Brandwatch

# Detecting events on the Web in real-time with Java, Kafka & ZooKeeper

Dr. James Stanier  |  brandwatch.com  |  jamess@brandwatch.com

# Coming Up/

- Me, Brandwatch and new problems

- Moving to Kafka

- Processing data

- Distributing work

- Finding meaning

# Who?

**Dr. James Stanier**

VP Engineering, Product (Backend) Brandwatch
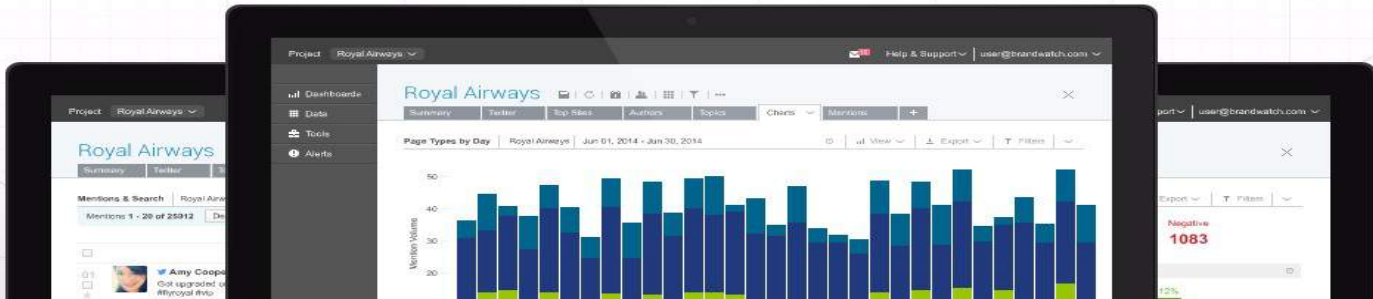
@jstanier | jamess@brandwatch.com

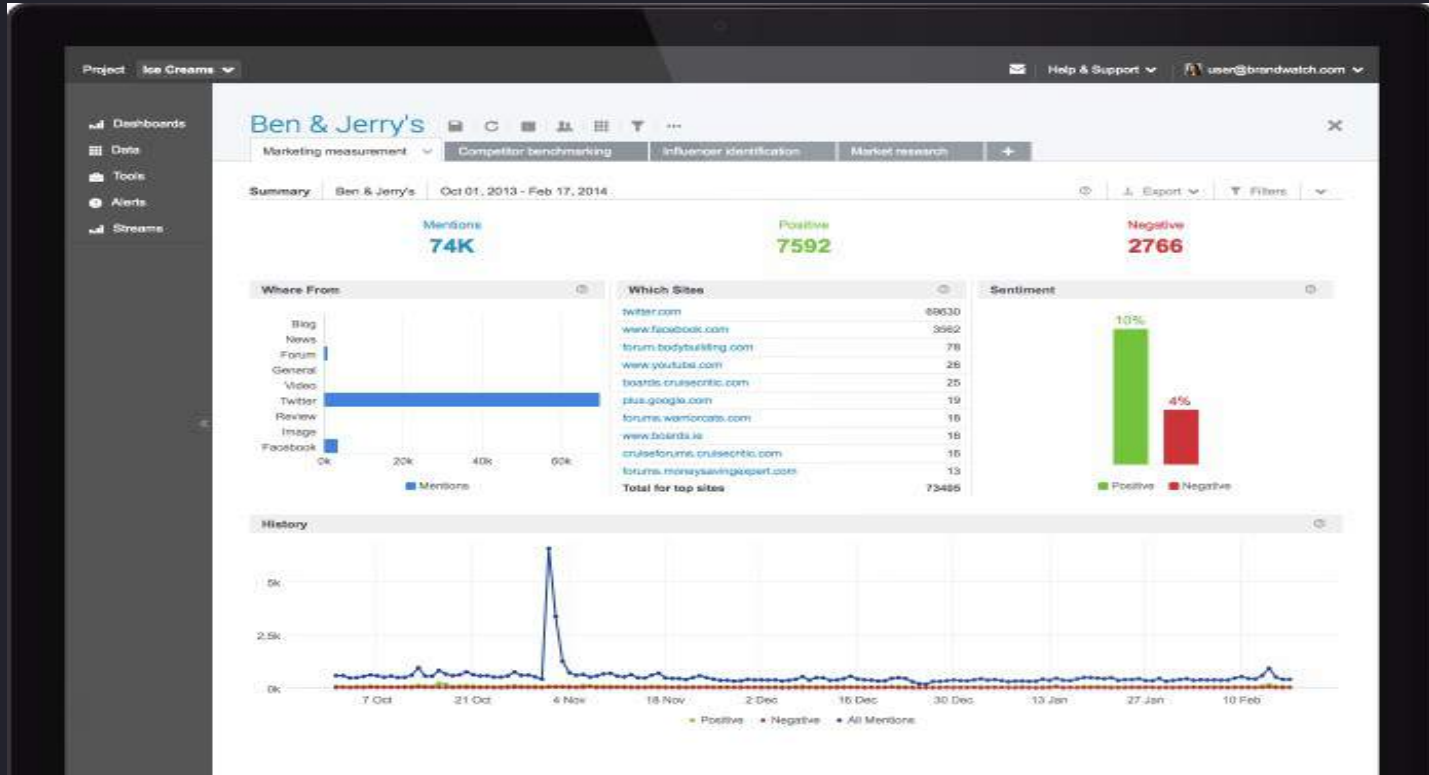# Data/ Presentation

# Data/ Aggregation

# Data/ Classification

# Data/ Not just top level metrics

# Data/ The numbers

- 50+ Java Web Crawlers

- 10+ Historical crawlers for new queries

- Twitter via GNIP (now Twitter), Weibo, Disqus and more

- 80M+ query matches per day

# A new challenge

# The challenge/ The signal from the noise

# The challenge/ at scale

- 130K+ user queries

- 80M+ mentions per day

- Polling the data stores for mentions for all queries takes 8hrs for one pass

# The Problem/ How we handled it…

Crawler 1

Crawler 2

Crawler n-1

Crawler N

Mentions

**Kafka**

Signals

Mentions

**Signals** Processing cluster

Signals

**Signals** grouper

DB

# Kafka

# Step 1/ Kafka

Crawler 1

Crawler 2

Crawler n-1

Crawler N

Mentions

**Kafka**
Cluster

# Kafka/ What is it?

- Apache Kafka is a publish-subscribe messaging system rethought as a distributed commit log

- Apache top level project November 2013

- Started at LinkedIn

# Kafka/ is…

- **Fast:** hundreds of MBs read/write per second from thousands of clients

- **Scalable:** clustered, partitioned over many machines, expanded without downtime

- **Durable:** messages persisted to disk and replicated in cluster

# Kafka/ Written to disk?



**FIGURE 3**

**Comparing Random and Sequential Access in Disk and Memory**

| | |
|---|---|
| Random, disk | 316 values/sec |
| Sequential, disk | 53.2M values/sec |
| Random, SSD | 1924 values/sec |
| Sequential, SSD | 42.2M values/sec |
| Random, memory | 36.7M values/sec |
| Sequential, memory | 358.2M values/sec |

Scale: 10, 100, 1000, $10^4$, $10^5$, $10^6$, $10^7$, $10^8$

Note: Disk tests were carried out on a freshly booted machine (a Windows 2003 server with 64-GB RAM and eight 15,000-RPM SAS disks in RAID5 configuration) to eliminate the effect of operating-system disk caching. SSD test used a latest-generation Intel high-performance SATA SSD.

http://q.acm.org/detail.cfm?id=1563874

# Kafka/ Bending, not breaking



Load Test - Tweets Per Minute For 1 Hour

Publisher Volume On Edge

Inbound Volume Internal component

http://engineering.gnip.com/tag/kafka/

# Kafka/ Sending from the crawlers

```
String message = toJson(...);

KeyedMessage<String, String> message = new
    KeyedMessage<String, String>("query.mentions",
queryId, message);

producer.send(message);
```

# Step 1/ Done

Crawler 1

Crawler 2

Crawler
n-1

Crawler
N

Mentions

**Kafka**
Cluster

# Processing

# Processing/ What's happening now?

# Step 2.1/ One processing JVM

Crawler 1

Crawler 2

Crawler n-1

Crawler N

Mentions

**Kafka**
Cluster

Signals
processor

# Processing/ A wild tweet appears!

Mention

date: 01/06/2015 16:05

pageType: twitter

author: @berlinperson

hashtags: [#berlinbuzzwords, #amazingtalk, #greatshoes]

mentionedTweeters: [@jstanier]

text: "@jstanier is at #berlinbuzzwords #amazingtalk #greatshoes"

# Processing/ Storing hashtags

`Map<Date, Multiset<String>>`

Initialise with the last 24 hours

# Processing/ Storing hashtags

`Map<Date, Multiset<String>>`

Mention

date: 01/06/2015 4:10PM

hashtags: [#berlinbuzzwords, #amazingtalk, #greatshoes]

# Processing/ Storing hashtags

`Map<Date, Multiset<String>>`

Mention

date: 01/06/2015 4:10PM

hashtags: [#berlinbuzzwords, #amazingtalk, #greatshoes]

```
add("#berlinbuzzwords")
add("#amazingtalk")
add("#greatshoes")
```

# Processing/ Cycling the buckets

```java
@Scheduled(cron = "0 0 * * * *")
public void cycleBuckets() {

    Date oldest = buckets.lastKey();

    removeBucket(oldest);

    DateTime newest = new

      DateTime(buckets.firstKey());

    addBucket(newest.plusHours(1).toDate());

  }
```

# Processing/ Detecting spikes

- At scheduled intervals

- For each #

  - Convert to a timeseries [5, …. 1002, 5499]
  - Compare previous hour to history
  - Give a score to it

- If score > threshold, it's interesting

- Send it on a new Kafka topic

# Processing/ What we just did

#hashtag
data
model

# Processing/ But we also track…

| #hashtag data model | country data model | author data model | page type data model | sentiment data model | volume data model | link share data model |
|---|---|---|---|---|---|---|

# Processing/ …for one query

**"Berlin Buzzwords"** query

| #hashtag data model | country data model | author data model | page type data model | sentiment data model | volume data model | link share data model |
|---|---|---|---|---|---|---|

# Processing/ 100K+ queries and rising

# Processing/ We need more JVMs

But how do we share the workload?

# Distribution of work

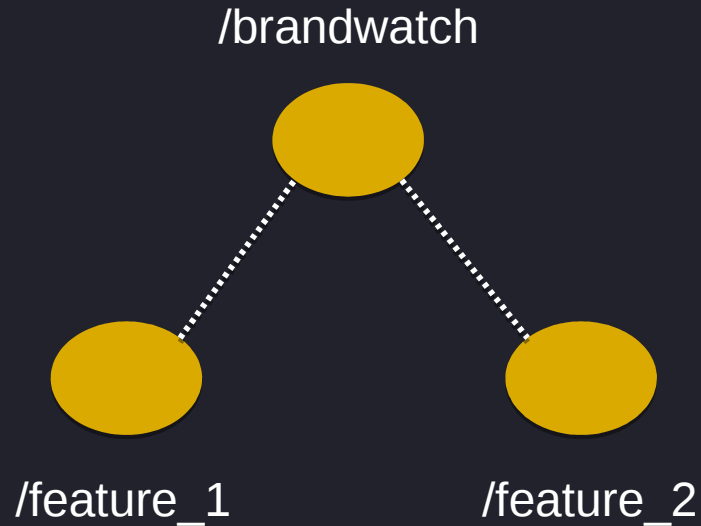# Distribution/ An atomic unit of work

?

Signals
Processing
cluster

# Distribution/ Leader election

A way of deciding who is the leader for a task in a group of distributed nodes

# Distribution/ Zookeeper

A way of coordinating and managing distributed applications

# Zookeeper/ It's like a file system

/brandwatch

/feature_1          /feature_2

# Zookeeper/ At the command line

```
[zk: localhost:12181(CONNECTED) 6] ls /
[zookeeper, admin, consumers, brandwatch, controller, brokers, controller_epoch]
[zk: localhost:12181(CONNECTED) 7] ls /brokers
[topics, ids]
[zk: localhost:12181(CONNECTED) 8]
```

# Distribution/ Recipes

# Distribution/ Offering jobs



/brandwatch

/signals

/queries

/15846

/1268589

**Manager** JVM

DB

# Distribution/ PGQ

# Distribution/ Leader election 101

/brandwatch

/signals

/queries

/15846

/1268589

Processing JVM 1

Processing JVM 2

Processing JVM 3

# Distribution/ Leader election 101



/brandwatch

/signals

/queries

/15846

/1268589

Processing JVM 1

Processing JVM 2

Processing JVM 3

1    2    3

# Distribution/ The leader dies

/brandwatch

/signals

/queries

/15846

/1268589

2    3

Processing JVM

Processing JVM

# Distribution/ The dead rises again



/brandwatch

/signals

/queries

/15846

/1268589

Processing JVM

Processing JVM

Processing JVM

2    3    4

# Distribution/ That's how we do it

Each time someone turns on the feature, we

leader elect for processing

# Distribution/ Almost there?

We are processing long running jobs

What about workers getting overloaded?

# Distribution/ After leader election

1. Take leadership

2. Hit max queries?

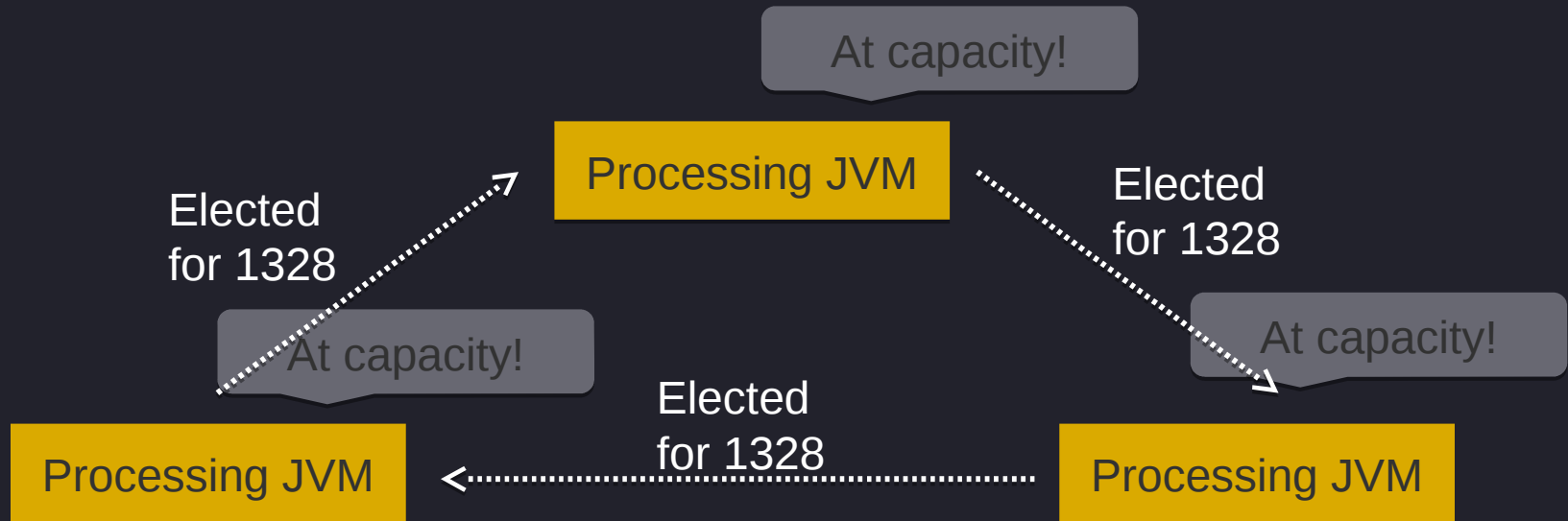   a. No – go to 3

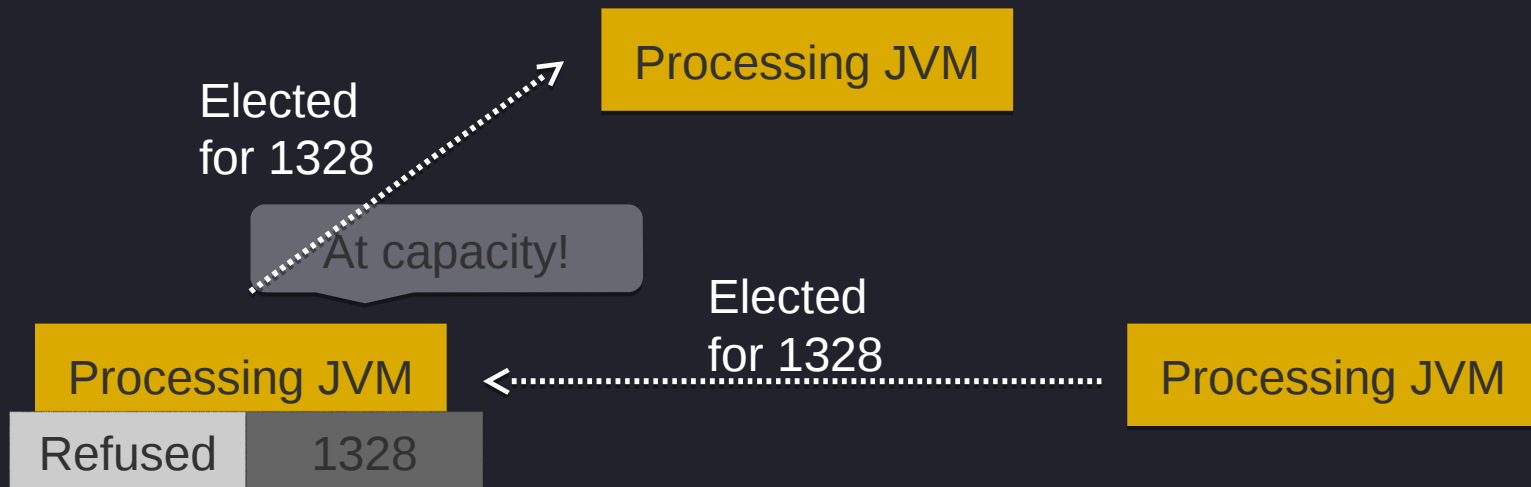   b. Yes – give up leadership, try again

3. Start working

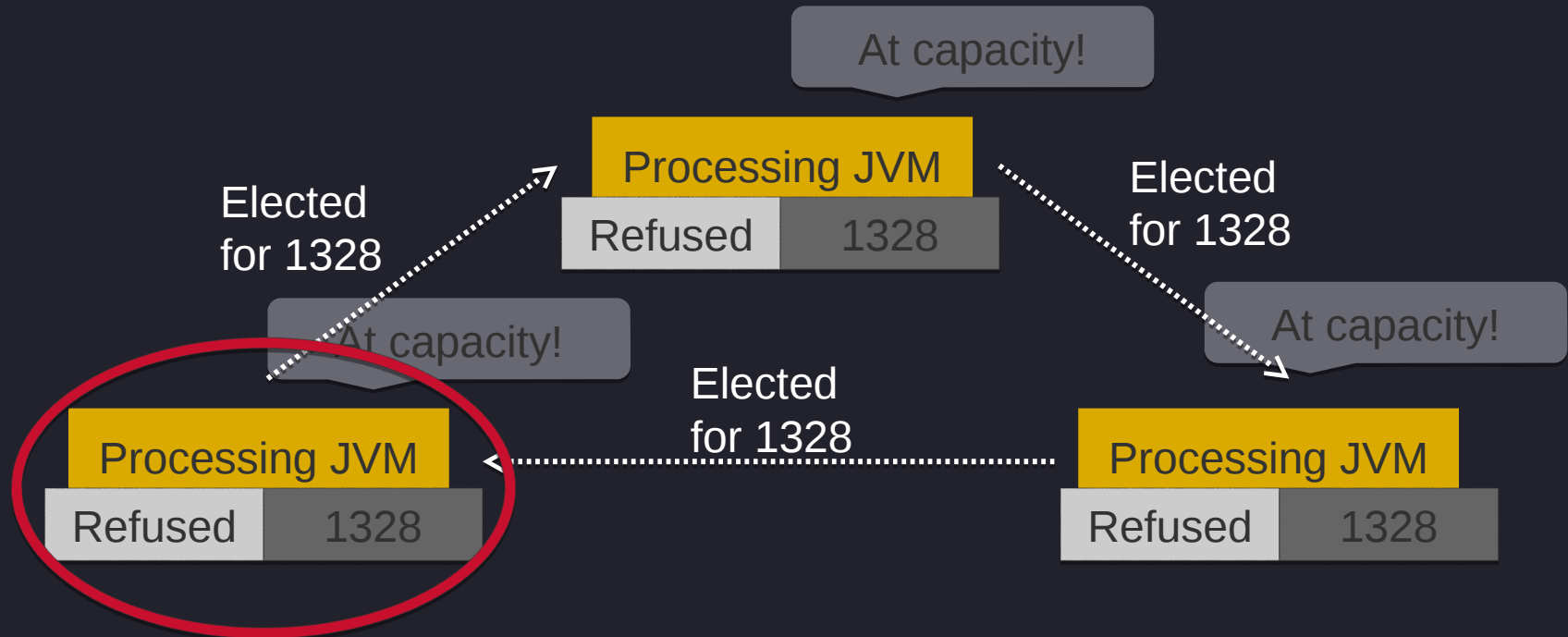# Distribution/ Now we're almost there?

Actually, no…

# Distribution/ Infinite election

# Distribution/ Solution

# Distribution/ Solution
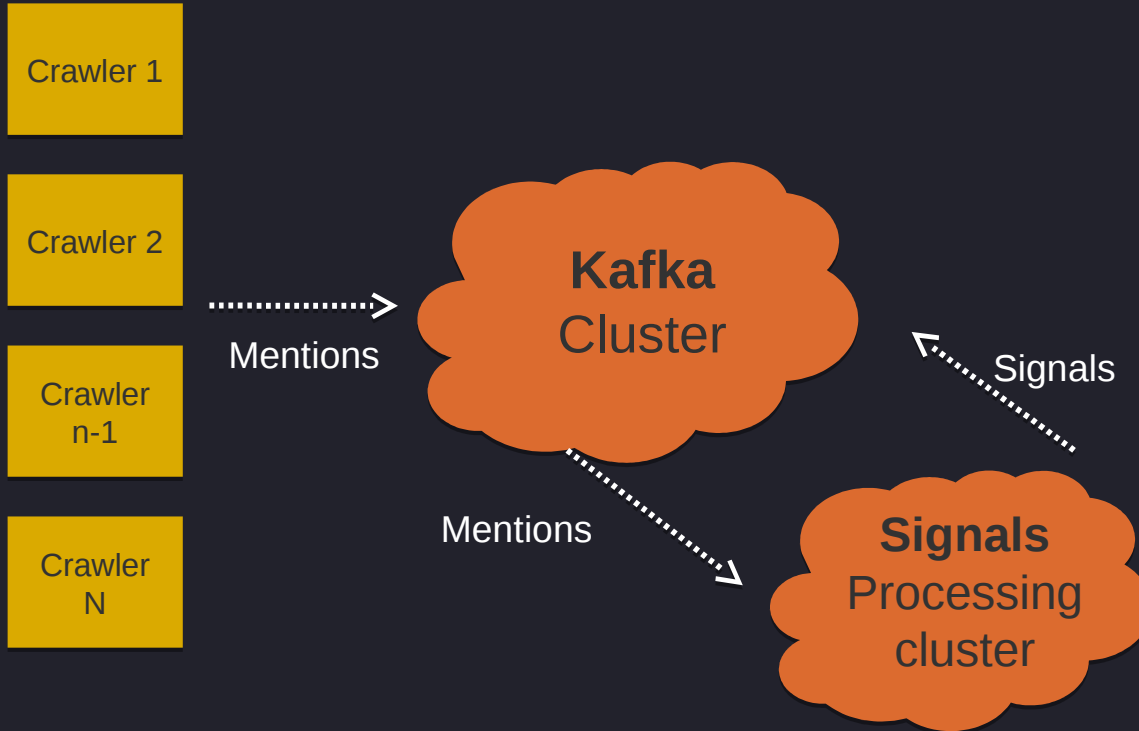
State

# State/ Snapshotting of worker data

If one worker dies, we want the other to pick up where it left off
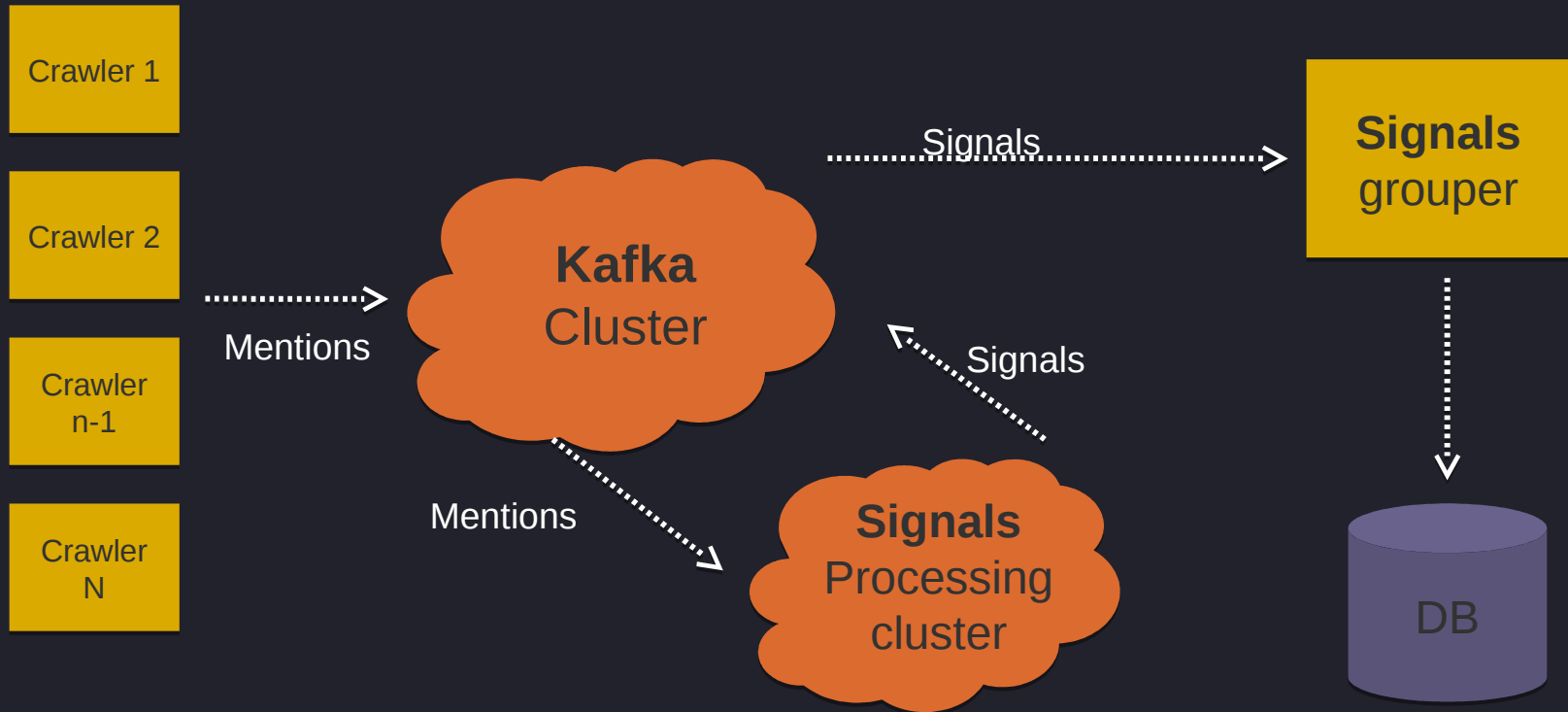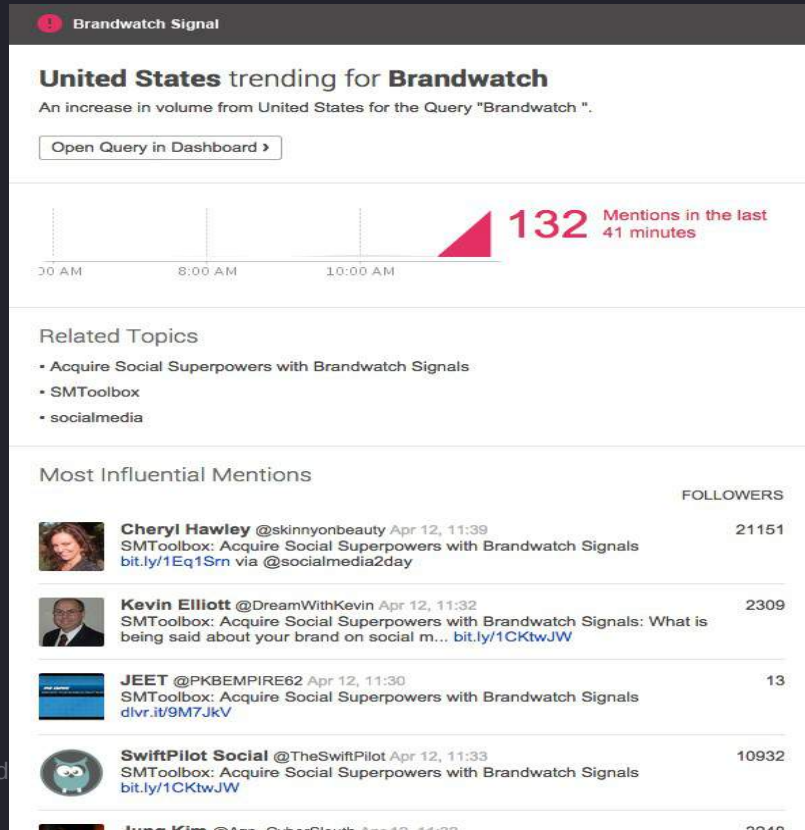
Regular snapshotting to HBase

# Step 2.2/ Done!

Crawler 1

Crawler 2

Crawler n-1

Crawler N

**Kafka** Cluster

Mentions

Mentions

**Signals** Processing cluster

Signals

# Finding meaning

# Step 3/ Meaning

Crawler 1

Crawler 2

Crawler n-1

Crawler N

Mentions

**Kafka** Cluster

Mentions

Signals

**Signals** Processing cluster

Signals

**Signals** grouper

DB

# Meaning/ Desired outcome

# Meaning/ 1. Topics

#bbuzz - - - - - - - - - - - - - - - > DB

{ "Hey @jstanier Your talk sucks! #bbuzz",
"I love #bbuzz",
"Where's the Club Mate? #bbuzz",

… ,
"I have consumed ALL THE FREE COFFEE!

#bbuzz" }

# Meaning/ 1. Topics
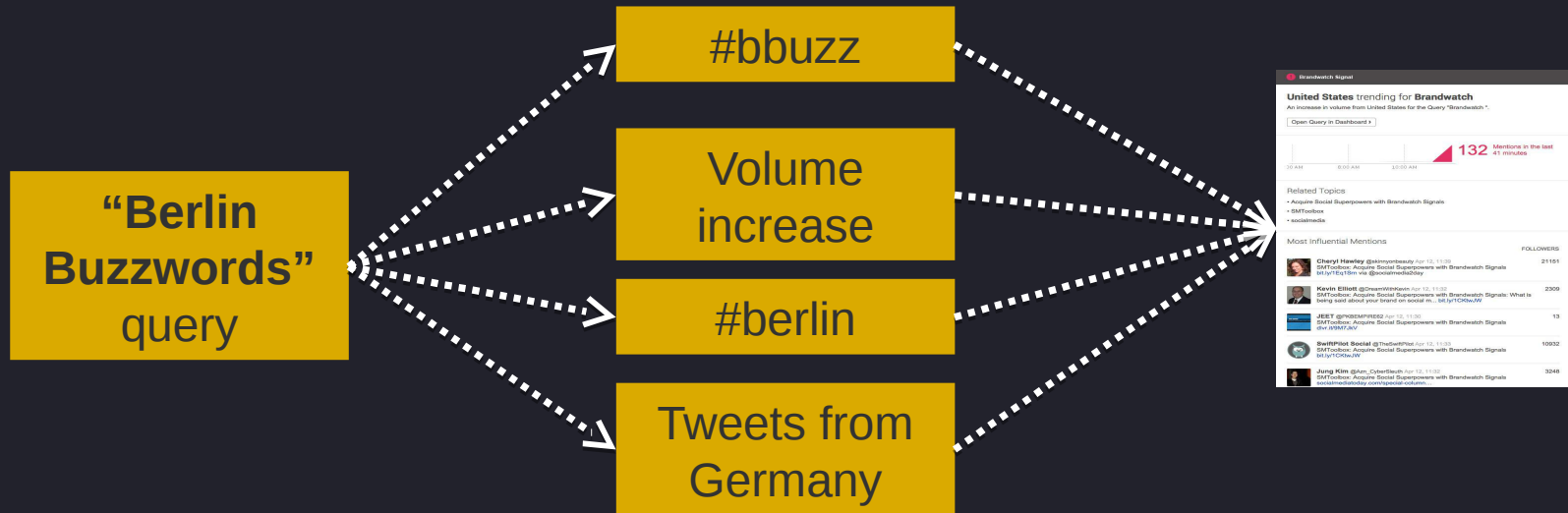
#bbuzz  - - - - - - - - - - - - - - - - - - >  DB

{

 Berlin Buzzwords
ElasticSearch
Scaling

}

# Meaning/ 2. Grouping events

"Berlin Buzzwords" query

#bbuzz

Volume increase

#berlin

Tweets from Germany
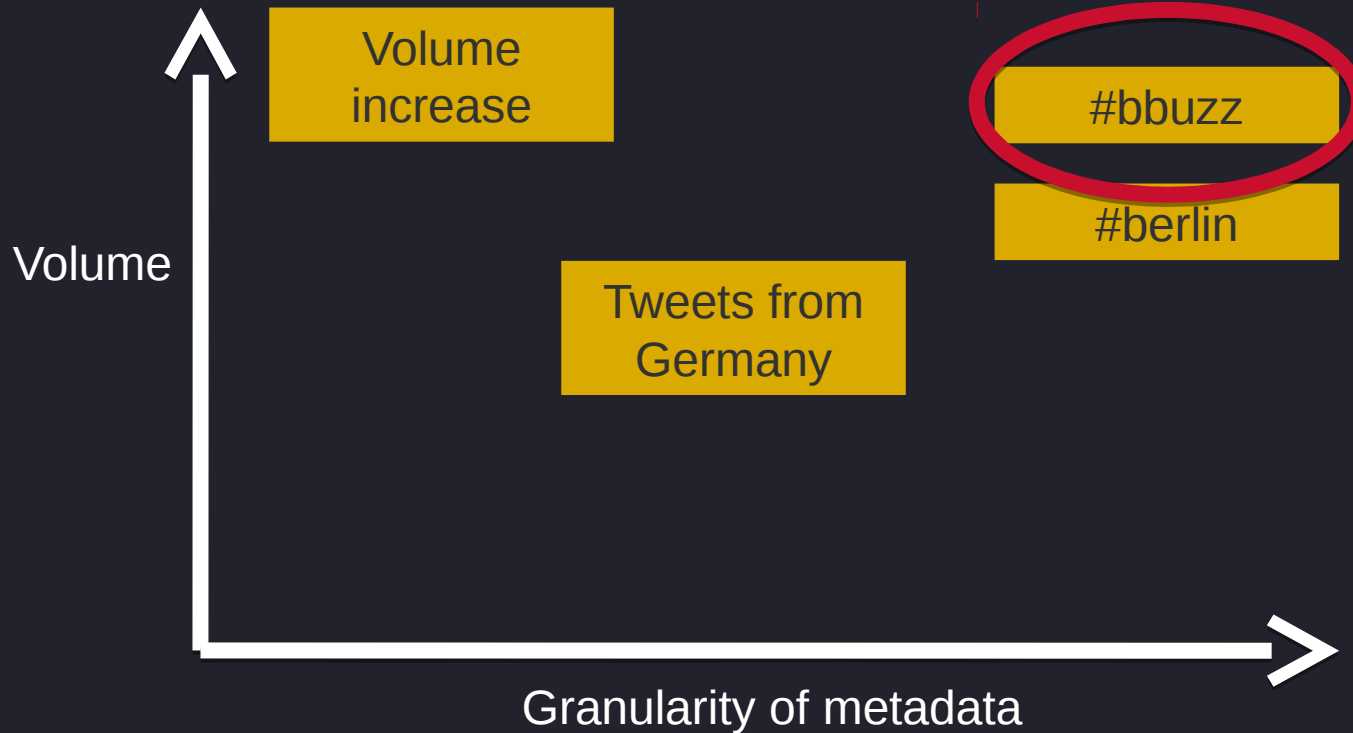
# Meaning/ 2. Grouping events

- Granularity

- Text similarity

- Shape of volume sparkline

# Meaning/ 2. Grouping events

Volume

Volume increase

#bbuzz

#berlin

Tweets from Germany

Granularity of metadata

# Meaning/ An example

Closing remarks

# Say hello/

# @jstanier

# Q&A