user/feeds/gdelt/events/*.export.CSV		2016-12-23 10		DateAdded	1100525	161	12351	<function>

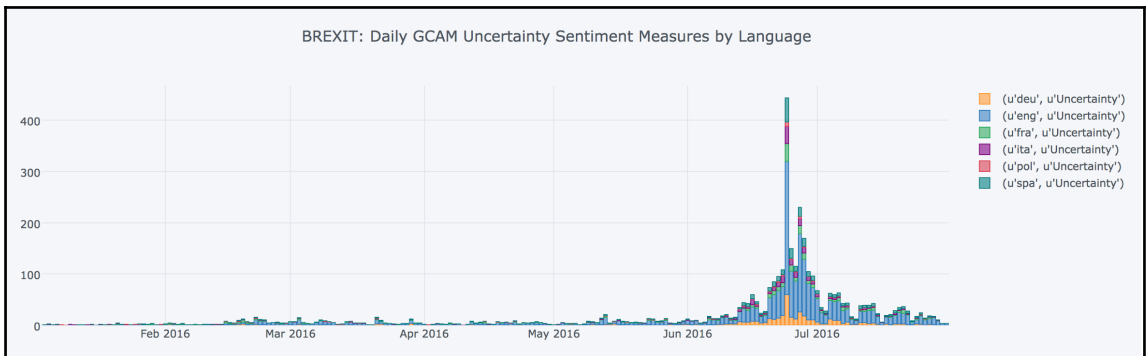
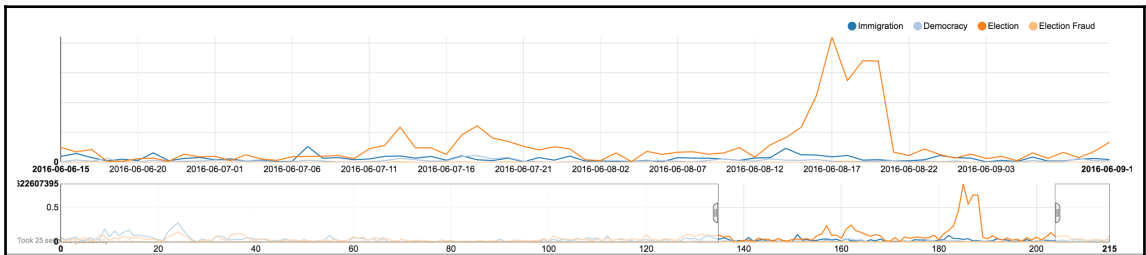
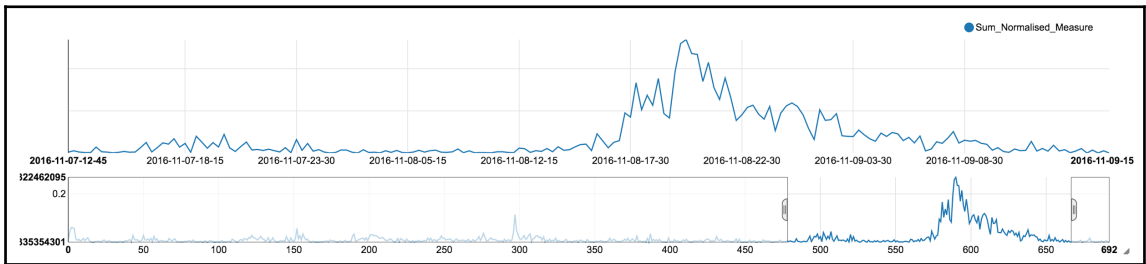
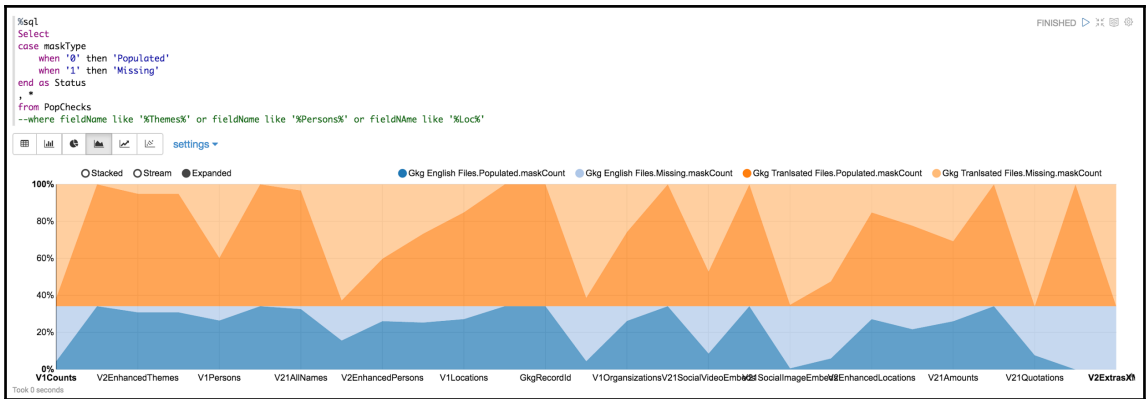
fieldname	Missing	Populated	PercentMissing	fileName
ActionGeo_ADM1Code	136	12315	1.53	/user/feeds/gdelt/events/*.export.CSV
ActionGeo_ADM2Code	1008	11343	42.88	/user/feeds/gdelt/events/*.export.CSV
ActionGeo_CountryCode	136	12315	1.53	/user/feeds/gdelt/events/*.export.CSV
ActionGeo_FeatureID	136	12315	1.53	/user/feeds/gdelt/events/*.export.CSV
ActionGeo_FullName	136	12315	1.53	/user/feeds/gdelt/events/*.export.CSV
ActionGeo_Lat	136	12315	1.53	/user/feeds/gdelt/events/*.export.CSV
ActionGeo_Long	136	12315	1.53	/user/feeds/gdelt/events/*.export.CSV
ActionGeo_Type	0	12351	0.0	/user/feeds/gdelt/events/*.export.CSV
Actor1Code	198	12153	18.42	/user/feeds/gdelt/events/*.export.CSV
Actor1CountryCode	1937	11414	39.86	/user/feeds/gdelt/events/*.export.CSV
Actor1EthnicCode	12332	19	99.19	/user/feeds/gdelt/events/*.export.CSV
Actor1Geo_ADM1Code	1229	12122	9.74	/user/feeds/gdelt/events/*.export.CSV
Actor1Geo_ADM2Code	1979	11372	41.64	/user/feeds/gdelt/events/*.export.CSV
Actor1Geo_CountryCode	1229	12122	9.74	/user/feeds/gdelt/events/*.export.CSV
Actor1Geo_FeatureID	1229	12122	9.74	/user/feeds/gdelt/events/*.export.CSV
Actor1Geo_FullName	1229	12122	9.74	/user/feeds/gdelt/events/*.export.CSV
Actor1Geo_Lat	1229	12122	9.74	/user/feeds/gdelt/events/*.export.CSV
Actor1Geo_Long	1229	12122	9.74	/user/feeds/gdelt/events/*.export.CSV
Actor1Geo_Type	0	12351	0.0	/user/feeds/gdelt/events/*.export.CSV



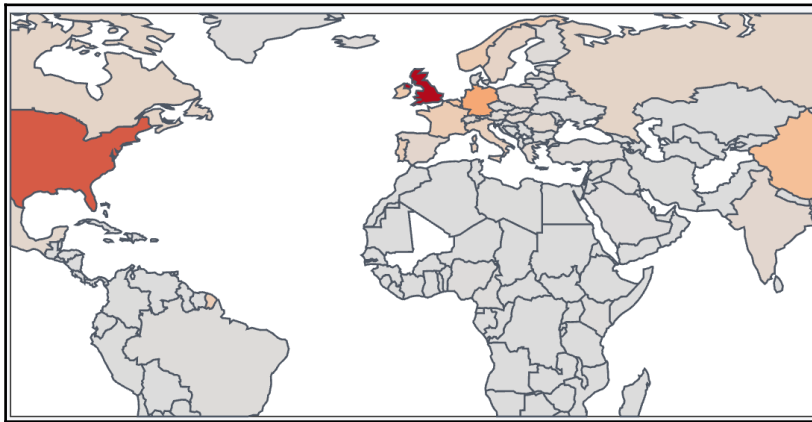
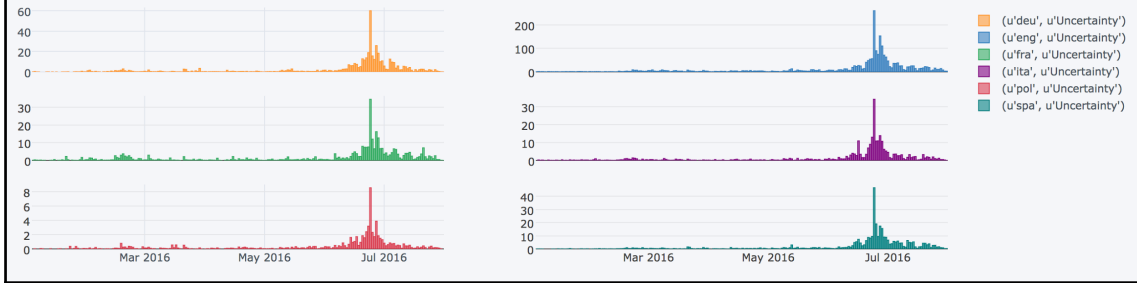
```

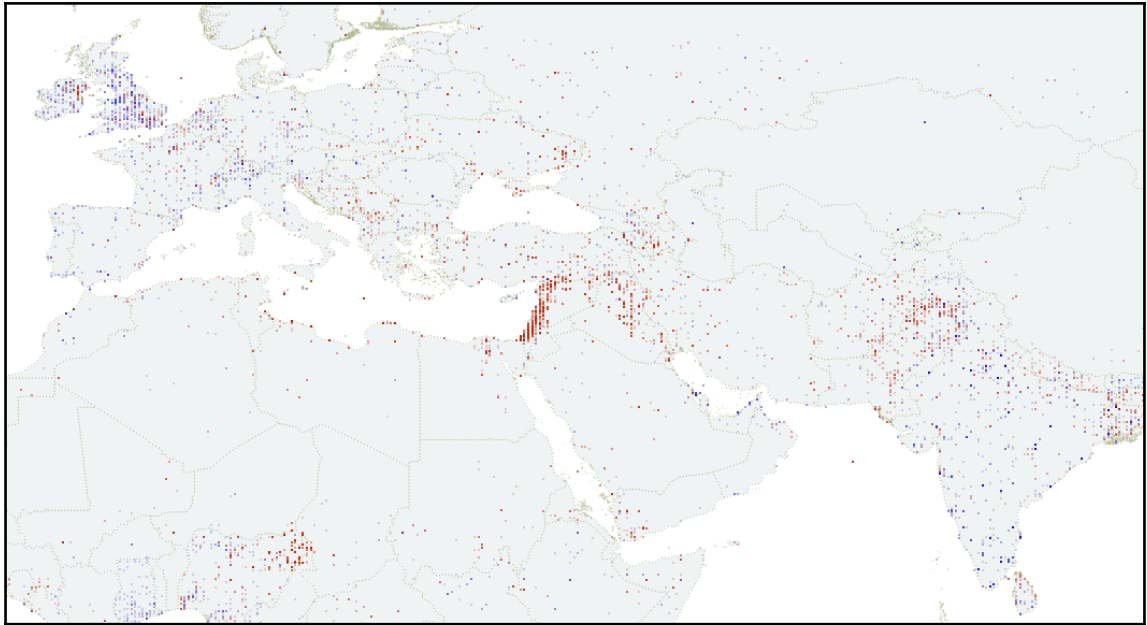
proReport: org.apache.spark.sql.DataFrame = [sourceStudied: string, metricDescriptor: string, fieldName: string, ingestTime: date, maskInstance: string, maskCount: bigint, description: st
-----
|sourceStudied|metricDescriptor|fieldName|ingestTime|maskInstance|maskCount|
-----
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_ADM1Code|2016-12-23|A9|1207||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_ADM1Code|2016-12-23|A|1108||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_ADM1Code|2016-12-23|1|136||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_ADM2Code|2016-12-23|19|1204||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_ADM2Code|2016-12-23|1|1088||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_ADM2Code|2016-12-23|A9|1139||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_CountryCode|2016-12-23|A|12315||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_CountryCode|2016-12-23|1|136||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_FeatureID|2016-12-23|1-9|1108||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_FeatureID|2016-12-23|A|1810||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_FeatureID|2016-12-23|19|1397||
|/user/feeds/gdelt/events/*.export.CSV|ASCII_LOWGRAIN|ActionGeo_FeatureID|2016-12-23|1|136||
  
```





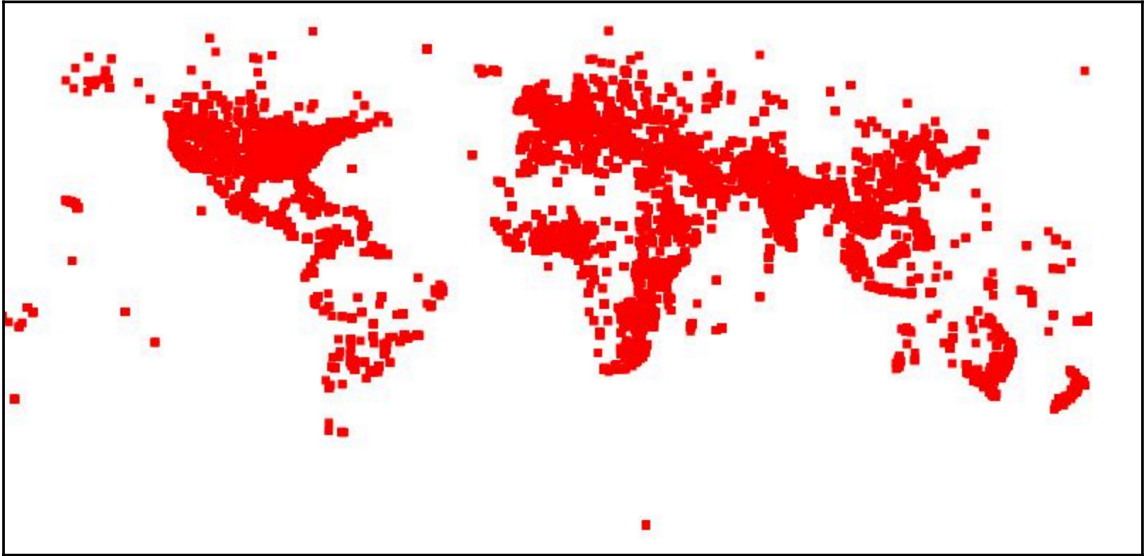
BREXIT: Daily GCAM Uncertainty by Language, 2016-01 through 2016-07

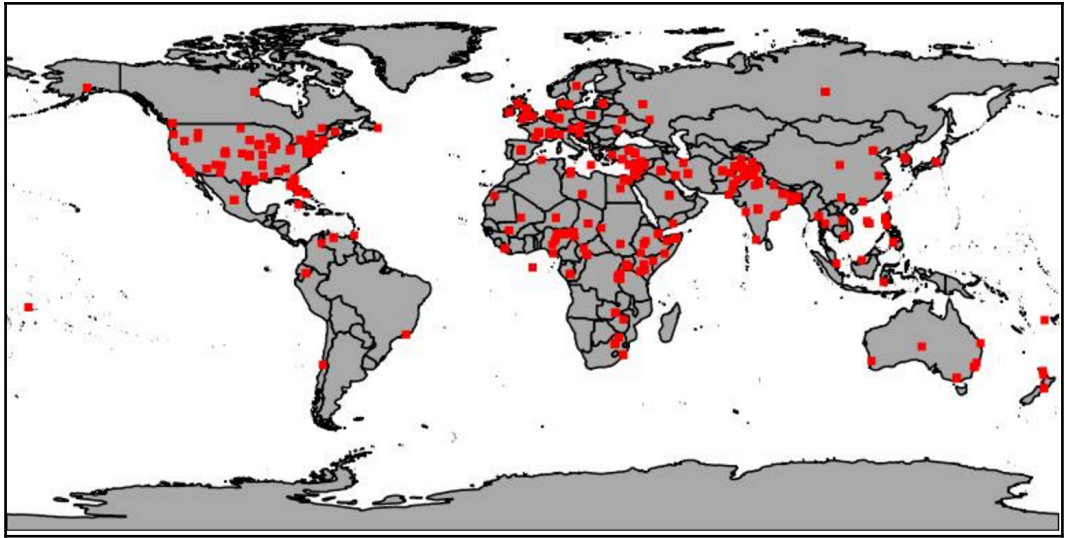




Chapter 05: Spark for Geographic Analysis

$$2r \arcsin \left(\sqrt{\sin^2 \left(\frac{\varphi_2 - \varphi_1}{2} \right) + \cos(\varphi_1) \cos(\varphi_2) \sin^2 \left(\frac{\lambda_2 - \lambda_1}{2} \right)} \right)$$





Chapter 06:Scraping Link-Based External Data

BBC Sign in News Sport More

Search


NEWS Find local news

Home UK World Business Politics Tech Science Health More

Entertainment & Arts

David Bowie dies of cancer aged 69

11 January 2016 Entertainment & Arts



David Bowie: "I'm a collector. I collect personalities"

Singer David Bowie, one of the most influential musicians of his era, has died of cancer at the age of 69.

A statement was issued on his social media accounts, saying he "died peacefully, surrounded by his family" after an "18-month battle with cancer".

Tributes have been paid from around the world to the "extraordinary artist" whose last album was released days ago.

Sir Paul McCartney **described him as a "great star"** who "played a very strong part in British musical history".

Bowie's son Duncan Jones, who is a Bafta-winning film director, **wrote on Twitter:** "Very sorry and sad to say it's true. I'll be offline for a while. Love to all."

Top Stories

Iran sanctions lifted over nuclear deal

International sanctions on Iran are lifted after its compliance with obligations under the nuclear agreement with world powers was certified.

1 hour ago


Corbyn warns firms over 'unfair' pay

6 hours ago

Broad gives England series win


9 hours ago

Features



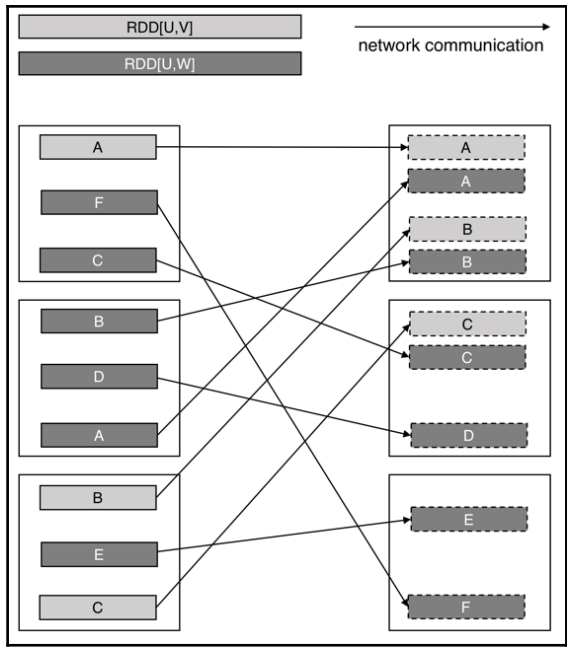
Change your tune

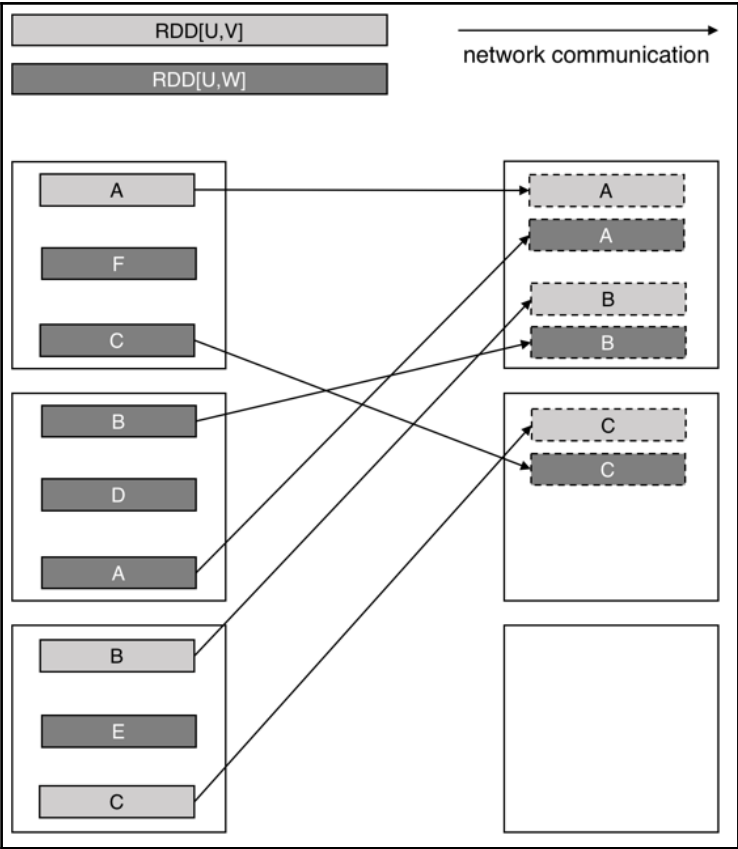
Is it time for a new British national anthem?

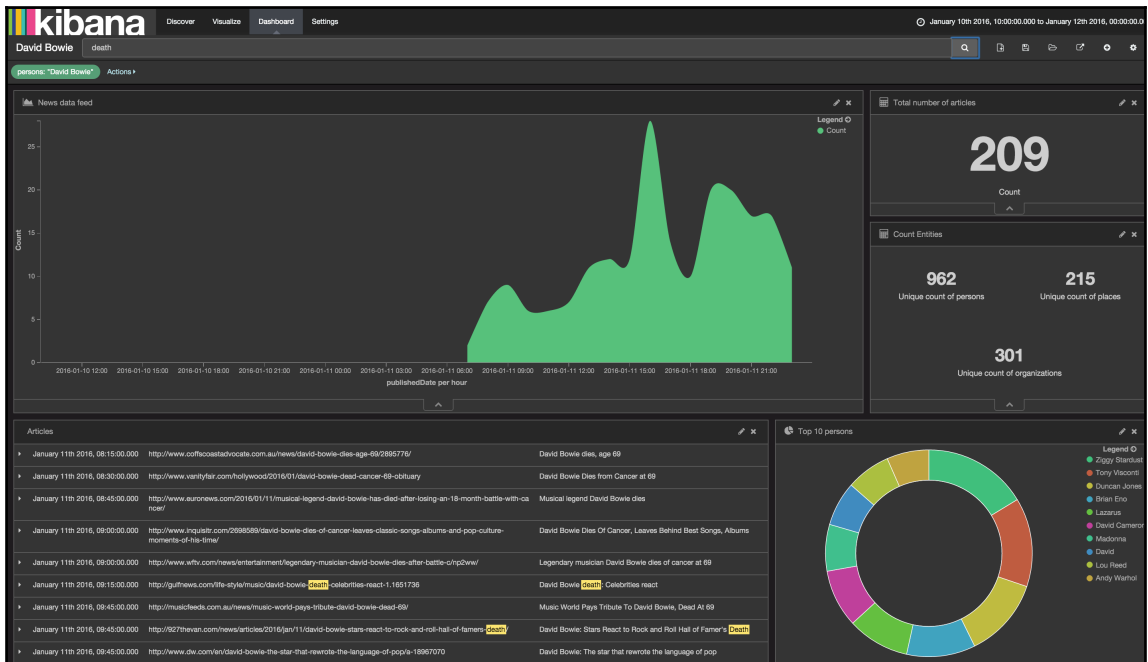
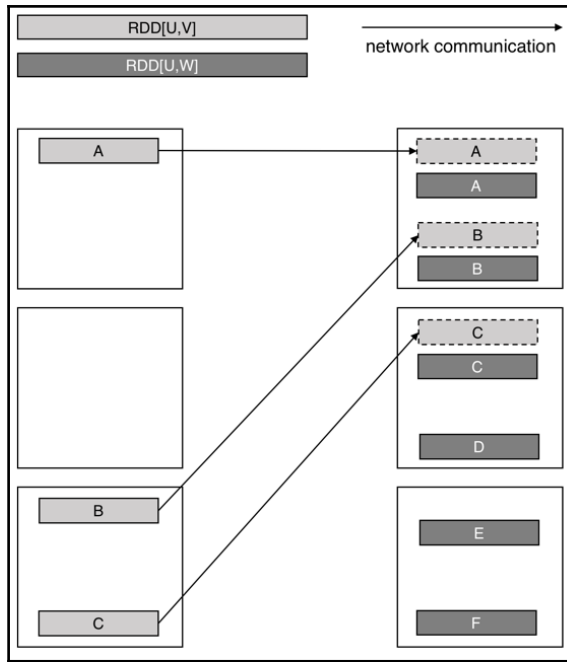


Choking up

The city in love with









Gianfranco Ravasi ✓

@CardRavasi

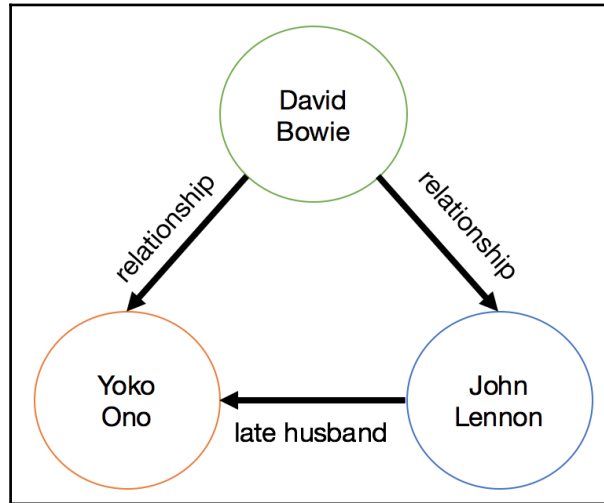
 Follow

Ground Control to Major Tom
Commencing countdown,
engines on
Check ignition
and may God's love be with you (David Bowie)

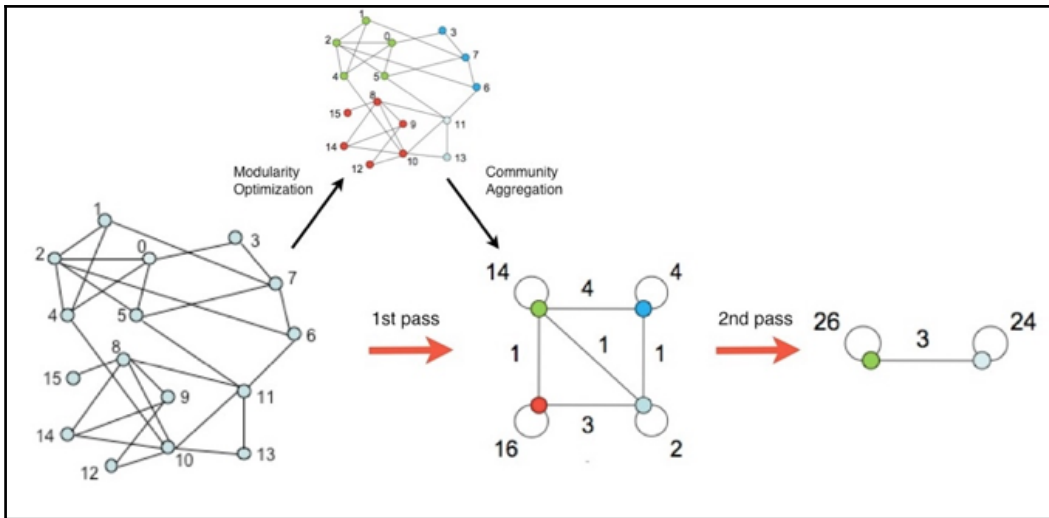
8:35 AM - 11 Jan 2016

  4,226  4,146

Chapter 07: Building Communities



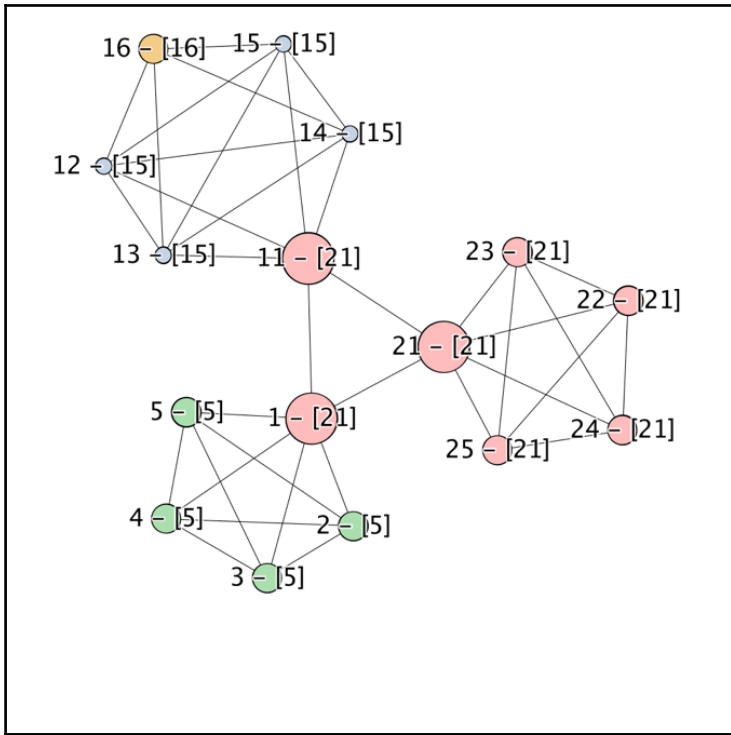
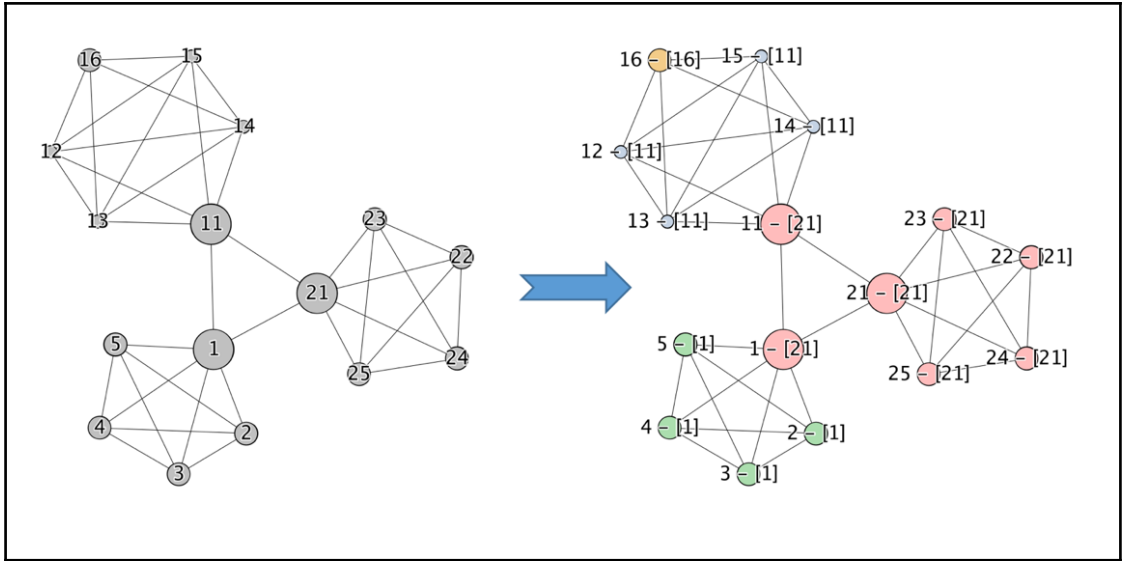
row key	column families					
	relationA			relationB		
	qualifier	value	visibility	qualifier	value	visibility
personA	personB	1	INTERNAL	personD	1	SECRET
	personC	1	CONFIDENTIAL			
personB	personC	1	INTERNAL	personD	1	CONFIDENTIAL



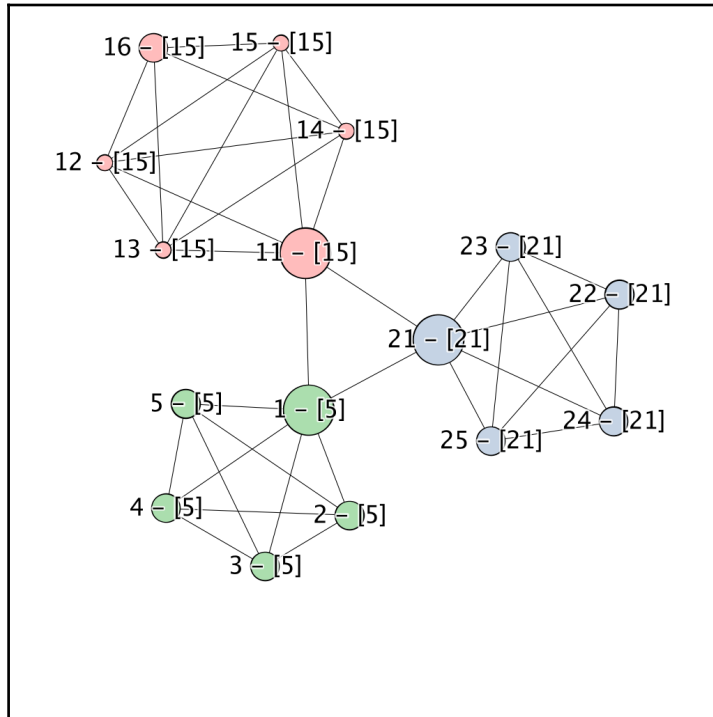
$$WCC(x, C) = \begin{cases} \frac{t(x, C)}{t(x, V)} \cdot \frac{vt(x, V)}{|C \setminus \{x\}| + vt(x, V \setminus C)} & \text{if } t(x, V) \neq 0; \\ 0 & \text{if } t(x, V) = 0. \end{cases}$$

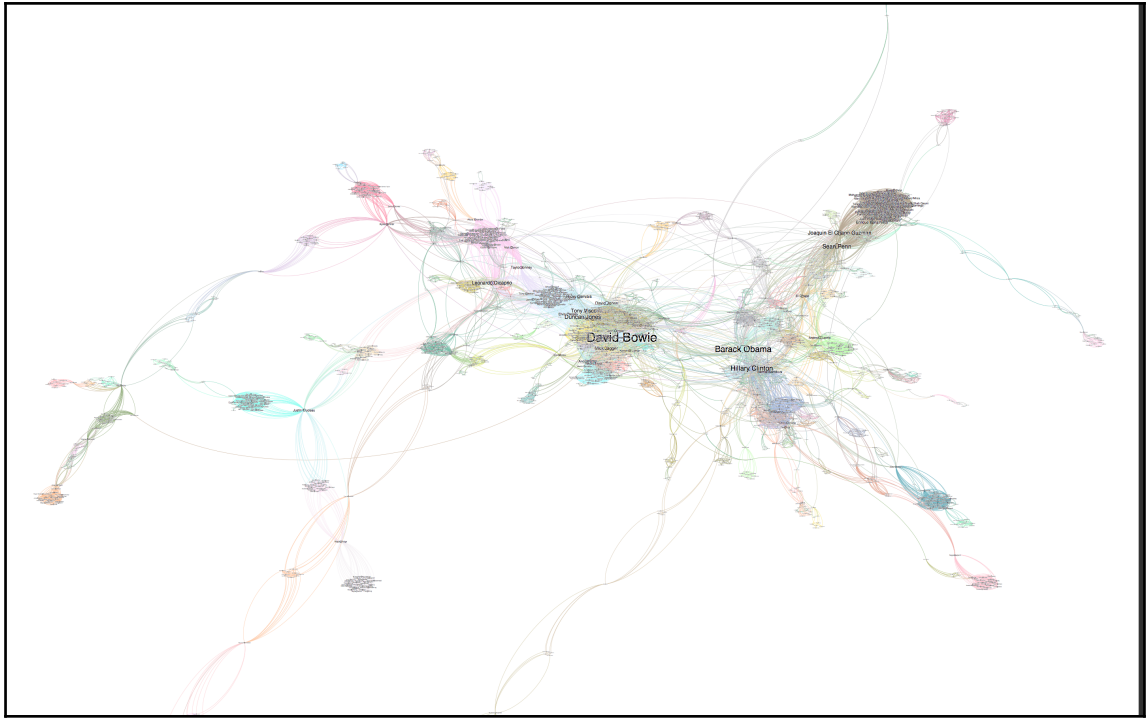
$$WCC(S) = \frac{1}{|C|} \sum_{x \in S} WCC(x, C).$$

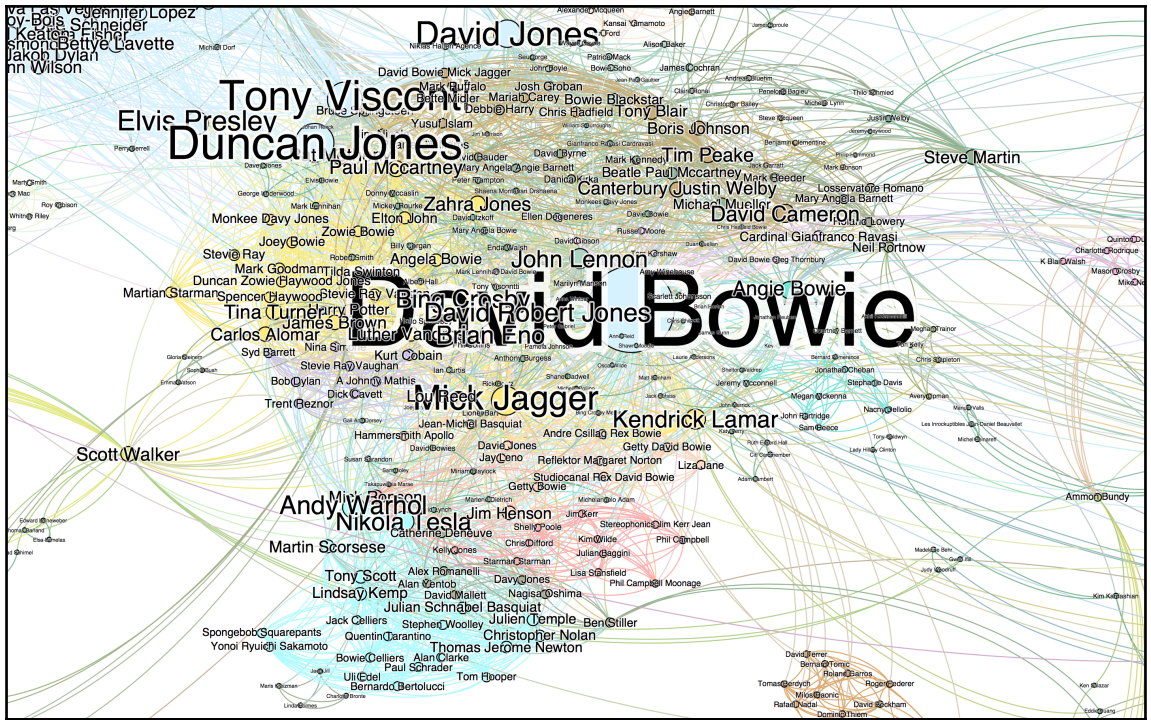
$$WCC(\mathcal{P}) = \frac{1}{|V|} \sum_{i=1}^n (|C_i| \cdot WCC(C_i)).$$

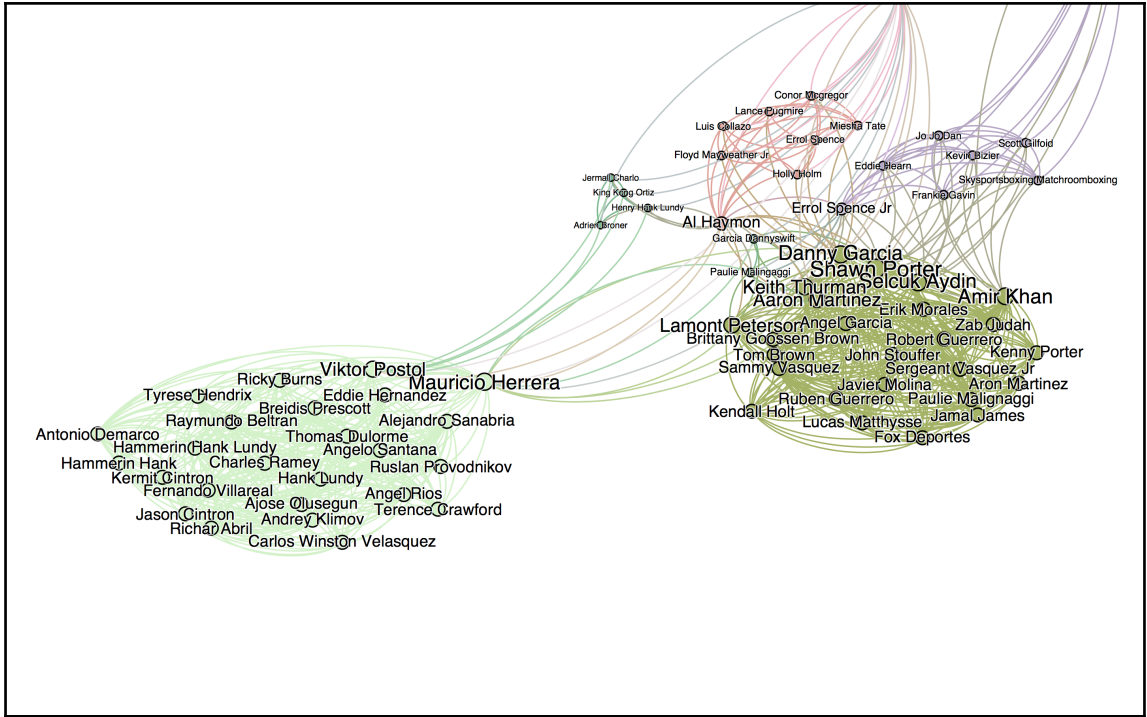


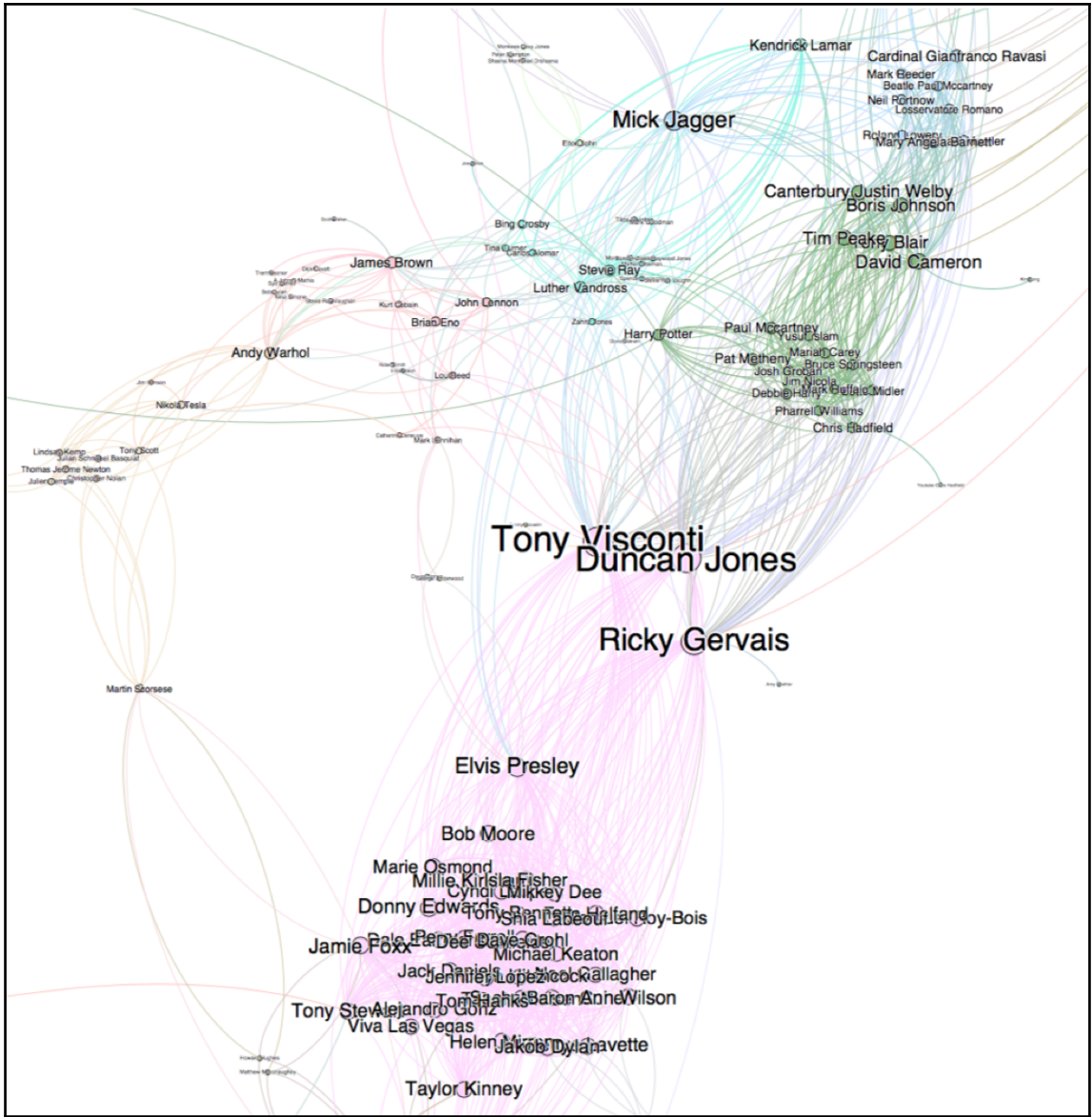
$$\begin{aligned}
 WCC(P') - WCC(P) &= WCC'_I(v, C) \\
 &= \frac{1}{V} \cdot (d_{in} \cdot \Theta_1 + (r - d_{in}) \cdot \Theta_2 + \Theta_3)
 \end{aligned}$$



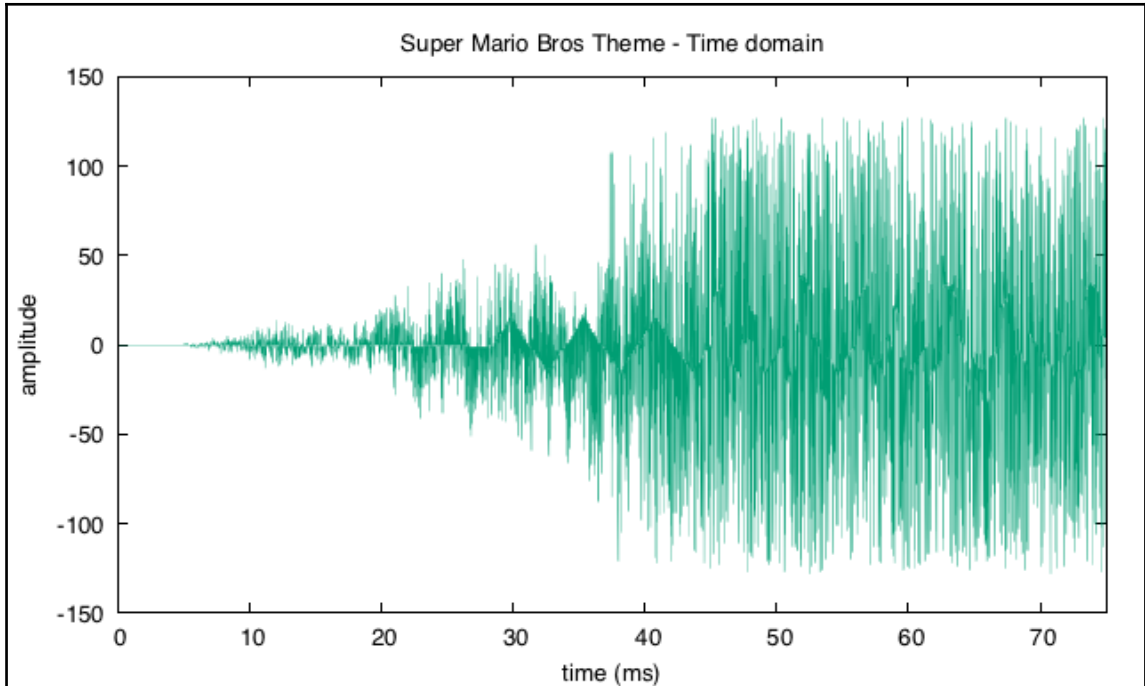




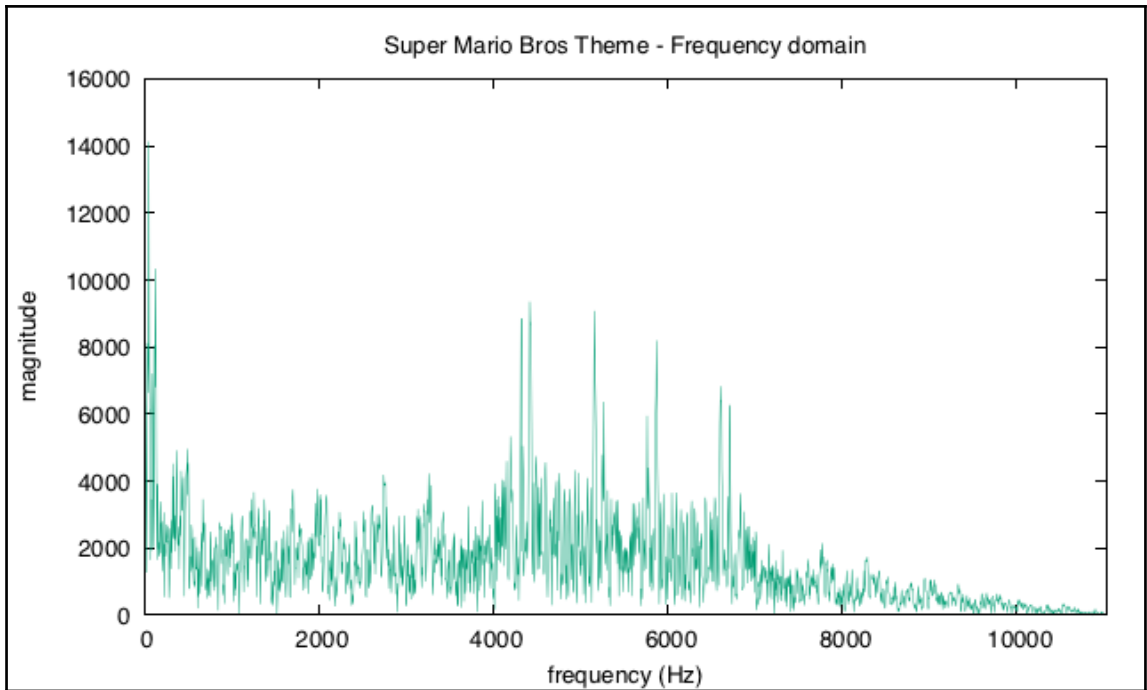




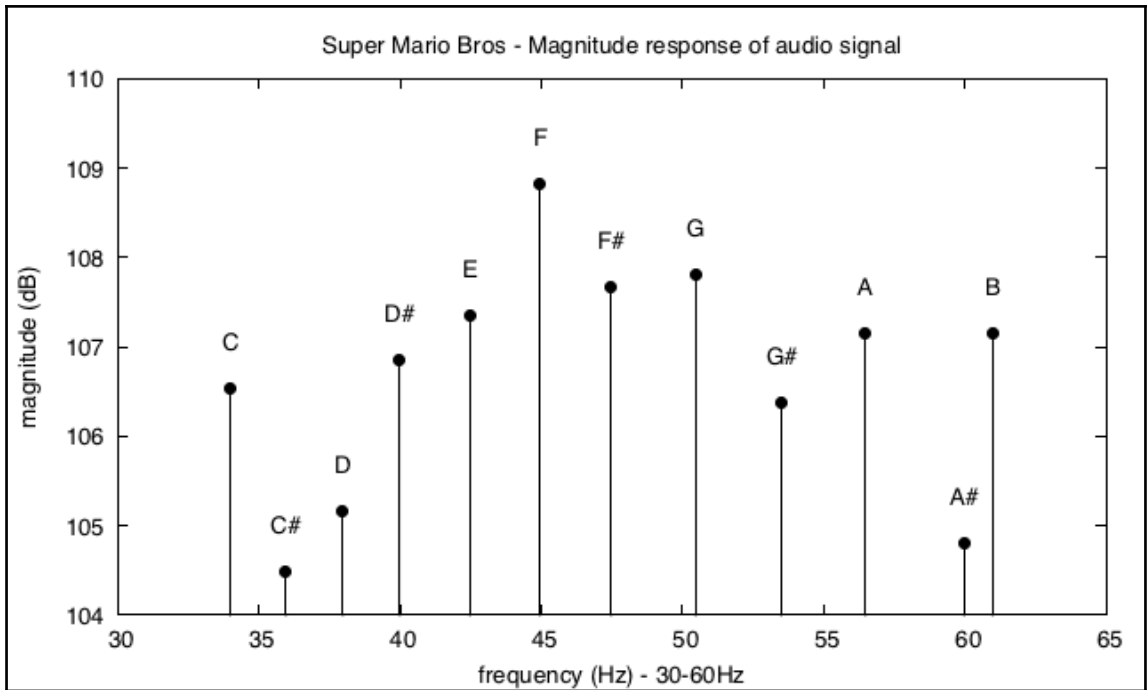
Chapter 08: Building a Recommendation System

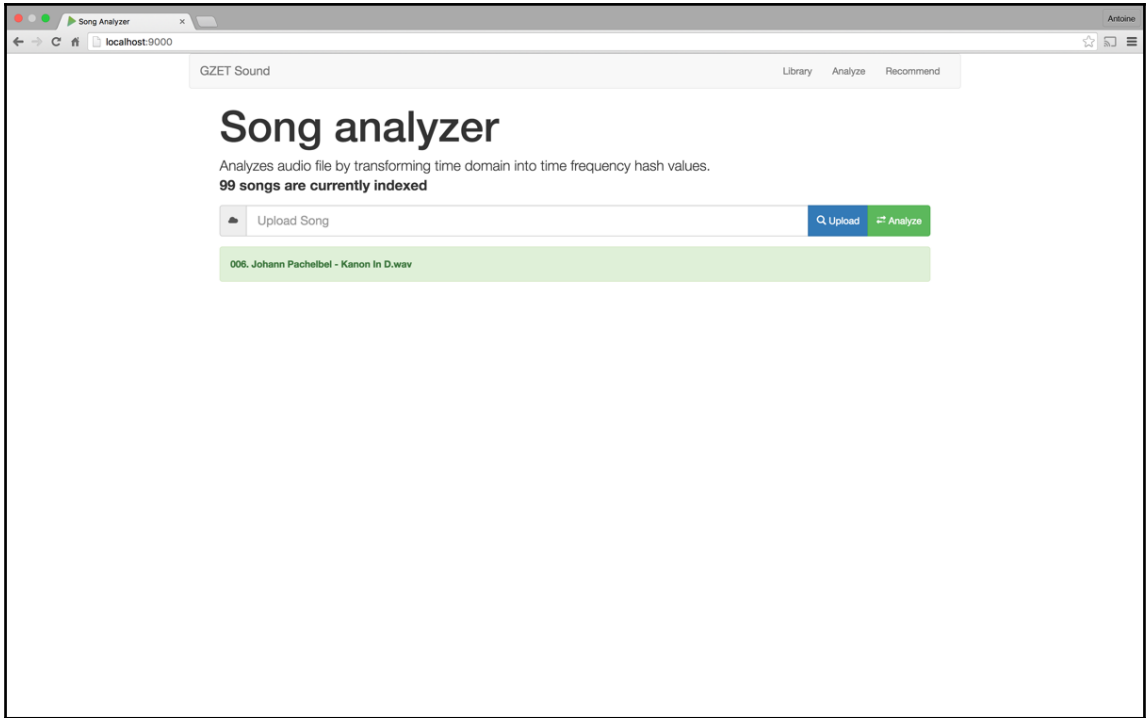


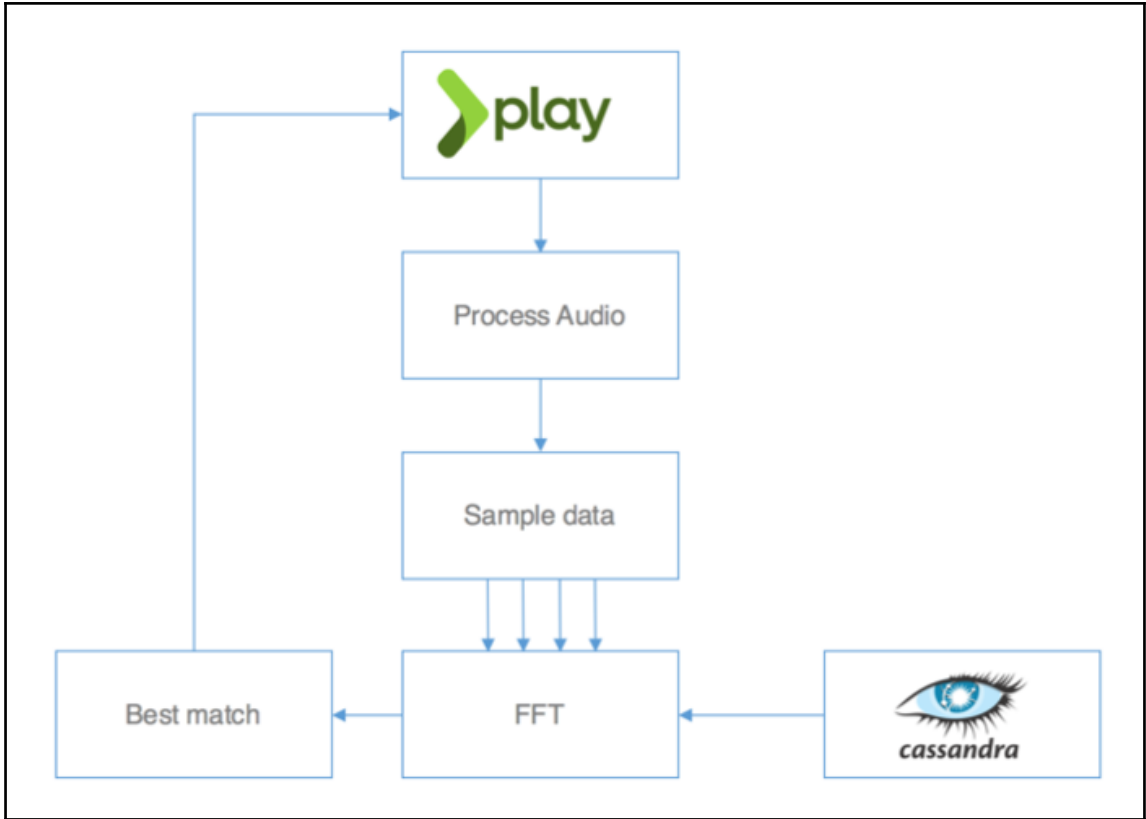
$$F(n) = \sum_{k=0}^{N-1} x(k) e^{\frac{-j2\pi kn}{N}}$$



	C	C#	D	Eb	E	F	F#	G	G#	A	Bb	B
0	16.35	17.32	18.35	19.45	20.60	21.83	23.12	24.50	25.96	27.50	29.14	30.87
1	32.70	34.65	36.71	38.89	41.20	43.65	46.25	49.00	51.91	55.00	58.27	61.74
2	65.41	69.30	73.42	77.78	82.41	87.31	92.50	98.00	103.8	110.0	116.5	123.5
3	130.8	138.6	146.8	155.6	164.8	174.6	185.0	196.0	207.7	220.0	233.1	246.9
4	261.6	277.2	293.7	311.1	329.6	349.2	370.0	392.0	415.3	440.0	466.2	493.9
5	523.3	554.4	587.3	622.3	659.3	698.5	740.0	784.0	830.6	880.0	932.3	987.8
6	1047	1109	1175	1245	1319	1397	1480	1568	1661	1760	1865	1976
7	2093	2217	2349	2489	2637	2794	2960	3136	3322	3520	3729	3951
8	4186	4435	4699	4978	5274	5588	5920	6272	6645	7040	7459	7902







Playlist

localhost:9000/playlist

GZET Sound

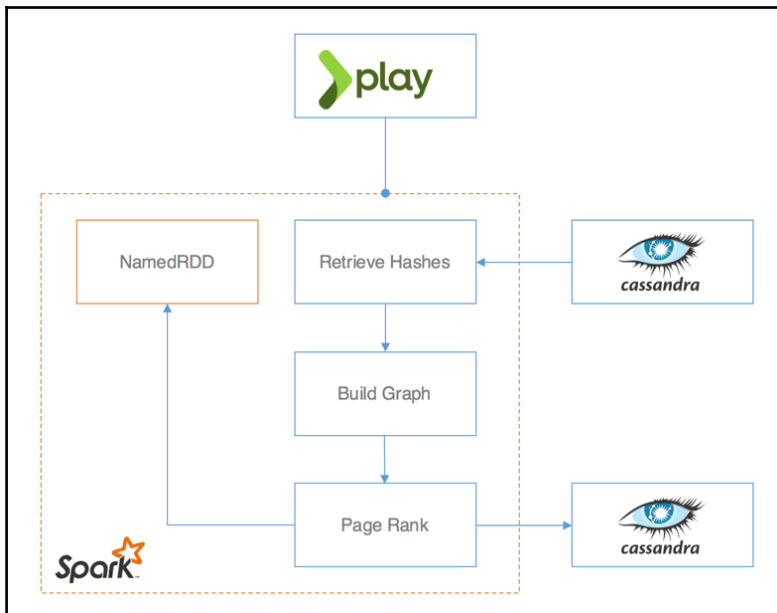
Library Analyze Recommend

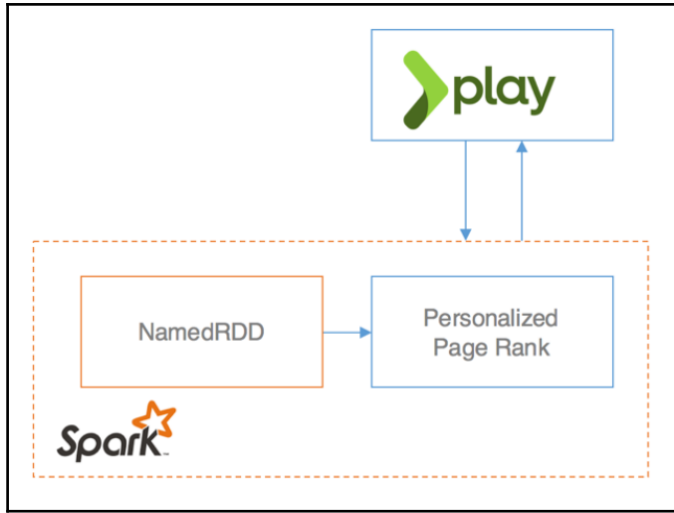
Playlist recommender

We run a page rank algorithm on frequency graph cooccurrence
 We rank playlist by chord frequency popularity with a 10% chance of random teleport
99 songs are currently ranked

Re-Analyze

▶ 015. Gustav Mahler - Symphony No. 5.wav	1.0618424685516701
▶ 010. Tomaso Giovanni Albinoni - Adagio In Sol Minore.wav	1.0519542177068021
▶ 029. Maurice Ravel - Boléro.wav	1.051923303138239
▶ 035. Ludwig van Beethoven - Für Elise (WoO 59).wav	1.0517504349088352
▶ 024. Bedřich Smetana - Má Vlast - Vltava.wav	1.0421080201856098
▶ 006. Johann Pachelbel - Canon In D.wav	1.0420447197782963
▶ 016. Ludwig van Beethoven - Symphony No. 9 (Op. 125).wav	1.0419668982009602
▶ 072. Pyotr Ilyich Tchaikovsky - The Nutcracker (Op. 71) - Dance Of The Sugar-Plum Fairy - $\text{C} \rightarrow \text{u} \rightarrow \text{f} - \text{E} - \text{D} - \text{A} - [\text{f}] - \text{J}$.wav	1.041928716835439
▶ 086. Wolfgang Amadeus Mozart - Konzert F \sharp F \sharp F \sharp , Harle und Orchester (K. 299) - Allegro.wav	1.0418360584591596
▶ 003. Ludwig van Beethoven - Piano Concerto No. 5 (Op. 73) - Adagio Un Poco Mosso.wav	1.041831082384432
▶ 014. Edvard Hagerup Grieg - Peer Gynt Suite No. 1 (Op. 46) - Morgenstemning.wav	1.0323285474135953
▶ 033. Charles-François Gounod - Ave Maria.wav	1.0322969309038053
▶ 038. Max Christian Friedrich Bruch - Violinkonzert Nr. 1 (Op. 26) - Allegro Moderato.wav	1.0322903467086827
▶ 022. Giulio Caccini - Ave Maria.wav	1.032287546540544
▶ 039. Wolfgang Amadeus Mozart - Requiem (K. 626) - Introitus.wav	1.0322714284371724
▶ 080. Joaquín Rodrigo Vidre - Concierto De Aranjuez - Adagio.wav	1.032263463334513





Playlist

localhost:9000/playlist/4

GZET Sound Library Analyze Recommend

Playlist recommender

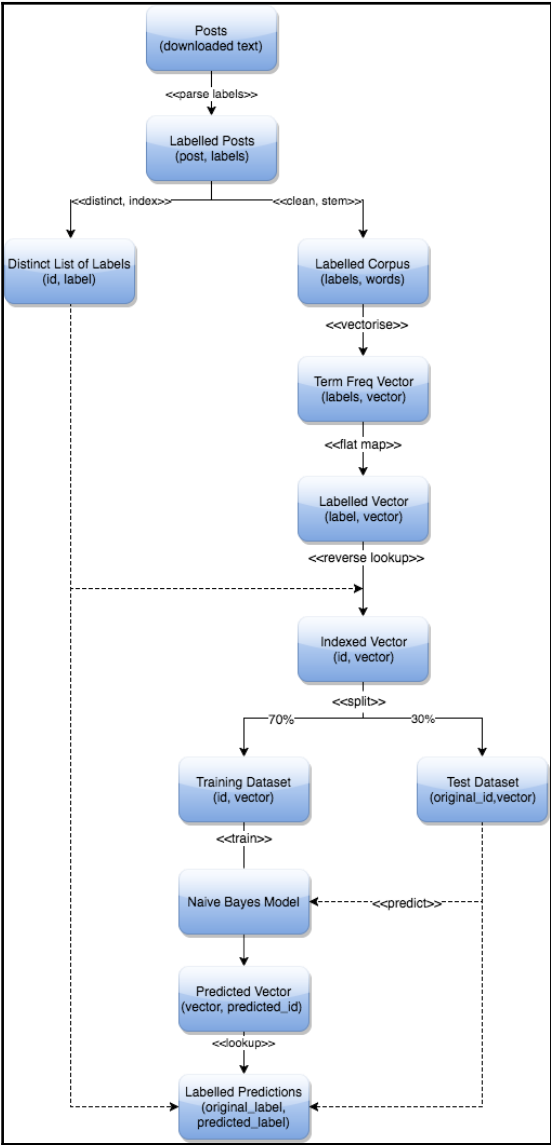
004. Wolfgang Amadeus Mozart - Klarinettenkonzert (K. 622) - Adagio.wav

We run a personalized page rank algorithm on frequency graph cooccurrence
 We rank playlist by chord frequency popularity with a 10% chance of teleport
99 songs are currently ranked

[Re-Analyze](#)

▶ 004. Wolfgang Amadeus Mozart - Klarinettenkonzert (K. 622) - Adagio.wav	0.100892919605565
▶ 015. Gustav Mahler - Symphony No. 5.wav	0.8363271122440324E-4
▶ 052. Gabriel Urbain Fauré - Requiem (Op. 48) - Pie Jesu.wav	0.728505730456107E-4
▶ 010. Tomaso Giovanni Albinoni - Adagio In Sol Minore.wav	0.724724912391168E-4
▶ 022. Giulio Caccini - Ave Maria.wav	0.722727191418038E-4
▶ 029. Maurice Ravel - Boléro.wav	0.72052290143364E-4
▶ 021. Carl Orff - Carmina Burana - O Fortuna.wav	0.71748442655657E-4
▶ 016. Ludwig van Beethoven - Symphony No. 9 (Op. 125).wav	0.716801140229047E-4
▶ 035. Ludwig van Beethoven - Für Elise (WoO 59).wav	0.70280932797553E-4
▶ 072. Pyotr Ilyich Tchaikovsky - The Nutcracker (Op. 71) - Dance Of The Sugar-Plum Fairy --♩-μ-♯-♯-É-Ω-♯-♯-♯.wav	0.701820867338293E-4
▶ 086. Wolfgang Amadeus Mozart - Konzert F,Fr F,ùte	0.68953098390811E-4
▶ 038. Max Christian Friedrich Bruch - Violinkonzert Nr. 1 (Op. 26) - Allegro Moderato.wav	0.627709848951405E-4
▶ 014. Edvard Hagerup Grieg - Peer Gynt Suite No. 1 (Op. 46) - Morgensterming.wav	0.605458323541484E-4
▶ 080. Joaquín Rodrigo Vidre - Concierto De Aranjuez - Adagio.wav	0.624126854039341E-4
▶ 034. Brahms Johannes - 1. Klavierkonzert - 1. Satz.wav	0.61488366489965E-4

Chapter 09: News Dictionary and Real-Time Tagging System





WFMY News 2  @WFMY · Apr 24

Adorable #Prince George misses bedtime, meets President #Obama --
ow.ly/4n2snm

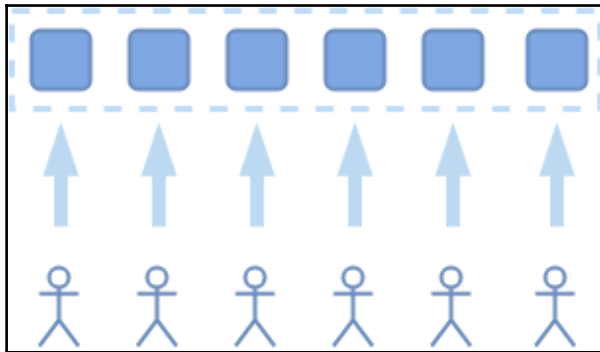


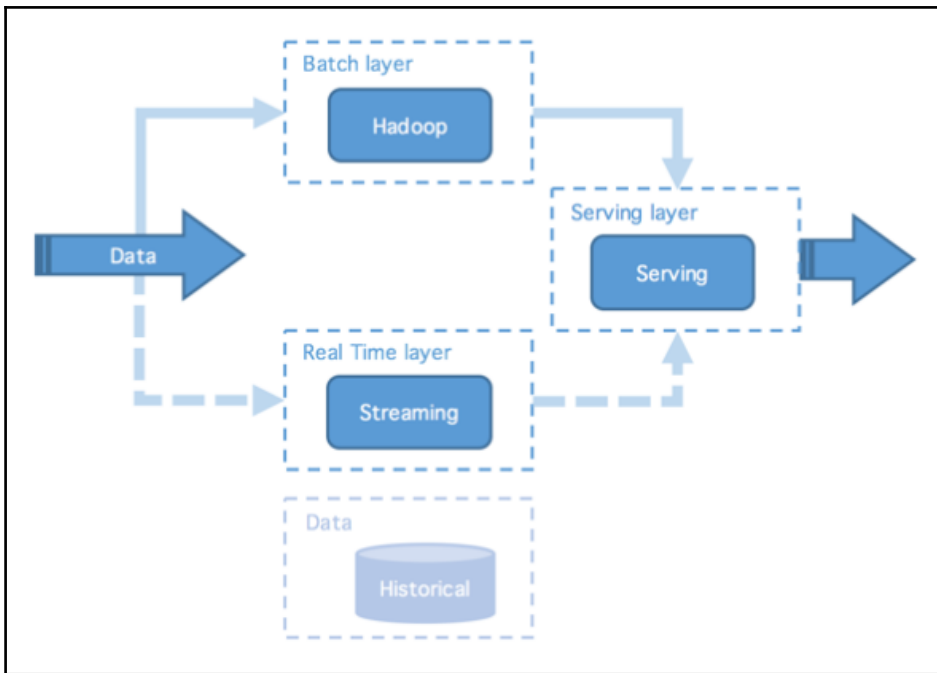
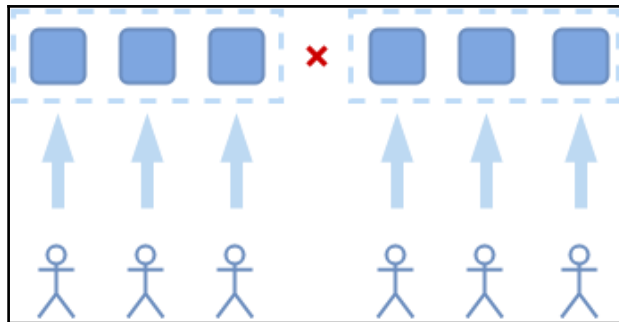


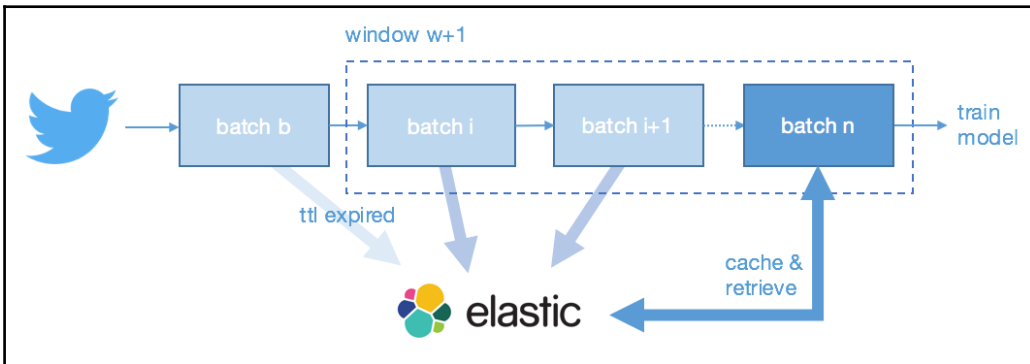
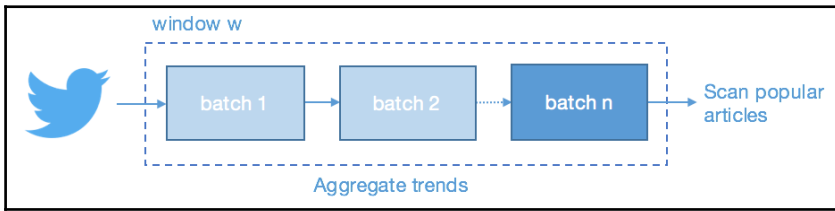
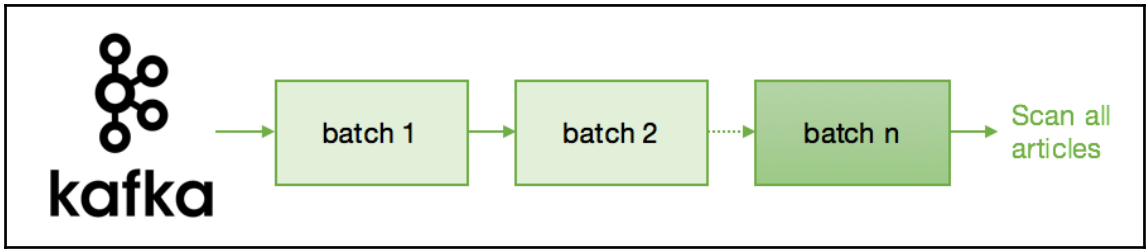
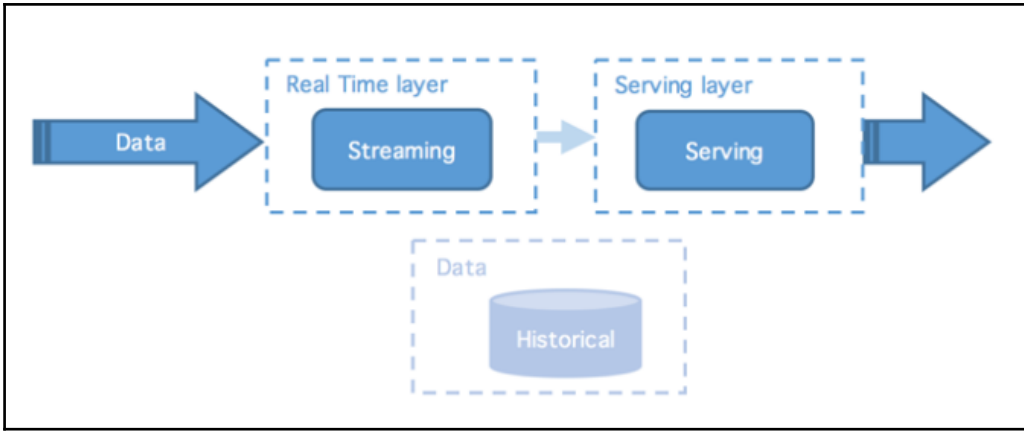
Kids of the Eighties @kidsofeighties · 30 Dec 2016

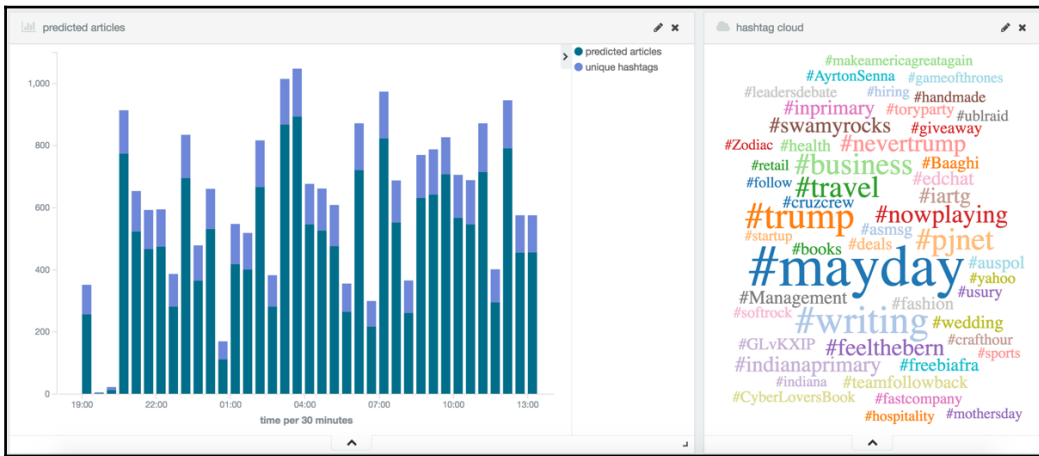
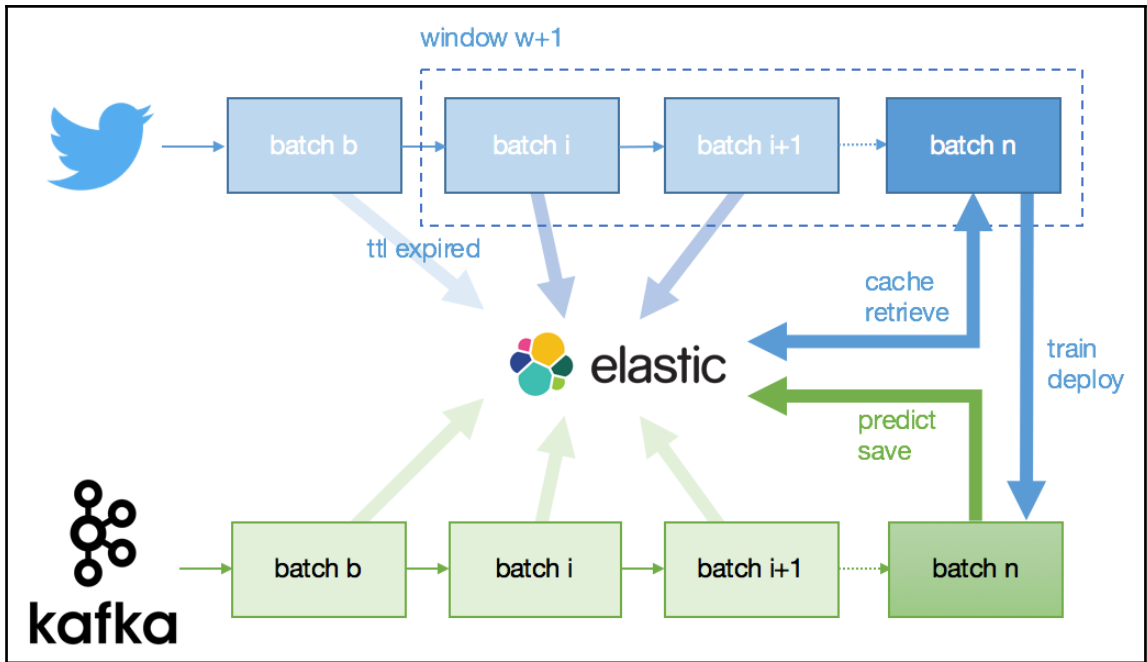
Death stars: music's great losses of 2016 via @guardian

theguardian.com/music/2016/dec... #DavidBowie #PRINCE #georgemichael #LeonardCohen









Chapter 10: Story De-duplication and Mutation



Antoine Amend

Just now · London · 🔒

[http://www.forbes.com/.../samsung-south-korean-government-i.../...](http://www.forbes.com/.../samsung-south-korean-government-i.../)



Galaxy Note 7 Fiasco: Samsung, South Korean Government Launches Investigation

Industry analysts say investigating why the Note 7 devices caught fire is crucial for the world's largest smartphone maker.

FORBES.COM | BY JOHN KANG

👍 Like 💬 Comment ➦ Share



Write a comment...



RELATED ARTICLES



Why Samsung Abandoned Its Galaxy Note 7 Flagship Phone

The unprecedented move by the South Korean electronics giant is an embarrassing reversal for a...

THE NEW YORK TIMES · 11,181 SHARES [Share](#) [Save](#)



Singapore Airlines bans Samsung Galaxy Note 7 on its flights

SINGAPORE Singapore Airlines said on Saturday it has banned Samsungs Galaxy Note 7 mobile phon...

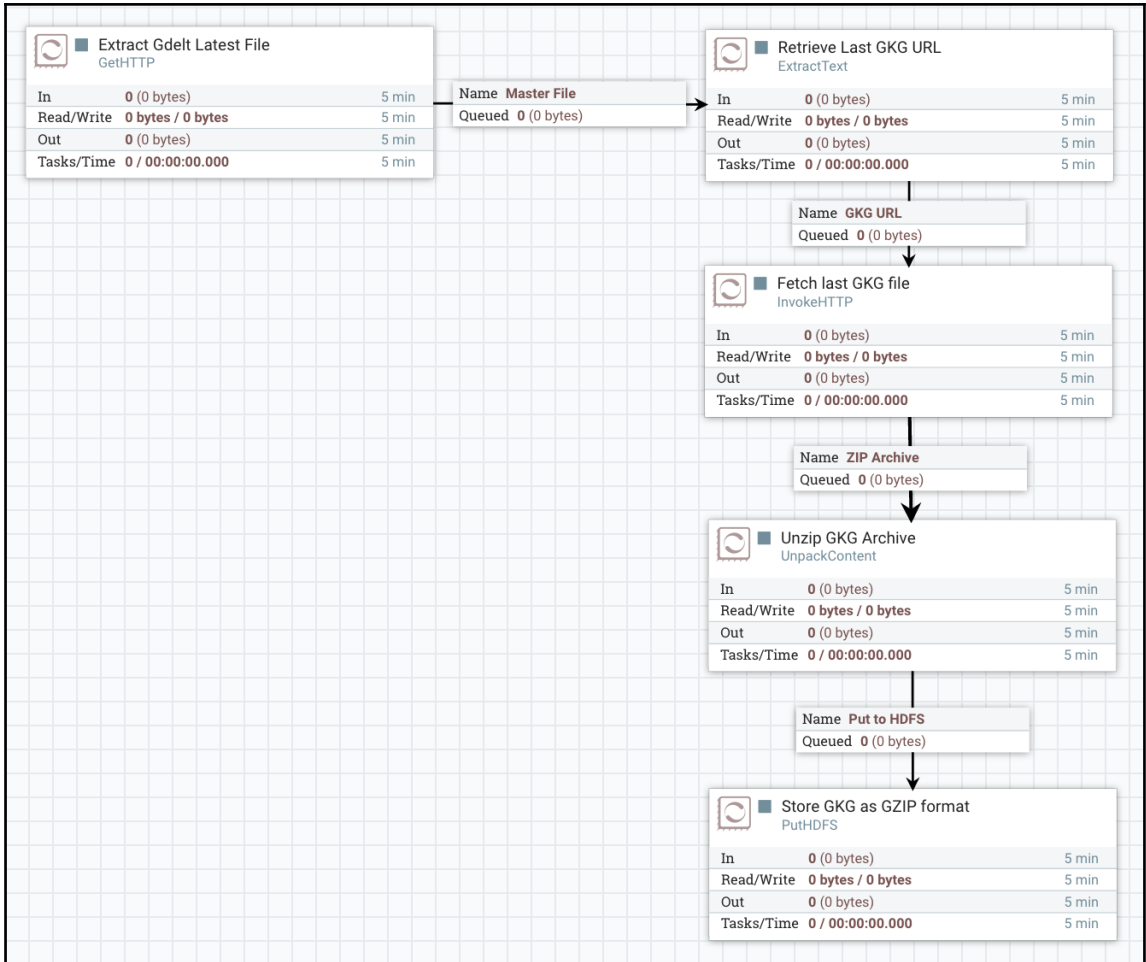
GMA NEWS · 4,386 SHARES [Share](#) [Save](#)

		hashcodes										
he	1	1	0	0	1	1	1	1	1	1	0	1
el	1	1	0	0	1	0	1	0	0	1	1	1
ll	1	1	0	1	1	0	0	0	0	0	0	0
lo	1	1	0	1	1	0	0	0	0	0	1	1
si	1	1	1	0	0	1	0	1	0	1	1	0
im	1	1	0	1	0	0	1	0	0	1	0	0
mh	1	1	0	1	1	0	0	1	1	0	1	1
ha	1	1	0	0	1	1	1	1	1	0	0	1
as	1	1	0	0	0	0	1	1	0	0	1	0
sh	1	1	1	0	0	1	0	1	0	1	0	1

	10	10	-6	-2	2	-2	0	2	-4	0	0	2
hello simhash	1	1	0	0	1	0	0	1	0	0	0	1

hello simhash	1	1	0	0	1	0	0	1	0	0	0	1
hello minhash	1	1	0	0	1	0	1	1	0	0	0	1
hello world	1	1	0	0	1	0	1	0	0	0	0	0

hello simhash	1	1	0	0	1	0	0	1	0	0	0	1
hello minhash	1	1	0	0	1	0	1	1	0	0	0	1
XOR	0	0	0	0	0	0	1	0	0	0	0	0



Note 7 fiasco leaves Samsung's smartphone brand in question

By [ASSOCIATED PRESS](#)

PUBLISHED: 20:02, 12 October 2016 | **UPDATED:** 20:02, 12 October 2016



SEOUL, South Korea (AP) — The fiasco of Samsung's fire-prone Galaxy Note 7 smartphones — and its stumbling response to the problem — has left consumers from Shanghai to New York reconsidering how they feel about the South Korean tech giant and its products.

Samsung Electronics said this week that it would stop making the Note 7 for good, after first recalling some devices and then recalling their replacements, too. Now, like the makers of Tylenol, Ford Pintos and other products that faced crises in the past, it must try to restore its relationship with customers as it repairs damage to its brand.

Samsung shares plunged as much as 8 percent in Seoul, their biggest one-day drop since the 2008 financial crisis, after the company apologized for halting sales of the Note 7.



Note 7 Fiasco Leaves Samsung's Smartphone Brand in Question



A visitor passes by an advertisement of the Samsung Electronics Galaxy Note 7 smartphone at its shop in Seoul, South Korea, Oct. 11, 2016.

SEOUL, SOUTH KOREA — The fiasco of Samsung's fire-prone Galaxy Note 7 smartphones — and its stumbling response to the problem — has left consumers from Shanghai to New York reconsidering how they feel about the South Korean tech giant and its products.

Related

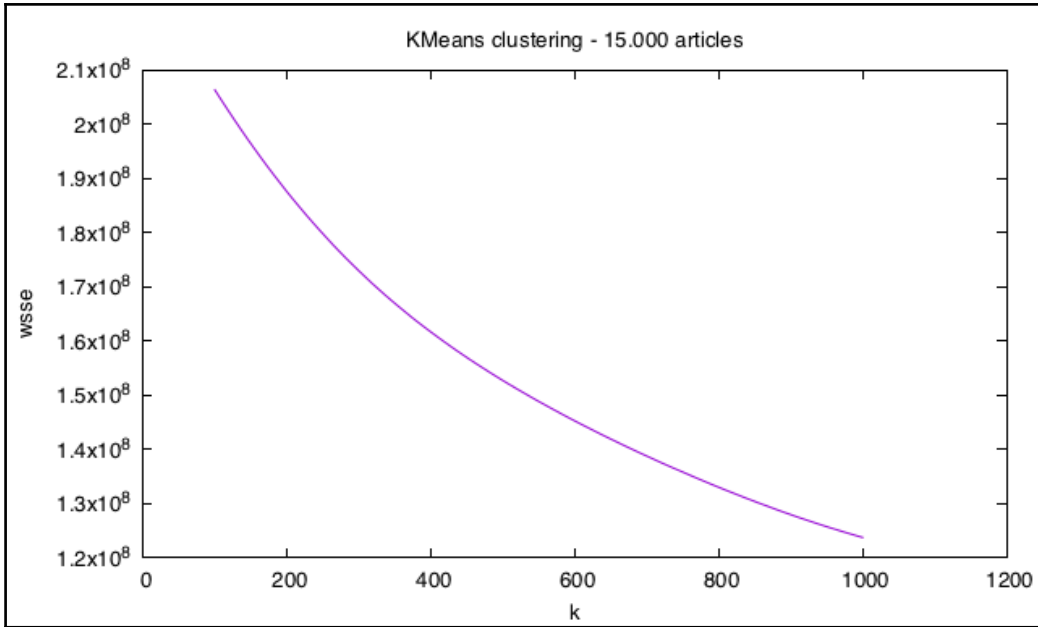


Samsung Drastically Cuts Profit Estimates As Malfunctions Continue

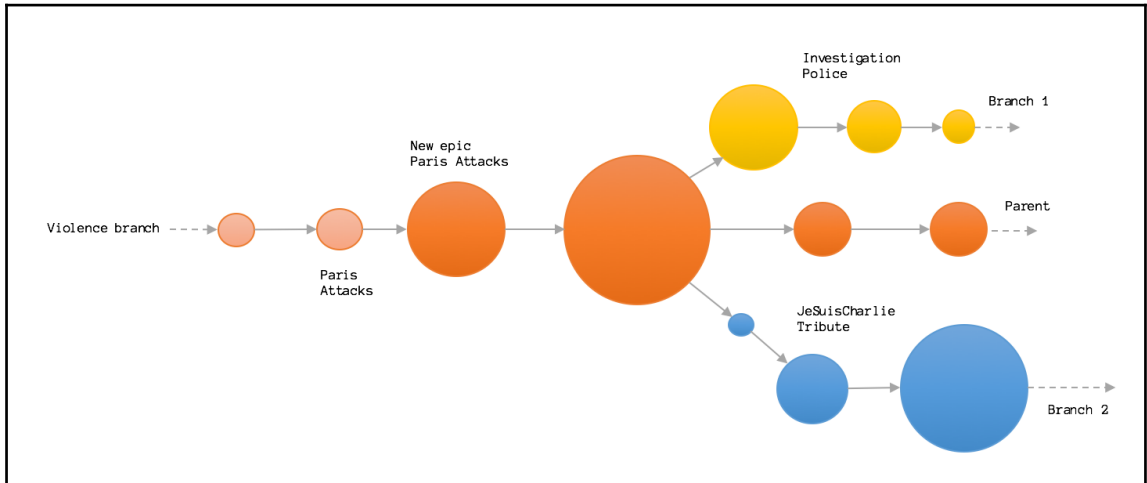
$$idf_i = \log\left(\frac{n+1}{df_i+1}\right)$$

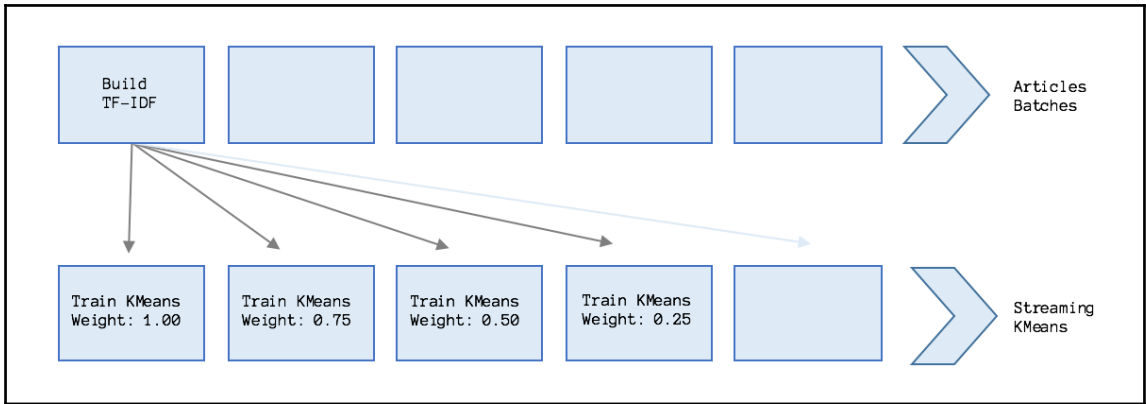
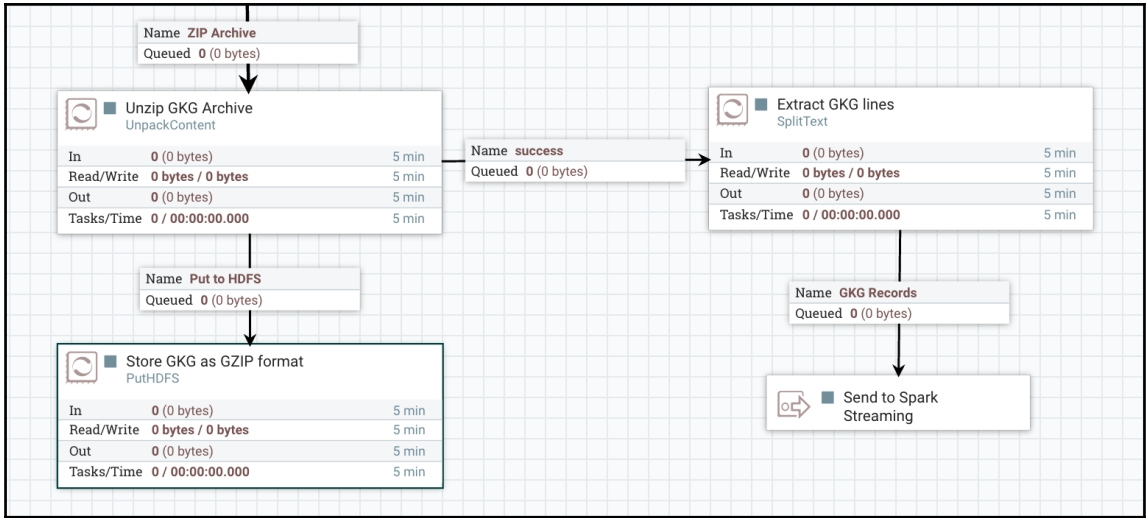
	Random Contexts								
mastering	1	0	0	0	0	1	1	0	-1
spark	0	1	0	0	1	1	1	-1	0
for	1	0	0	0	1	1	-1	0	1
data	0	0	0	0	0	0	0	1	0
science	-1	1	1	0	1	1	0	0	-1

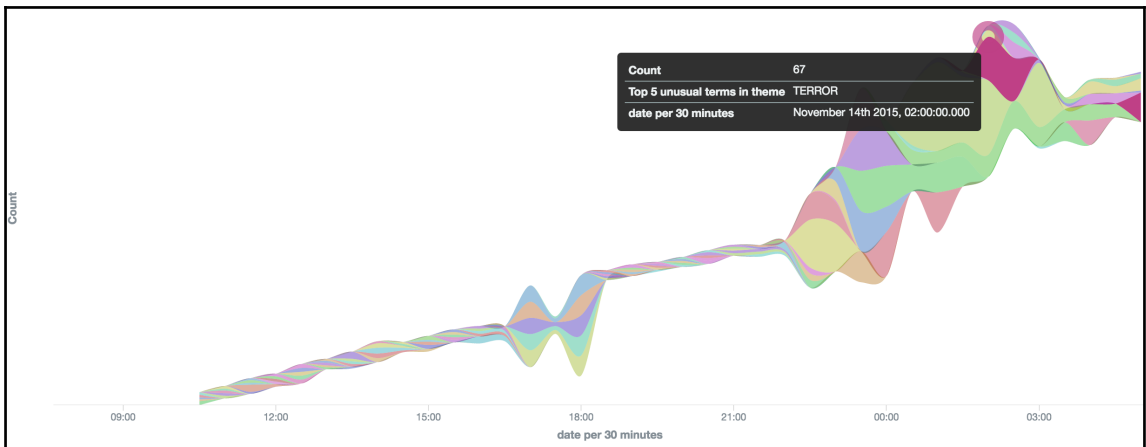
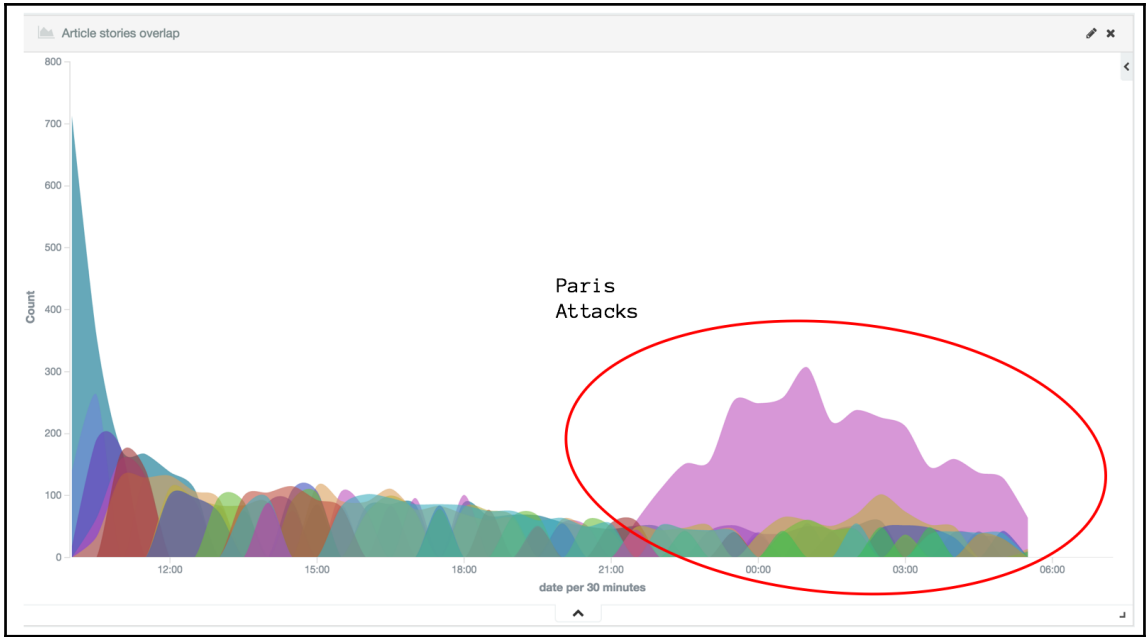
	Aggregated Context								
mastering spark for data science	1	2	1	0	3	4	1	0	-1



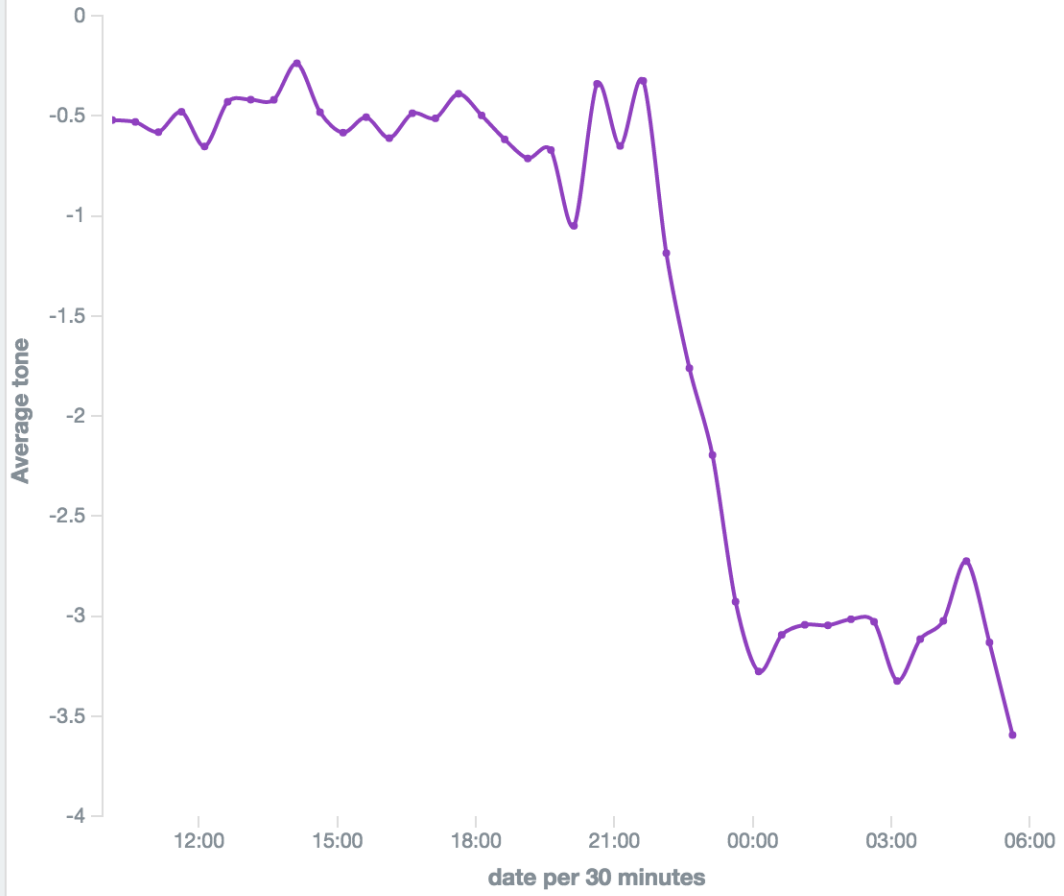
$$k \approx \sqrt{\frac{n}{2}}$$

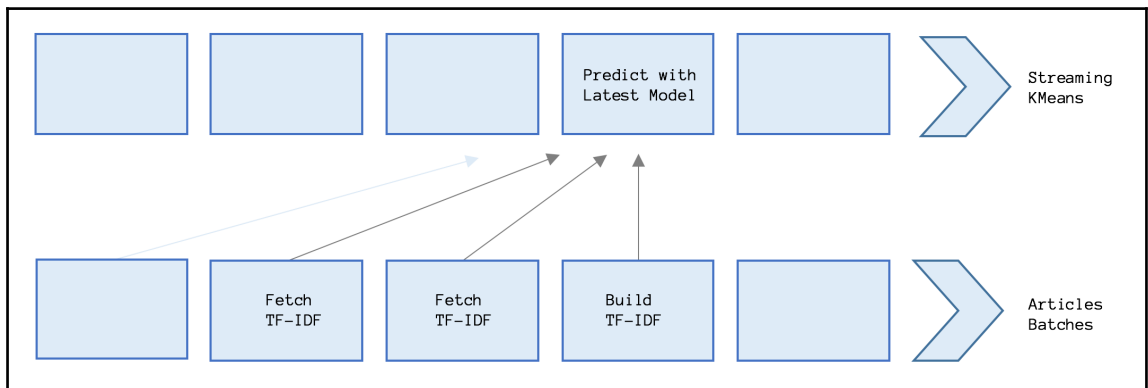
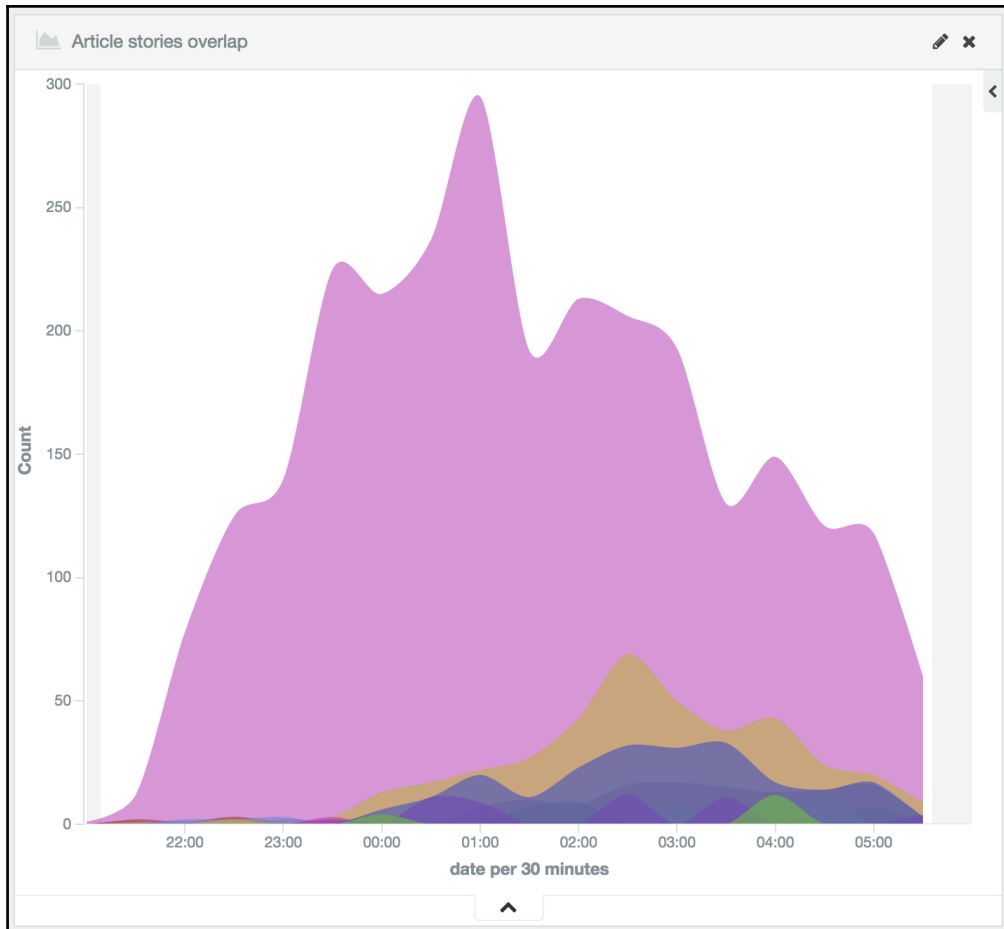


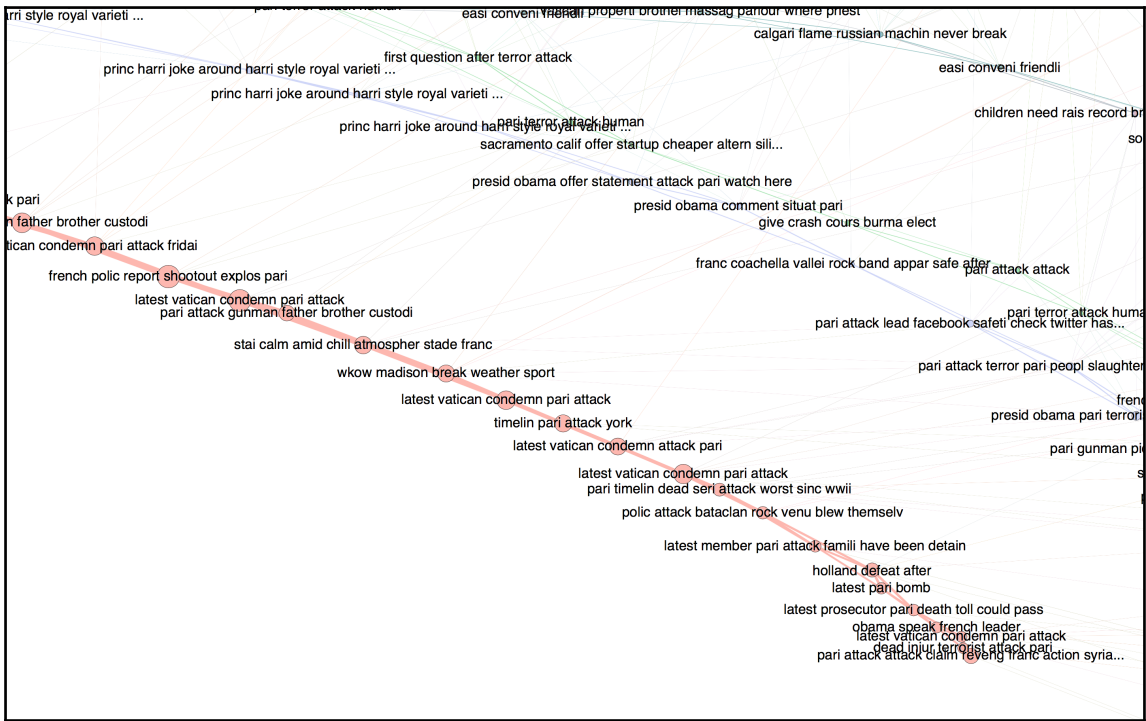
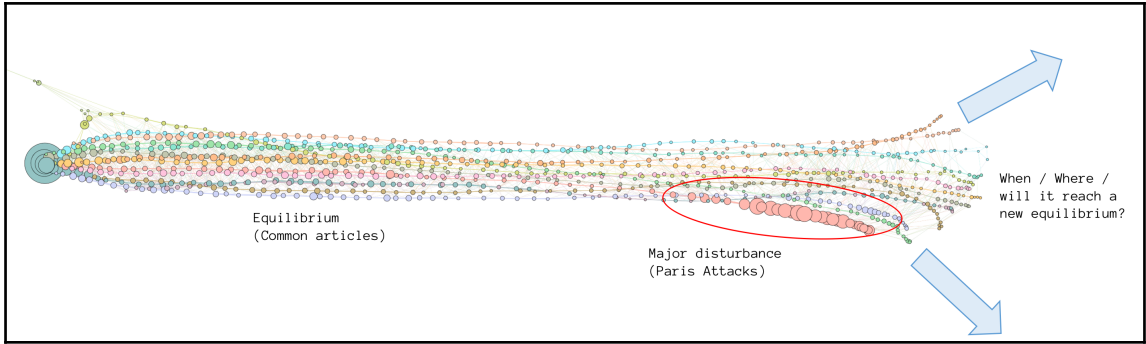




Average tone







Chapter 11: Anomaly Detection on Sentiment Analysis

```
scala> textWithEmojis
res1: String = A 🐱 and a 🐭 became friends❤️ For 🐶's birthday party, they all had 🍔, 🍟, 🍰 and 🍰

scala> EmojiUtils.shortCodify(textWithEmojis)
res2: String = A :cat:, :dog: and a :mouse: became friends:heart:. For :dog:'s birthday party, they all had :hamburger:s, :fries:s, :cookie:s and :cake:.

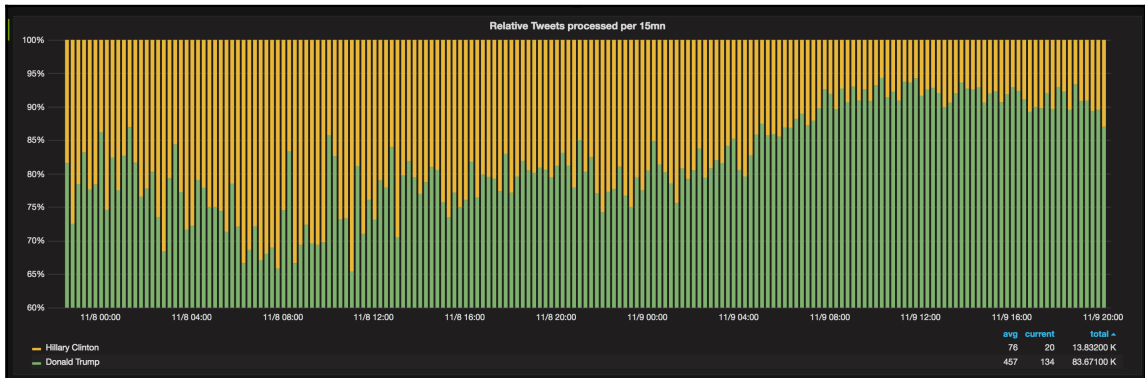
scala> EmojiUtils.removeAllEmojis(textWithEmojis)
res3: String = A , and a became friends. For 's birthday party, they all had s, s, s and .
```

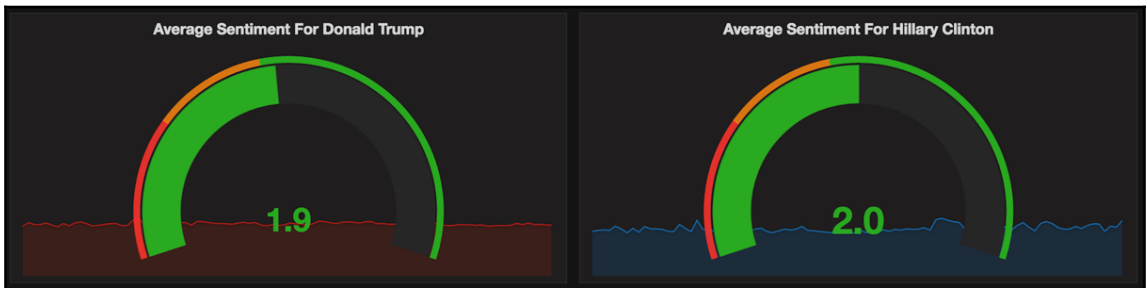
```
scala> import io.gzet.timeseries.twitter.Twitter._
import io.gzet.timeseries.twitter.Twitter._

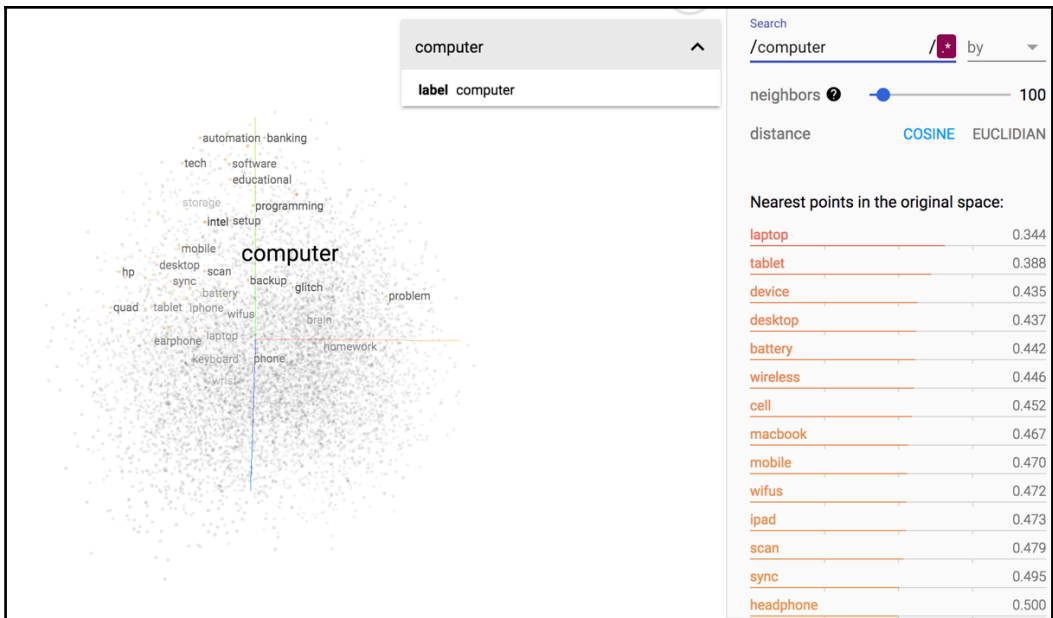
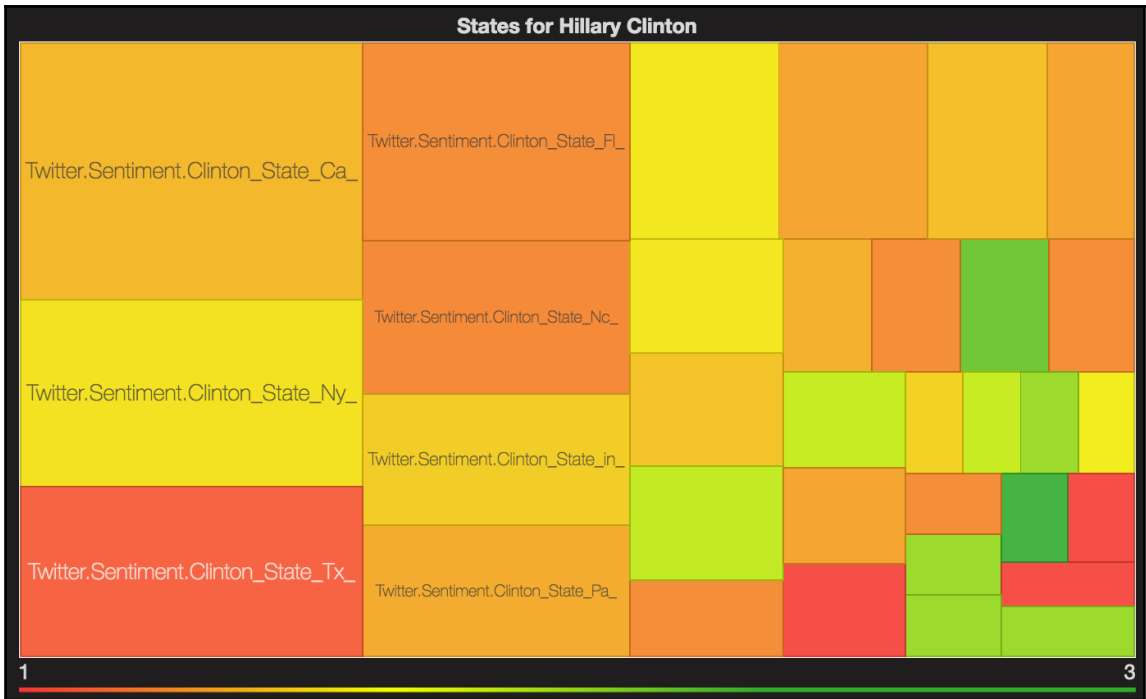
scala> println(text)
RT @johnnanjos: Michelle Obama for president in 2020 🇺🇸
Michelle Obama for president in 2020 🇺🇸
Michelle Obama for president in 2020 🇺🇸
#Not...

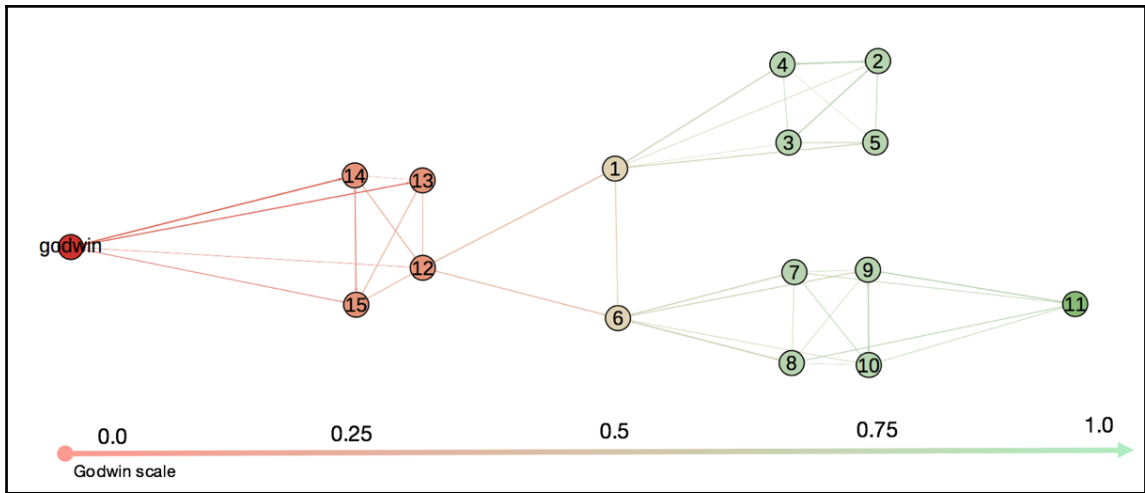
scala> println(text.clean)
michelle obama for president in 2020 michelle obama for president in 2020 michelle obama for president in 2020 #not

scala> println(text.emojis.mkString(" "))
us us us
```









$$M = USV^T$$

$$M^{2k} = US^{2k}U^T$$

$$M^{2k+1} = US^{2k+1}V^T$$

Retrieve anomalies for each cluster FINISHED

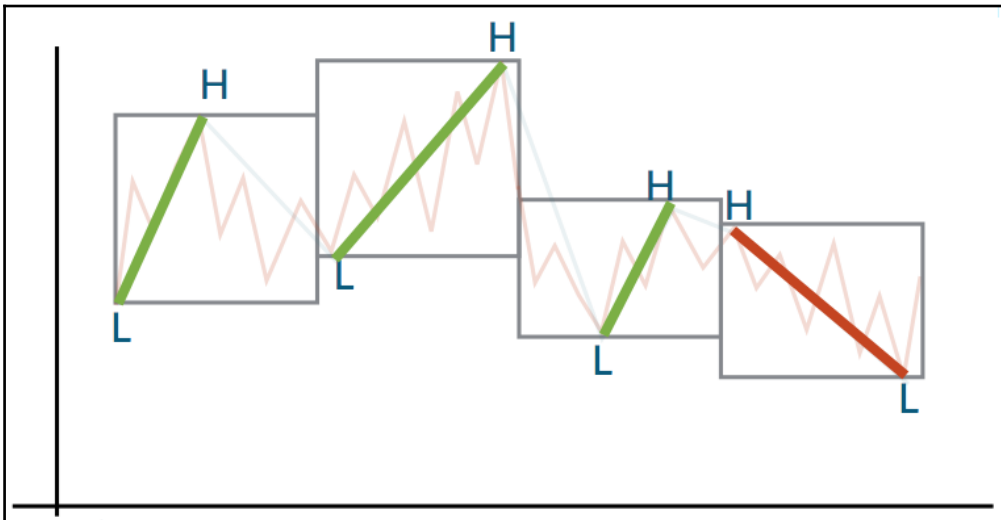
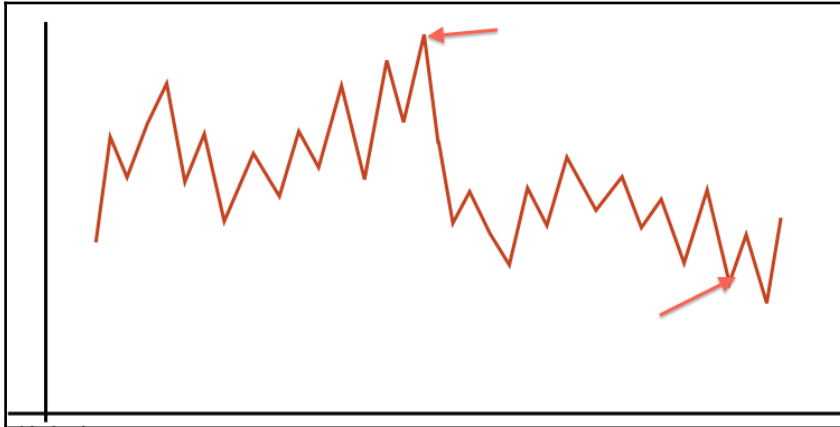
```
%sql
SELECT body, trump, clinton, love, hate, sentiment, emojis FROM tweet
WHERE cluster = ${cluster}
ORDER BY distance DESC
```

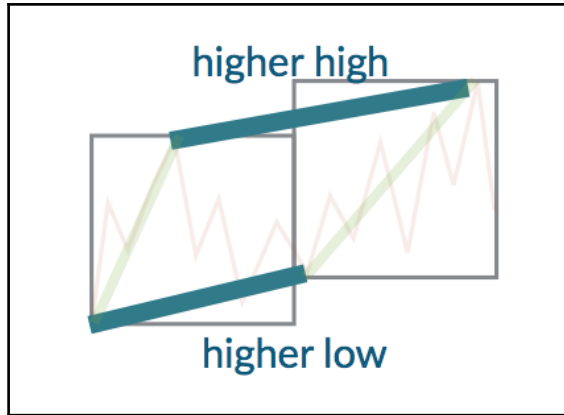
cluster: 0

body	trump	clinton	love	hate	sentiment	emojis
please negan please take trump and give we back glenn #trump #thewalkingdead	0.51166	0.45357	0.00572	0.06561	3	WrappedArray(cry, pray)
even we movie title hint the celebrated #trump age ago #achhedin	0.44397	0.37753	-0.01636	-0.00007	3	WrappedArray(smiley)
soydeandaya #soyandayer good morning friend #happyday #electionday	0.4067	0.38863	0.11563	0.06771	3	WrappedArray(grinning)
i love the fact that #toblerone be trend above #uselection where i live #lovegreatbritain	0.47246	0.39321	0.13134	0.09768	3	WrappedArray(joy)
get we good champagne glass out towel and the brut #maga x enjoy the show #schadenfreude all around the #msm	0.51371	0.39599	0.08136	0.07178	3	WrappedArray(sunglasses, sweat_smile, joy)
well have a permanente patime job for the next year #electionfinalthoughts #electionnight #heartbroken	0.58486	0.52501	0.0372	0.07604	3	WrappedArray(broken_heart)
guy so what be the back up plan if #drump win #electionnight #imwithher	0.49184	0.54833	0.1178	0.10663	3	WrappedArray(sob)
to have a super power nation with a leader that do not believe in #climatechange fear for the future of #biodiversity #uselection	0.26681	0.25984	0.03077	0.07656	1.5	WrappedArray(broken_heart)
nevada latino ur beautiful #bluewave grow into a huge tidal wave #voteblueballot #votehrc #strongertogether	0.52024	0.52676	0.02254	0.02089	3	WrappedArray(clap, clap, ocean, ocean, ocean)

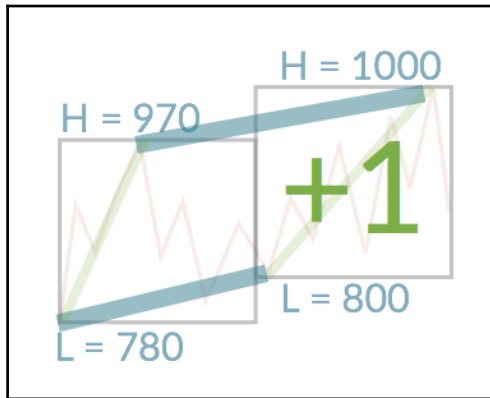
Took 2 sec. Last updated by anonymous at December 18 2016, 1:08:42 AM. (outdated)

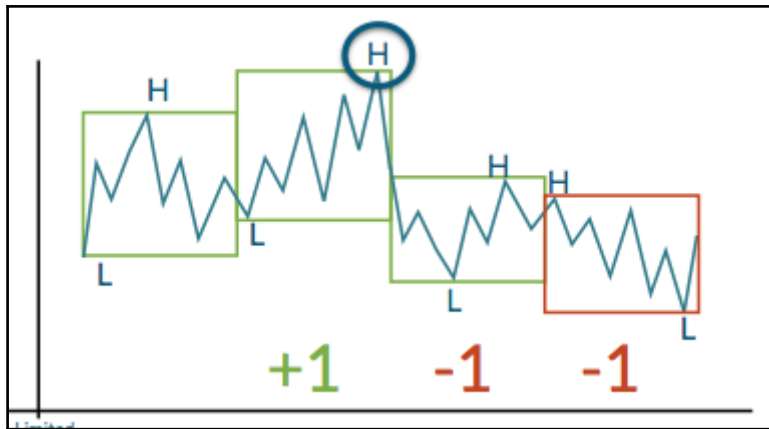
Chapter 12: TrendCalculus

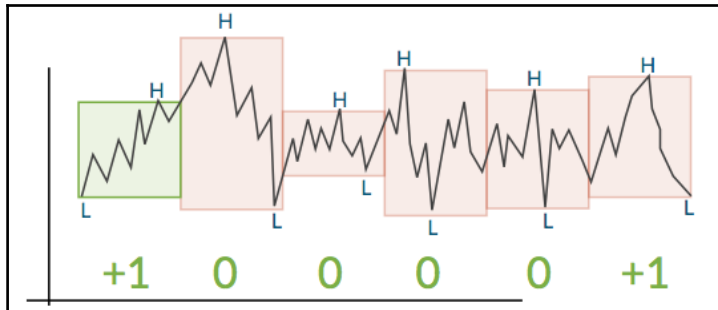
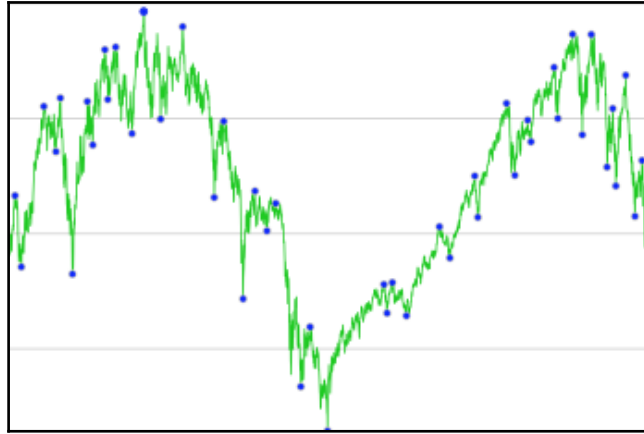
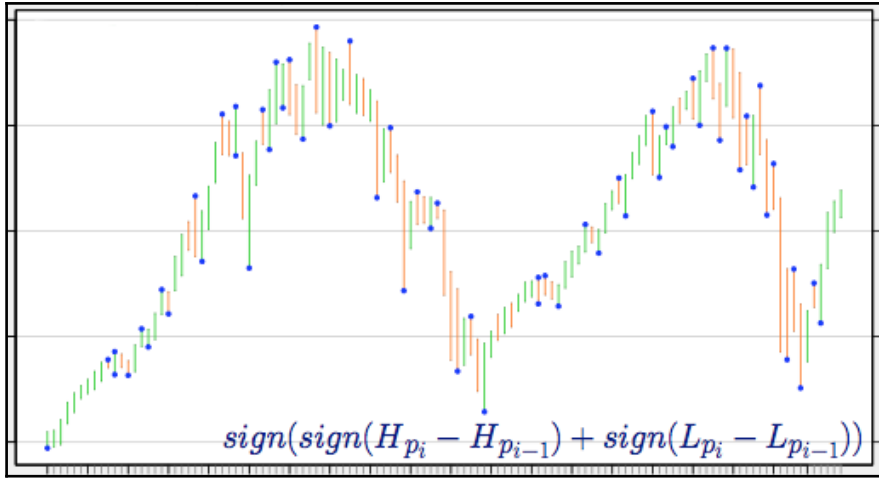


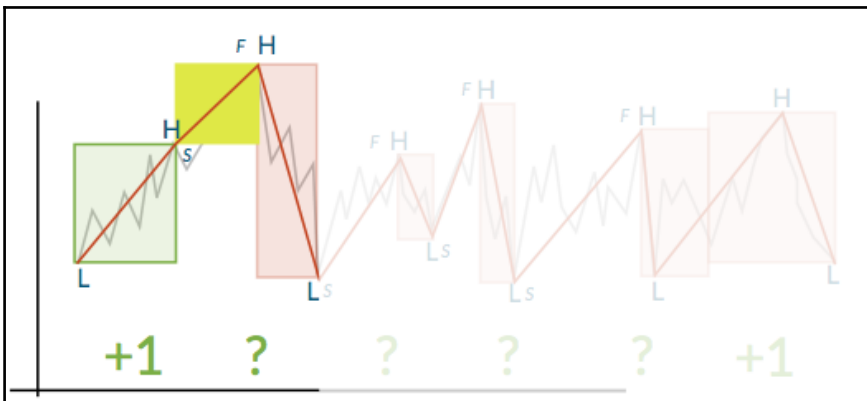
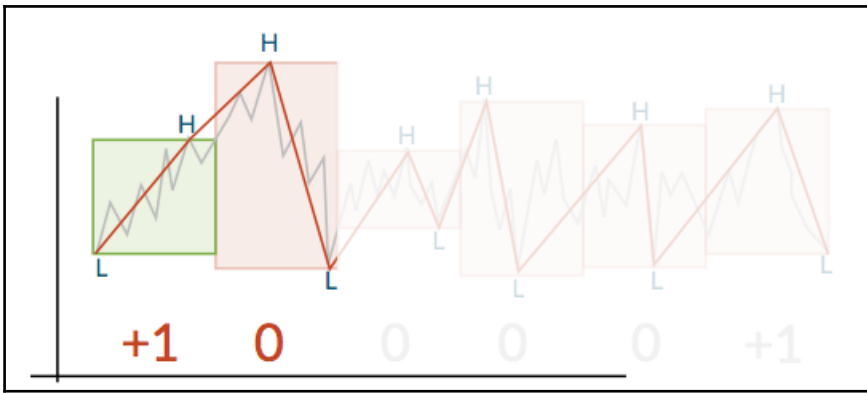
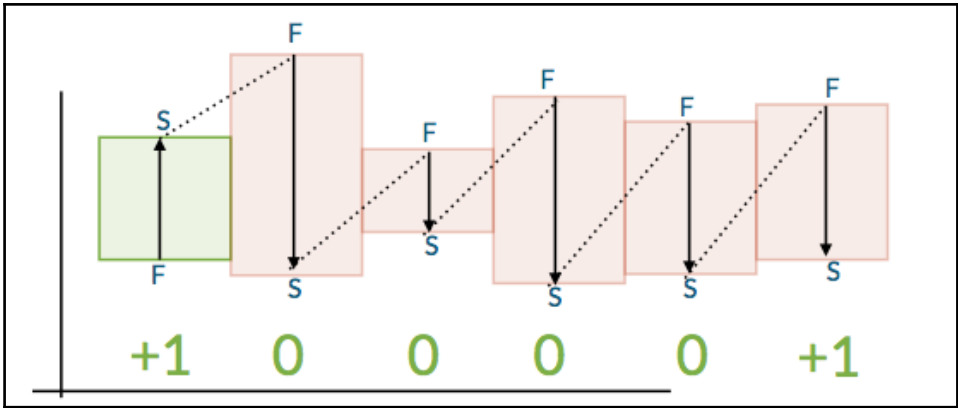


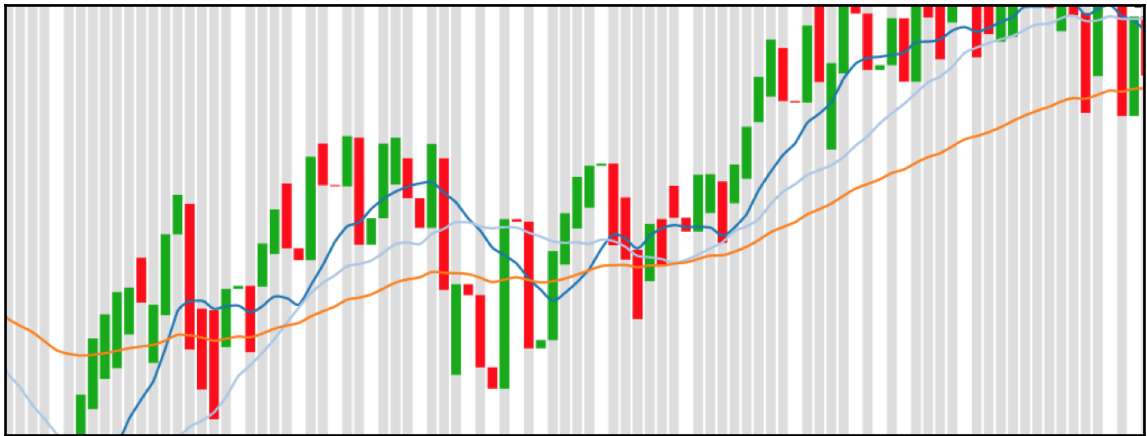
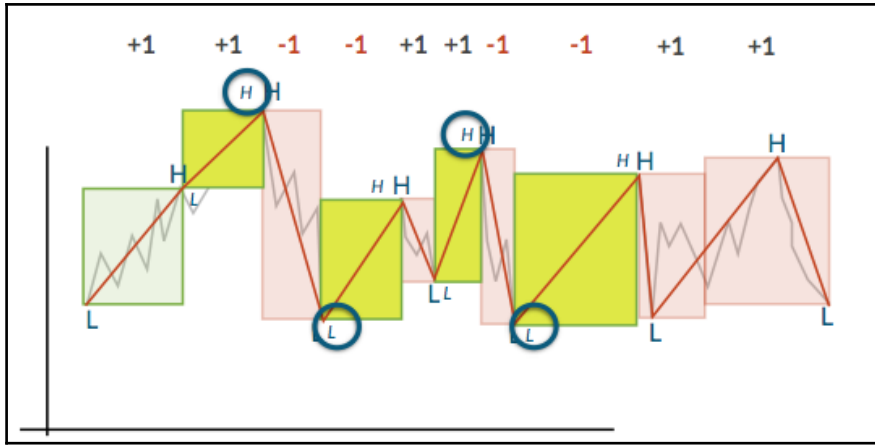
$$\text{sign}(\text{sign}(H_{p_i} - H_{p_{i-1}}) + \text{sign}(L_{p_i} - L_{p_{i-1}}))$$



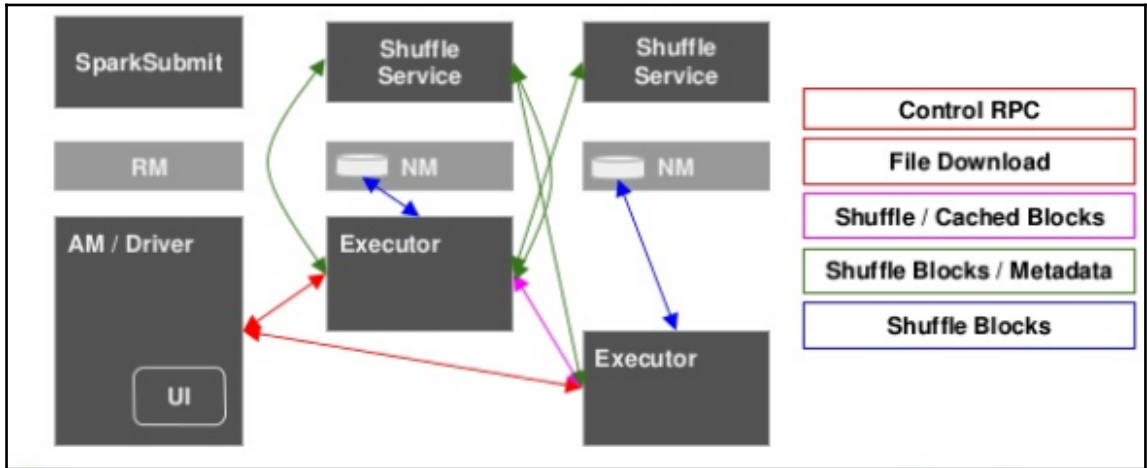




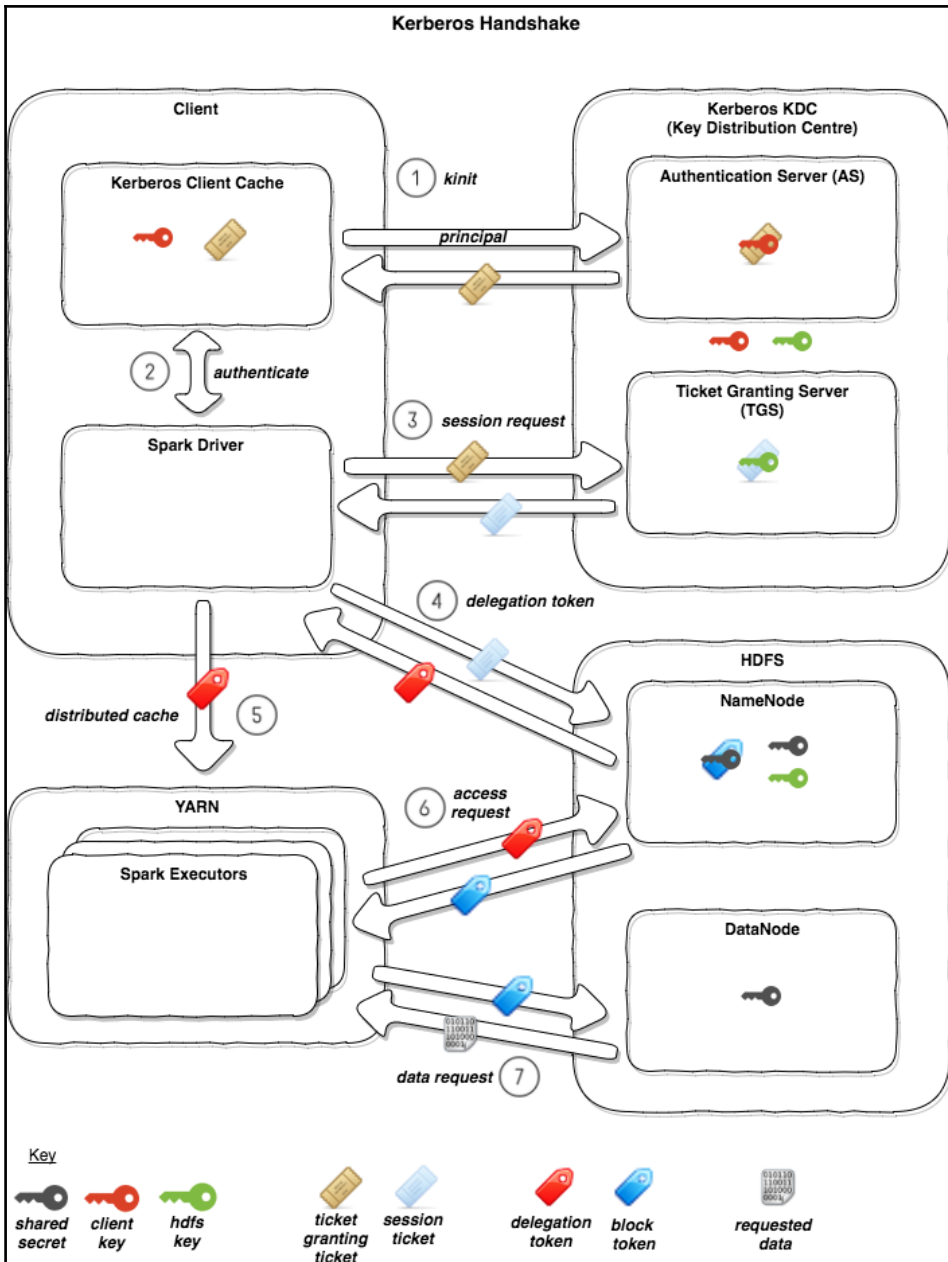


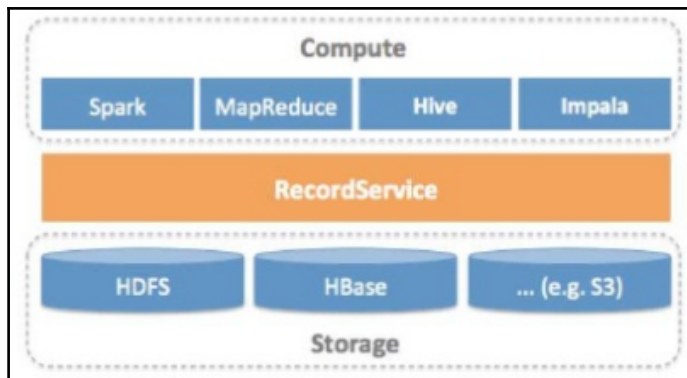


Chapter 13: Secure Data

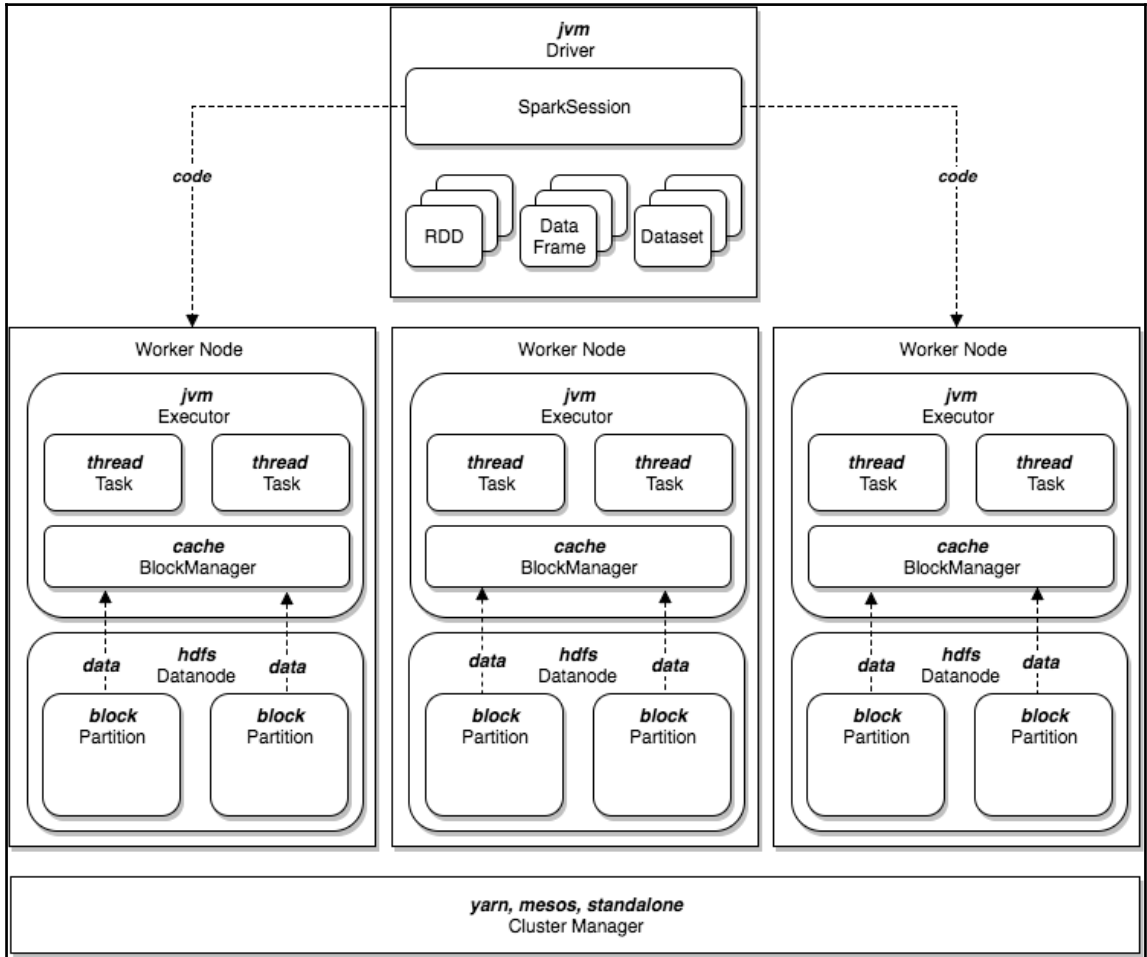


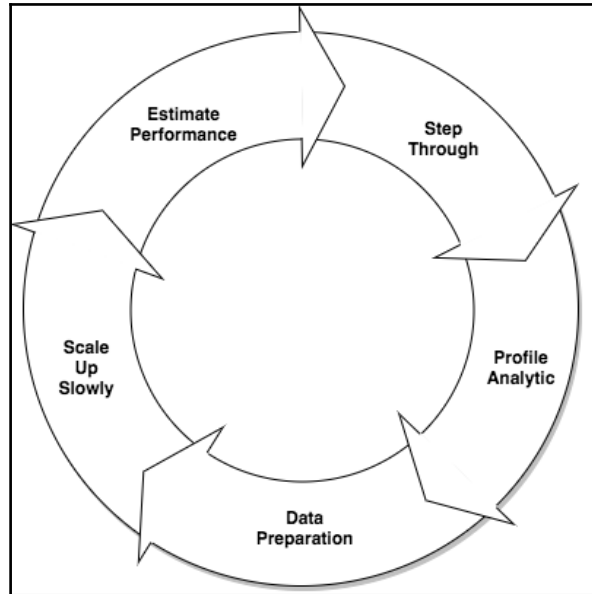
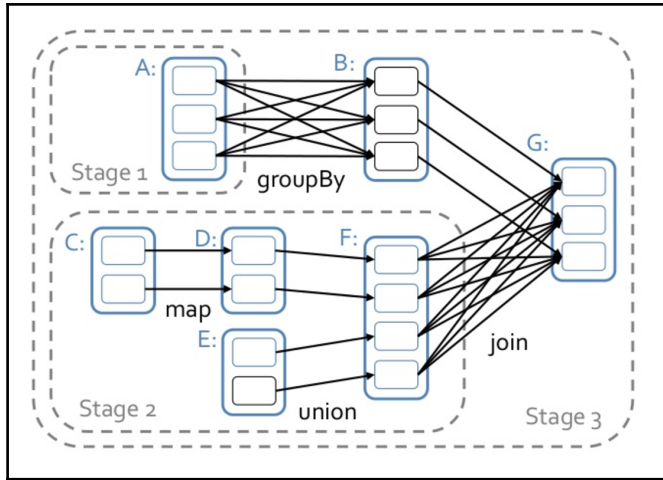
Kerberos Handshake





Chapter 14: Scalable Algorithms





Summary Metrics for 168 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	71 ms	7 s	12 s	13 s	17 s
Scheduler Delay	1 ms	2 ms	3 ms	3 ms	5 ms
Task Deserialization Time	0 ms	1 ms	1 ms	2 ms	24 ms
GC Time	3 ms	0.2 s	0.4 s	0.4 s	0.6 s
Result Serialization Time	0 ms	0 ms	0 ms	0 ms	1 ms
Getting Result Time	0 ms	0 ms	0 ms	0 ms	0 ms
Peak Execution Memory	0.0 B	0.0 B	0.0 B	0.0 B	0.0 B
Input Size / Records	119.3 KB / 11	16.5 MB / 1397	32.1 MB / 2576	32.1 MB / 2726	35.2 MB / 2869
Shuffle Write Size / Records	1591.0 B / 32	61.9 KB / 6445	98.4 KB / 11595	104.9 KB / 12861	120.2 KB / 17005

Aggregated Metrics by Executor


Executor ID ▲	Address	Task Time	Total Tasks	Failed Tasks	Succeeded Tasks	Input Size / Records	Shuffle Write Size / Records
0	CANNOT FIND ADDRESS	14 min	88	0	88	2035.1 MB / 170712	6.7 MB / 809037
1	CANNOT FIND ADDRESS	14 min	80	0	80	2.0 GB / 174489	6.9 MB / 832305

Tasks (168)

Page: >

2 Pages. Jump to . Show items in a page.

Index ▲	ID	Attempt	Status	Locality Level	Executor ID / Host	Launch Time	Duration	Scheduler Delay	Task Deserialization Time	GC Time	Result Serialization Time	Getting Result Time	Peak Execution Memory	Input Size / Records	Write Time	Shuffle Write Size / Records	Errors
0	336	0	SUCCESS	PROCESS_LOCAL	1 / 192.168.1.67	2016/12/22 06:07:19	12 s	4 ms	24 ms	0.4 s	0 ms	0 ms	0.0 B	32.1 MB / 2644	48 ms	104.9 KB / 15620	
1	337	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.1.67	2016/12/22 06:07:19	7 s	4 ms	24 ms	0.3 s	0 ms	0 ms	0.0 B	15.7 MB / 1286	44 ms	56.4 KB / 6727	
2	338	0	SUCCESS	PROCESS_LOCAL	1 / 192.168.1.67	2016/12/22 06:07:19	12 s	4 ms	23 ms	0.4 s	0 ms	0 ms	0.0 B	32.1 MB / 2627	45 ms	101.6 KB / 13084	
3	339	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.1.67	2016/12/22 06:07:19	7 s	3 ms	24 ms	0.3 s	0 ms	0 ms	0.0 B	17.0 MB / 1397	37 ms	62.1 KB / 7047	
4	340	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.1.67	2016/12/22 06:07:26	12 s	3 ms	2 ms	0.4 s	0 ms	0 ms	0.0 B	32.1 MB / 2720	35 ms	105.4 KB / 13534	
5	341	0	SUCCESS	PROCESS_LOCAL	0 / 192.168.1.67	2016/12/22 06:07:27	6 s	3 ms	1 ms	0.2 s	0 ms	0 ms	0.0 B	15.5 MB / 1306	39 ms	57.5 KB / 6523	
6	342	0	SUCCESS	PROCESS_LOCAL	1 / 192.168.1.67	2016/12/22 06:07:32	13 s	3 ms	1 ms	0.4 s	0 ms	0 ms	0.0 B	32.1 MB / 2717	35 ms	110.0 KB / 13248	


Jobs
Stages
Storage
Environment
Executors
SQL
Spark shell application UI

Thread dump for executor 1

Updated at 2016/12/22 06:00:41

[Expand All](#)

Search:


Thread ID	Thread Name	Thread State
58	Executor task launch worker-0	RUNNABLE
59	Executor task launch worker-1	RUNNABLE

```

org.json.simple.JSONValue.escape(Unknown Source)
org.json.simple.JSONValue.escape(Unknown Source)
org.json.simple.JSONValue.toJSONString(Unknown Source)
org.json.simple.JSONObject.toJSONString(Unknown Source)
org.json.simple.JSONObject.toJSONString(Unknown Source)
org.json.simple.JSONObject.toJSONString(Unknown Source)
org.json.simple.JSONObject.toJSONString(Unknown Source)
org.json.simple.JSONArray.toJSONString(Unknown Source)
org.json.simple.JSONArray.toJSONString(Unknown Source)
org.json.simple.JSONValue.toJSONString(Unknown Source)
org.json.simple.JSONObject.toJSONString(Unknown Source)
org.json.simple.JSONObject.toJSONString(Unknown Source)
org.json.simple.JSONObject.toJSONString(Unknown Source)
io.gzelt.util.gdelt.GkgJSON.parseV2ToJson(GkgJSON, java:248)
io.gzelt.util.spark.gdelt.GkgParser$.toJsonGkgV2(GkgParser, scala:23)
$line20.$read$$iw$$iw$$iw$$iw$$iw$$iw$$iw$$iw$$iw$anonfun$1.apply(<<console>:31)
$line20.$read$$iw$$iw$$iw$$iw$$iw$$iw$$iw$$iw$$iw$anonfun$1.apply(<<console>:31)
scala.collection.Iterator$$anon$11.next(Iterator, scala:409)
scala.collection.Iterator$$anon$11.next(Iterator, scala:409)
scala.collection.Iterator$$anon$12.nextCur(Iterator, scala:434)
scala.collection.Iterator$$anon$12.hasNext(Iterator, scala:440)
scala.collection.Iterator$$anon$11.hasNext(Iterator, scala:408)
org.apache.spark.shuffle.sort.ByPassMergeSortShuffleWriter.write(ByPassMergeSortShuffleWriter, java:147)
org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask, scala:79)
org.apache.spark.scheduler.ShuffleMapTask.runTask(ShuffleMapTask, scala:47)
org.apache.spark.scheduler.Task.run(Task, scala:86)
org.apache.spark.executor.Executor$TaskRunner.run(Executor, scala:274)
java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor, java:1142)
java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor, java:617)
java.lang.Thread.run(Thread, java:745)

```

31	dispatcher-event-loop-0	WAITING
----	-------------------------	---------


Jobs
Stages
Storage
Environment
Executors
SQL
Spark shell application UI

Details for Stage 0 (Attempt 0)

Total Time Across All Tasks: 18 min

Locality Level Summary: Process local: 114

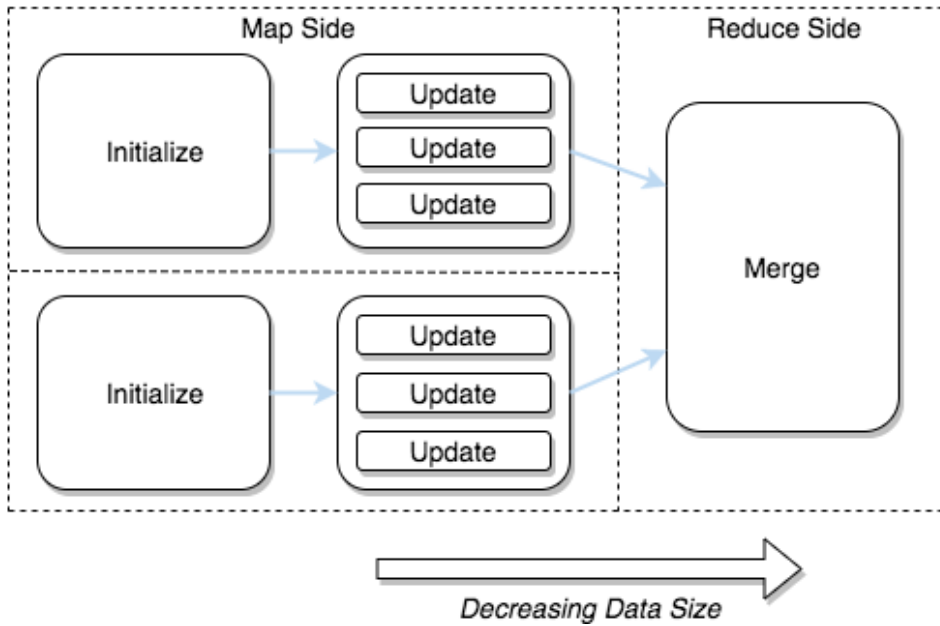
Input Size / Records: 2.6 GB / 216023

Shuffle Write: 8.7 MB / 1070998

- [DAG Visualization](#)
- [Show Additional Metrics](#)
- [Event Timeline](#)
- Enable zooming

Task ID	Host	Scheduler Delay	Task Deserialization Time	Shuffle Read Time	Executor Computing Time	Shuffle Write Time	Result Serialization Time	Getting Result Time
0 / 192.168.1.67	192.168.1.67	~10ms	~10ms	~10ms	~10ms	~10ms	~10ms	~10ms
1 / 192.168.1.67	192.168.1.67	~10ms	~10ms	~10ms	~10ms	~10ms	~10ms	~10ms

Summary Metrics for 110 Completed Tasks



Spark 2.0.1 Jobs Stages Storage Environment Executors SQL Spark shell application UI

Storage

RDDs

RDD Name	Storage Level	Cached Partitions	Fraction Cached	Size in Memory	Size on Disk
*BatchedScan parquet [url#144,body#145] Format: ParquetFormat, InputPaths: file:/data/gdel/gkg/enth.parquet, PartitionFilters: [], PushedFilters: [], ReadSchema: struct<url:string,body:string>	Memory Deserialized 1x Replicated	2	100%	37.9 MB	0.0 B

