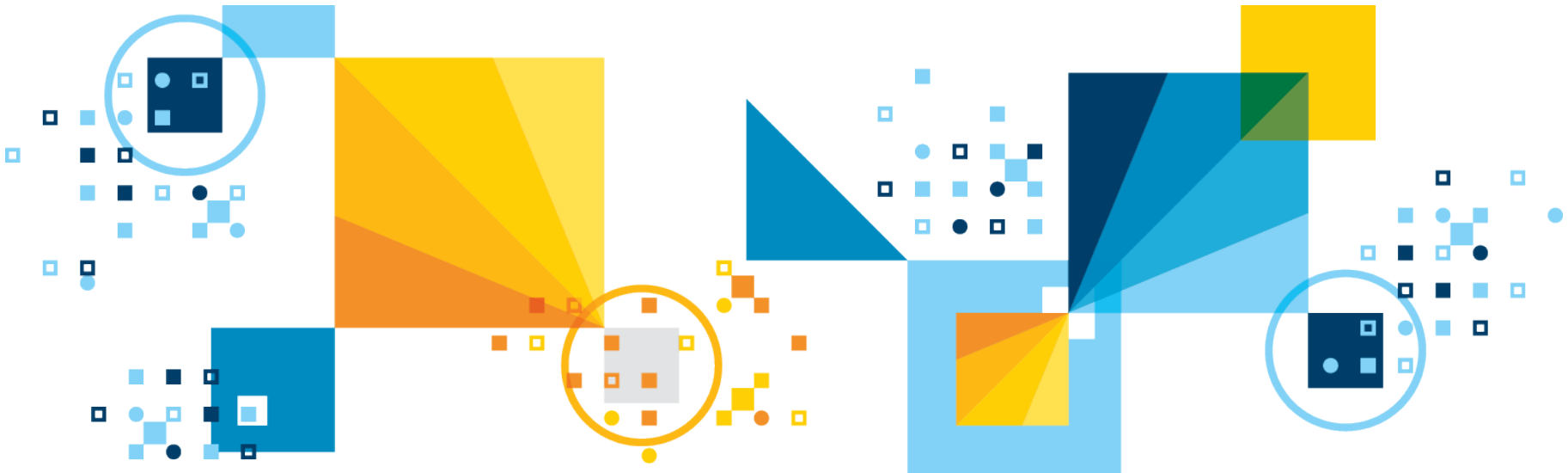


Analyzing and Searching Streams of Social Media Using Spark, Kafka, and Elasticsearch

Markus Lorch, mlorch@de.ibm.com, [@MarkusLorch](https://twitter.com/MarkusLorch)



Outline

- **Introduction and Scope**
 - IBM and Twitter Partnership
 - IBM Insights for Twitter on IBM Bluemix

- **Technology and Experiences**
 - Apache Spark in Streaming Mode as the Processing Engine
 - Apache Kafka as a distributed Messaging Queue
 - Elasticsearch as an “Index-based Repository”
 - Hardware Hosted on IBM Softlayer

Introduction and Scope

Related: TECH, DEAL!

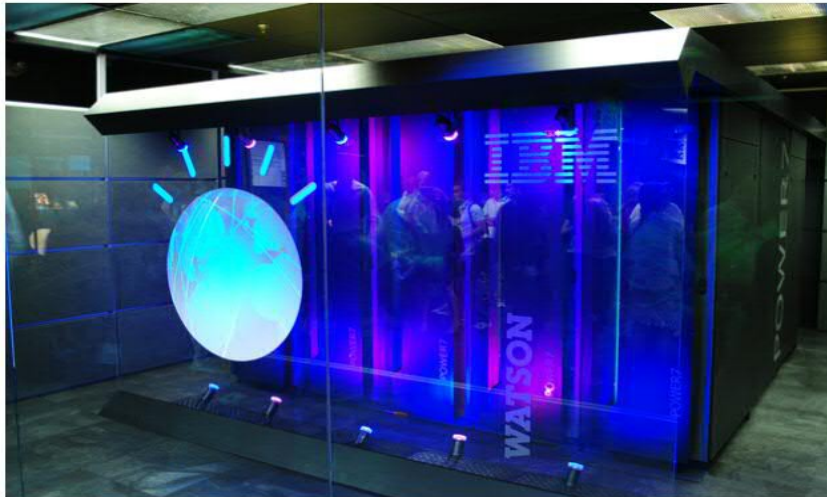
Technology | Wed Oct 29, 2014 4:27pm EDT

IBM, Twitter to partner on business data analytics

WASHINGTON | BY MARINA LOPES

BIG DATA

IBM Introduces Twitter-Fueled Data Services for Business

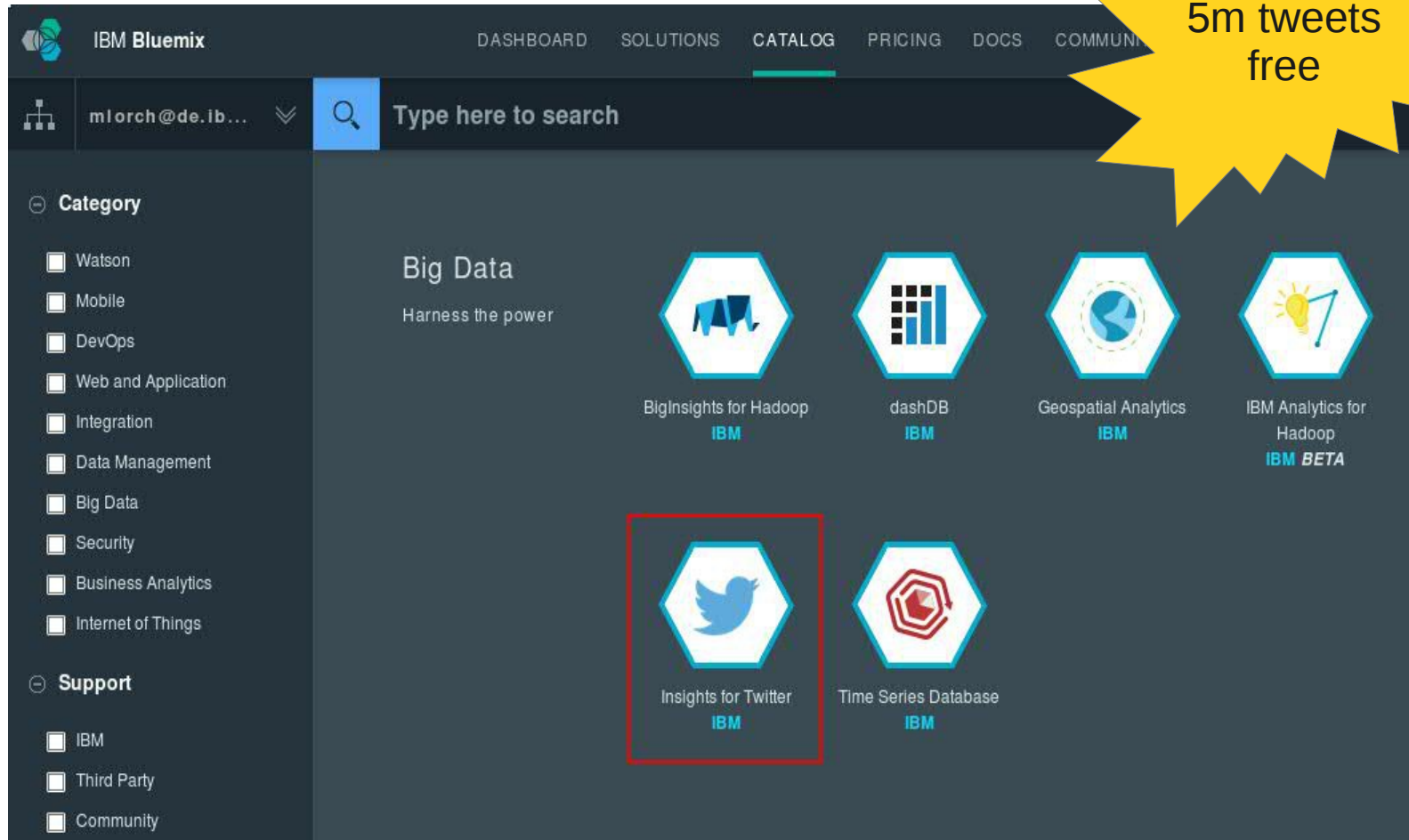


Credit: [Wikimedia/Clockready](#)

IBM Watson mines Twitter for sentiments

IBM's new Insights service harvests data from millions of tweets and uses Watson to analyze them for sentiment and behavior

Insights for Twitter Service on IBM Bluemix



Use it to build your own service leveraging Twitter Tweets and IBM Analytics

Sample Application



Twitter Search

Twitter Count

kids toys

45723



Christopher @teesang

This needs to come to Canada. #parenthood #kids - Flinto's fun learning toys for kids <http://t.co/kbVGismngi> #flintobox

IBM Insights for Twitter



kids toys has:children

Twitter Search

Twitter Count

kids toys has:children

1471



TraderRLH @TraderRLH

RT @Todd_Kincannon: I'm preparing mental toughening exercises for my future kids. Yelling "THIS IS SPARTA!" a lot and throwing their toys away in front of them.

[IBM Insights for Twitter](#)



leilan mc @bgkahuna

The airport gave me a great idea. I'm going to collect all the toys my kids don't put back where they belong and auction them off.

[IBM Insights for Twitter](#)



kids toys has:children is:married

Twitter Search

Twitter Count

kids toys has:children is:married

29



Jonny @Jonny5isalive1

Great to see the kids play with all of their new toys today. Absolutely awesome.

IBM [Insights](#) for [Twitter](#)



Mat Owsley @Mowsley

I only drink crown royal so my kids have bags for their toys.

IBM [Insights](#) for [Twitter](#)



kids toys has:children is:married sentiment:positive

Twitter Search

Twitter Count

kids toys has:children is:married sentiment:positive

11



Jonny @Jonny5isalive1

Great to see the kids play with all of their new toys today. Absolutely awesome.

IBM Insights for Twitter



Samantha Esmeralda @sam_esmeralda

@Culligan27 hahaha...I think we are having more fun with our kids toys too..hes only 11 weeks old!

IBM Insights for Twitter

#BUZZ EXAMPLE



#bbuzz sentiment:positive followers_count:1000

Twitter Search

Twitter Count

#bbuzz sentiment:positive followers_count:1000

13



Berlin Buzzwords @berlinbuzzwords

Get ready for the sixth edition of Berlin Buzzwords and save your limited "Trust Us"-ticket <http://t.co/e0g6LkGLAz> #bbuzz

IBM Insights for Twitter



Simon Willnauer @s1m0nw

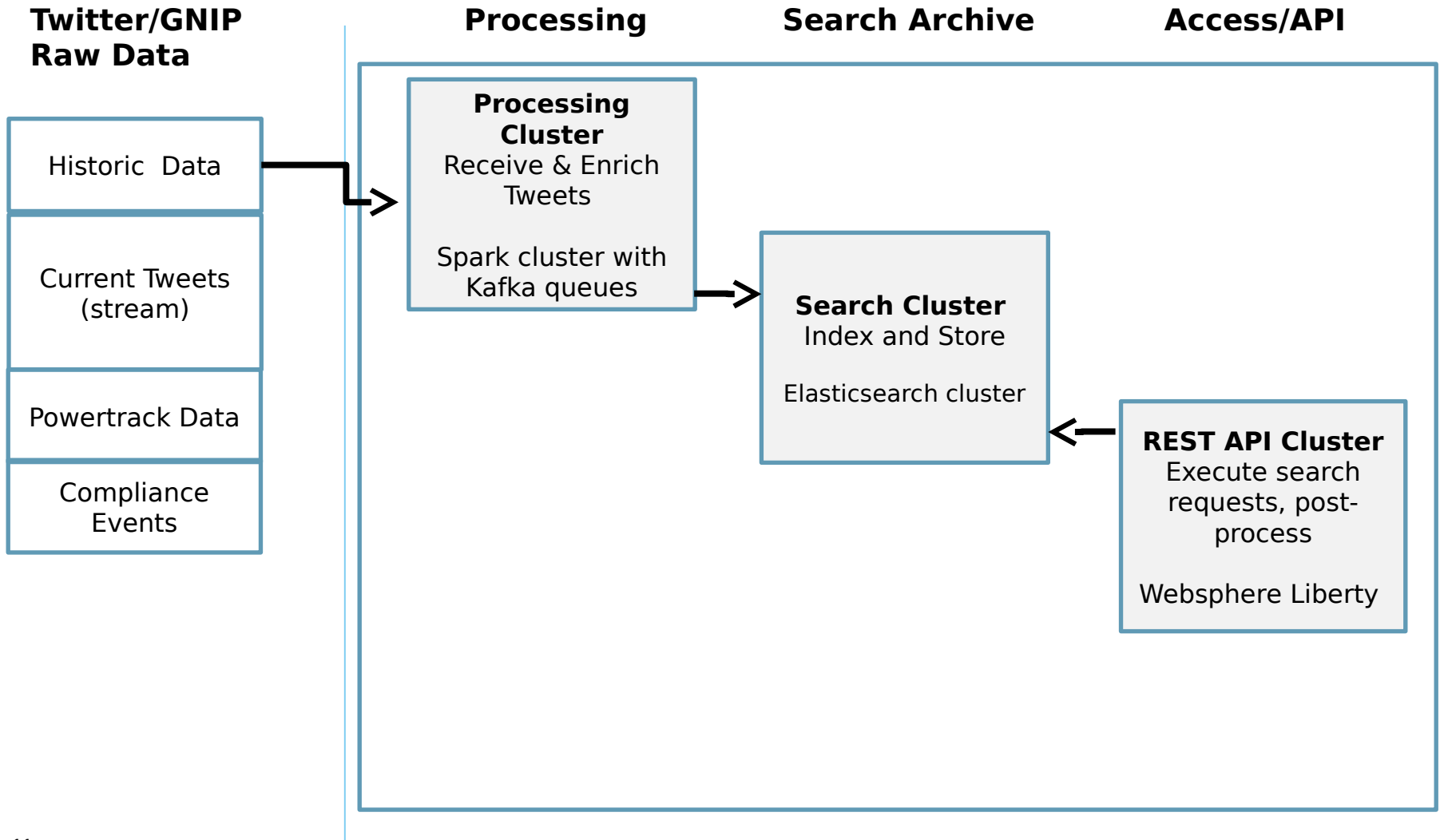
RT @Ellen_Friedman: I'll speak at @berlinbuzzwords 2015 Great #OSS conference: good tech, super people. A real community. <http://t.co/sb156lUpS1> #bbuzz #bigdata

IBM Insights for Twitter

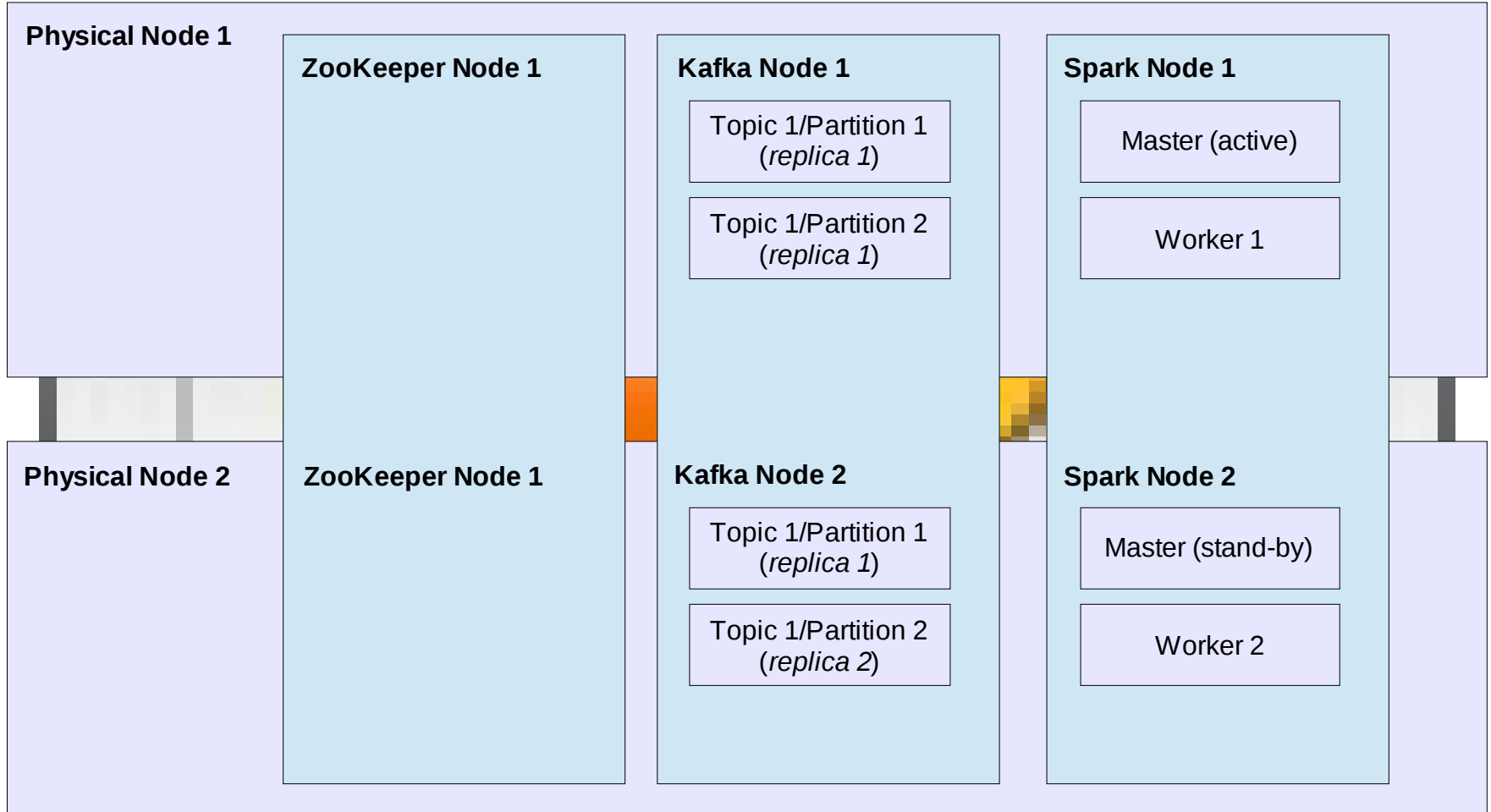
Outline

- Introduction and Scope
 - IBM and Twitter Partnership
 - IBM Insights for Twitter on IBM Bluemix
- **Technology and Experiences**
 - Apache Spark in Streaming Mode as the Processing Engine
 - Apache Kafka as a distributed Messaging Queue
 - Elasticsearch as an “Index-based Repository”
 - Hardware Hosted on IBM Softlayer

High-level Architecture



Use of Kafka and Spark



System Design

Hardware

- System running on IBM Softlayer bare-metal servers, use many (relatively) small servers which no hardware redundancy.
- Smaller servers → faster recovery from failure and higher redundancy
- each function has a minimum of 3 servers to ensure HA even in the case of maintenance
- Continuous availability (rolling restarts)

Software

- All redundancy provided by software stack
- Use Spark as processing engine
- leverage Spark streaming with micro-batching
 - future: direct streams with better Kafka integration
- Use Kafka as distributed messaging / queueing system with message persistence
- Leveraging a large Elasticsearch cluster as an index-based repository optimized for low query-time
- Leveraging IBM Websphere Liberty for REST API implementation

Experiences

- Kafka helps to decouple processing and queue messages
→ ability to delay incoming processing
- Kafka also allows us to read raw-data as well as analyzed data with multiple consumers (e.g. index but also write to files)
- Spark streaming with micro-batching adds about 1 sec delay, creates very small RDDs
- Spark streaming causes inefficient copying of data from Kafka, and issues with locally stored RDDs (spark scratch)
- Existing analytics code (java) was easy to get to run in Spark, much new analytics code is being written for Spark
- Elasticsearch provided a solid and scalable search engine, but with larger cluster size maintenance is not effortless
- Storing the Tweets only in the index avoids joins with DB storage