

# RDF Stream Processing on Big Data

**Oscar Orlando CEBALLOS ARGOTE, M.Sc.**

[oscar.ceballos@correounivalle.edu.co](mailto:oscar.ceballos@correounivalle.edu.co)

**John SANABRIA, Ph.D.**

Universidad del Valle  
Director

**Oscar CORCHO, Ph.D.**

Universidad Politécnica de Madrid  
Codirector

**María-Constanza PABÓN, Ph.D.**

Pontificia Universidad Javeriana  
Asesora

Doctorado en Ingeniería  
Énfasis en Ciencias de la Computación  
Escuela de Ingeniería de Sistemas y Computación  
Universidad del Valle, Cali - Colombia

# Contents

---

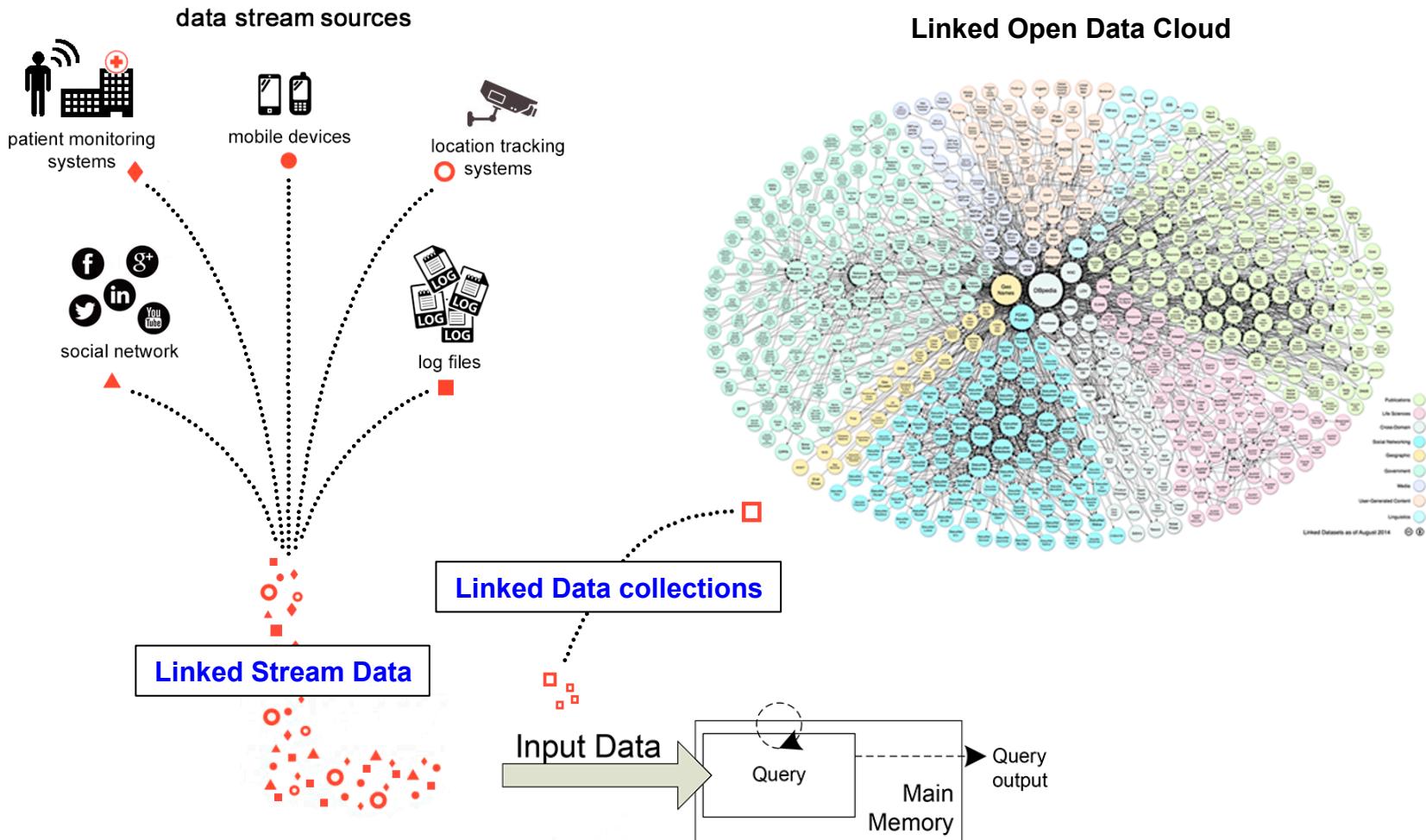
1. RDF Stream Processing
2. Big Data Technologies
3. Current work
  - 3.1 Queries on batch (static)
  - 3.2 Queries on streaming
4. To do

# Contents

---

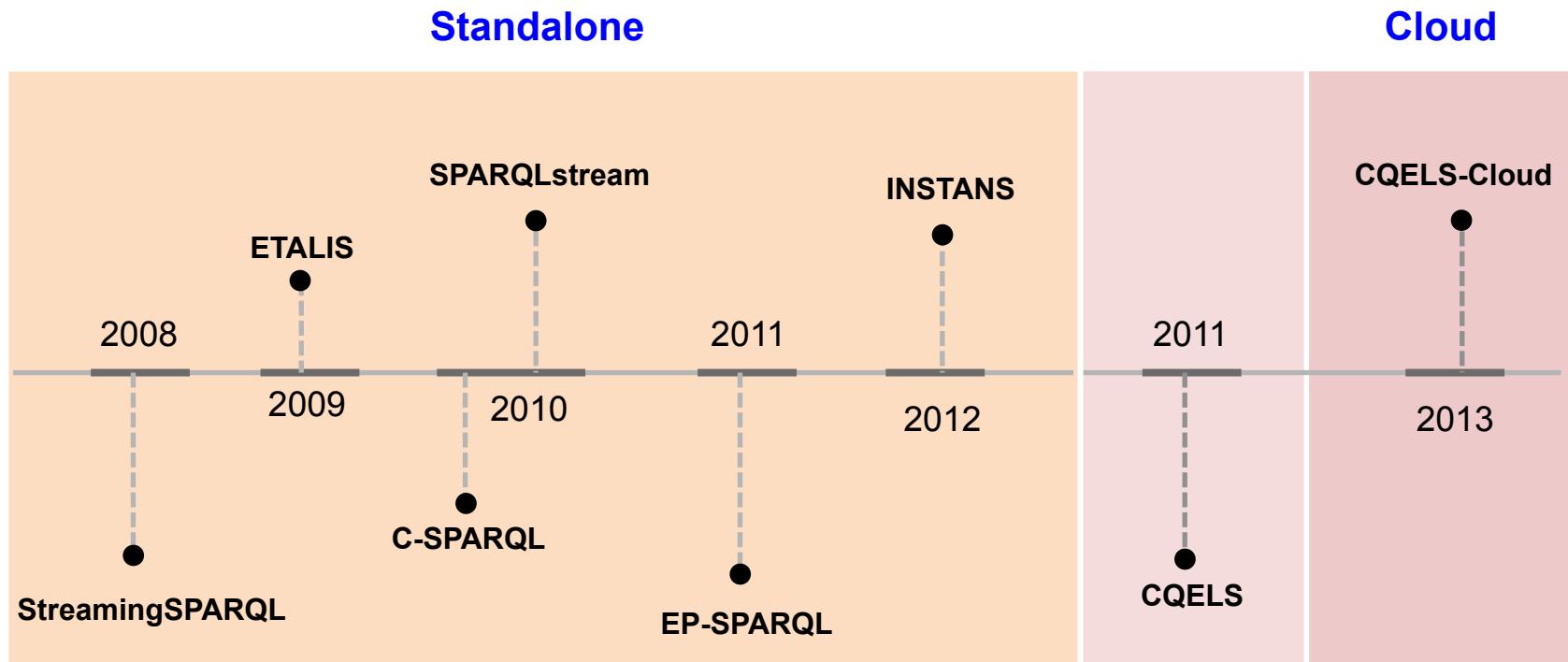
- 1. RDF Stream Processing**
2. Big Data Technologies
3. Current work
  - 3.1 Queries on batch (static)
  - 3.2 Queries on streaming
4. To do

# 1. RDF Stream Processing



# 1. RDF Stream Processing

Some **RDF Stream Processing Systems** (**RSP Systems**) support unified query processing over **Linked Stream Data** (**RDF Stream**) and **Linked Data collections** (**static RDF datasets**)



# Problem Statement

---

RSP Systems **performance** and **scalability** still need to be improved in terms of  
**query execution time, data size (dynamic and static), number**  
**of concurrent queries, and number of streams**

# Contents

---

1. RDF Stream Processing

## 2. Big Data Technologies

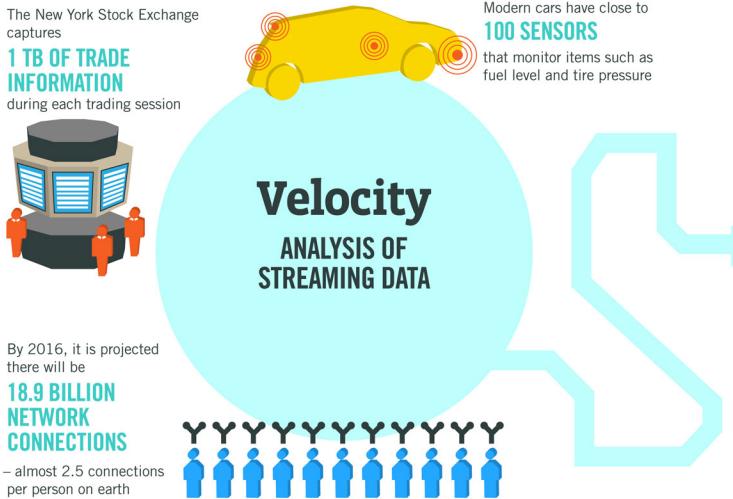
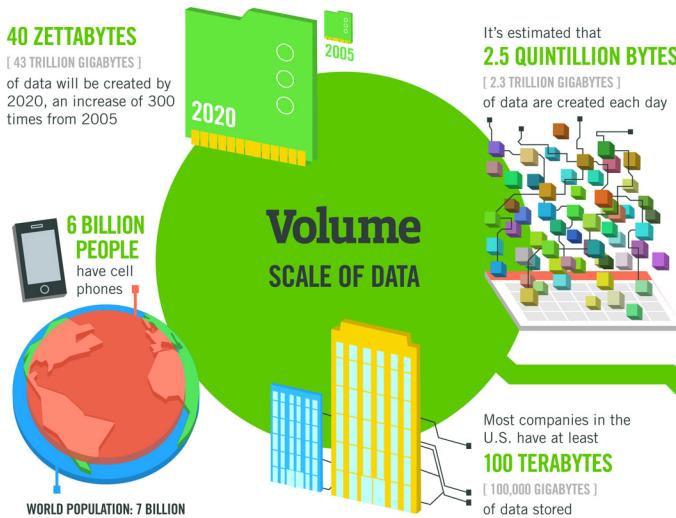
3. Current work

3.1 Queries on batch (static)

3.2 Queries on streaming

4. To do

## 2. Big Data Technologies



### The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume**, **Velocity**, **Variety** and **Veracity**.

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015 **4.4 MILLION IT JOBS** will be created globally to support big data, with 1.9 million in the United States



As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES** [ 161 BILLION GIGABYTES ]



**30 BILLION PIECES OF CONTENT**

are shared on Facebook every month



By 2014, it's anticipated there will be

**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**

are watched on YouTube each month



**Variety DIFFERENT FORMS OF DATA**



**400 MILLION TWEETS** are sent per day by about 200 million monthly active users

**1 IN 3 BUSINESS LEADERS**

don't trust the information they use to make decisions



Poor data quality costs the US economy around

**\$3.1 TRILLION A YEAR**

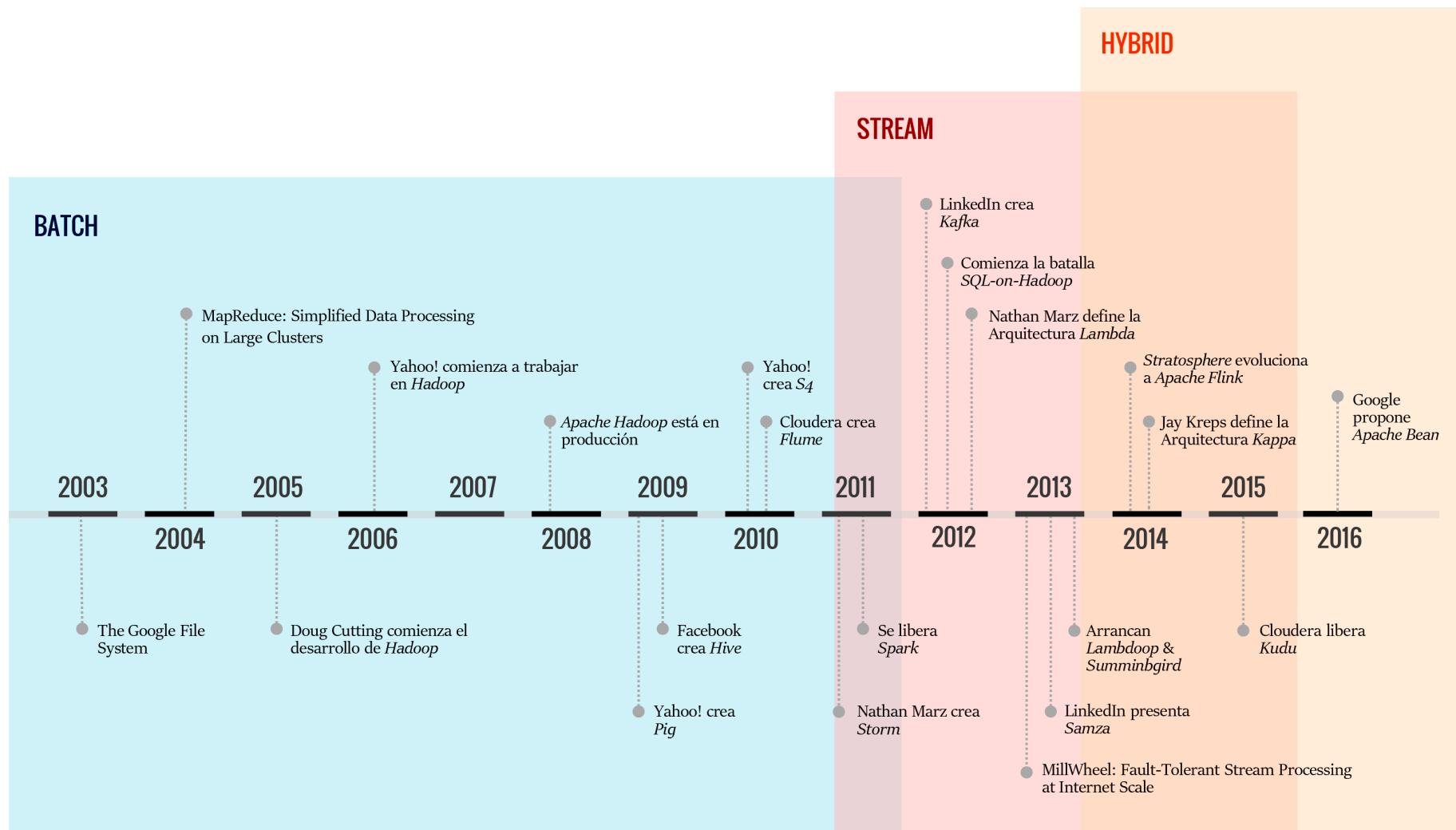


**Veracity UNCERTAINTY OF DATA**

in one survey were unsure of how much of their data was inaccurate



# 2. Big Data Technologies



# Hypothesis

---

When extending **Big Data Technologies** to process real time **unified queries** over **RDF streams** and **static RDF datasets**, it is possible to **improve performance and scalability** limits of current **RSP Systems**

# Research Questions

---

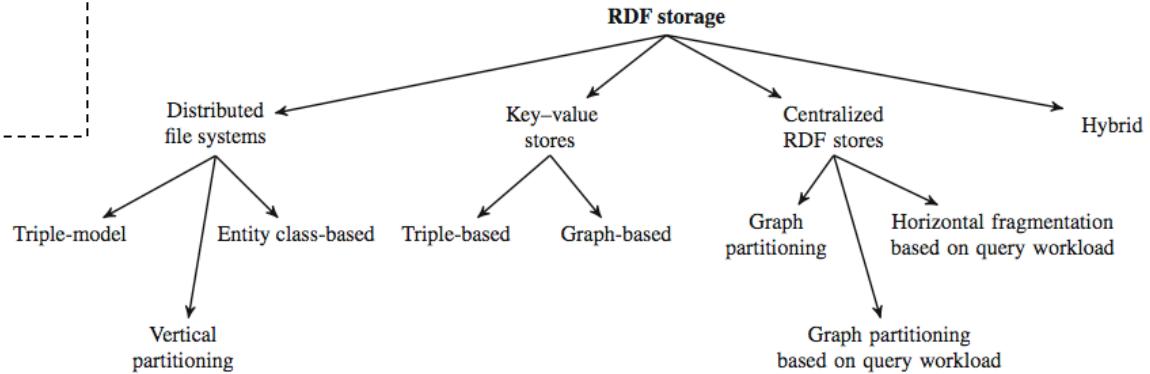
1. Which of the existing technologies in Big Data are appropriate?
  
2. Based on such technology:
  - Is it possible to process continuous queries on RDF streams?
  - Is the set of operations available (e.g., map, reduce, filter, join) sufficient?
  - What other operators are required and necessary to be implemented?
  - What query optimization techniques are necessary to be adapted?
  
3. What tests can be used to evaluate performance, scalability, completeness and correctness properties on the proposed extensions?

# State of Art

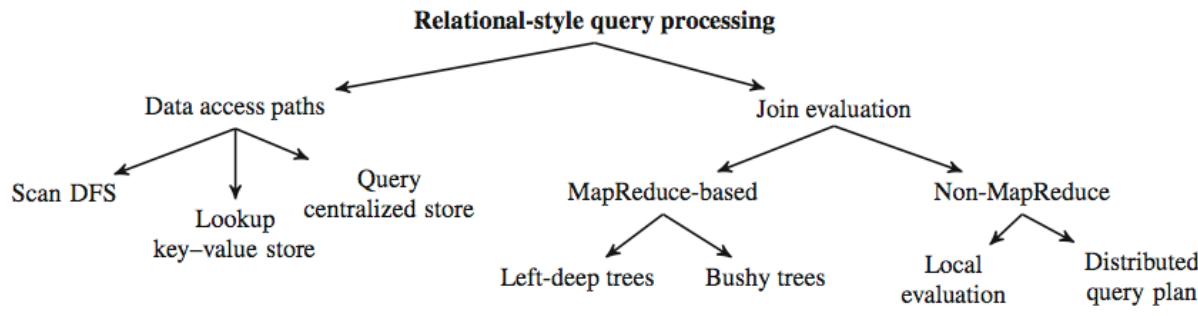
The VLDB Journal (2015) 24:67–91  
DOI 10.1007/s00778-014-0364-z

## RDF in the clouds: a survey

Zoi Kaoudi · Ioana Manolescu



Taxonomy of storage back-ends and partitioning schemes used by the systems



Taxonomy of relational-style query processing strategies

# State of Art

*The evolution of massive-scale data processing*

(<https://goo.gl/5k0xal> - 2016)

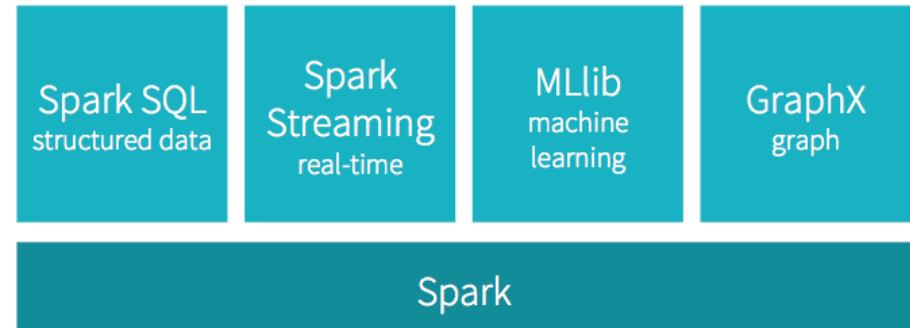
	Open Source	Managed Service	Auto-Awesome	Batch	Streaming	Iterative	Pipelines	High-level API	Optimizer	Unified Engine	Unified API	No Lambda	Exactly Once	Unified Engine	State	Timers	Watermarks	Windowing	Triggers
MapReduce				X								X							
Hadoop	X	X		X								X	X						
Flume			X	X								X	X	X					
Storm	X				X														X
Spark	X	X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	
MillWheel					X			X								X	X	X	X
Flink	X			X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Cloud Dataflow	X	X	X	X	X			X	X	X	X	X	X	X	X	X	X	X	X

# Apache Spark

***Resilient Distributed Datasets: A fault-tolerant abstraction for in-memory cluster computing***  
*(Zaharia et al., 2012)*

## Resilient Distributed Datasets (RDDs)

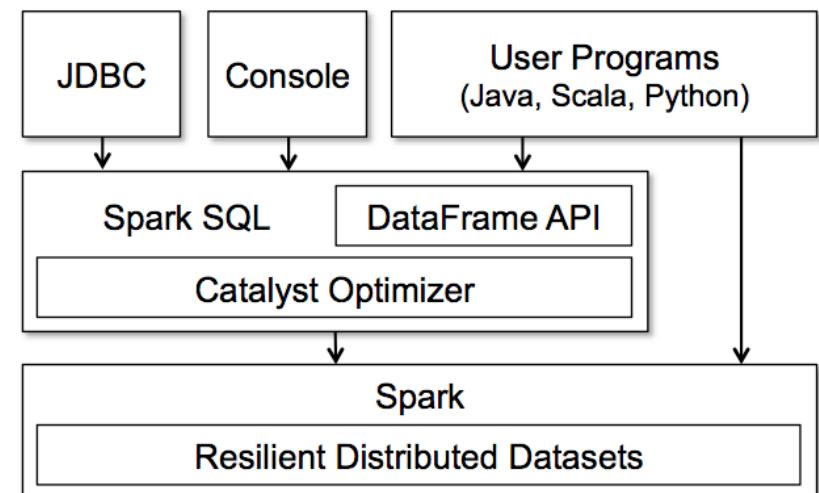
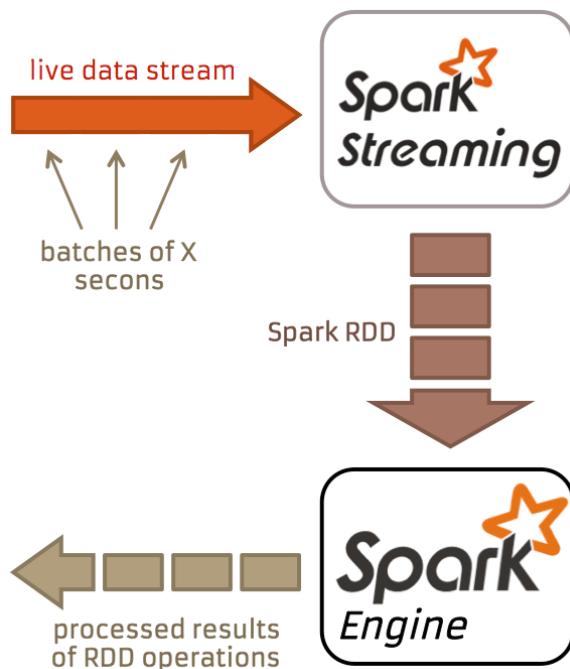
- Write programs in terms of operations on distributed datasets
- Partitioned collections of objects spread across a cluster; stored in memory or on disk
- RDDs built and manipulated through a diverse set of parallel transformations (map, filter, join) and actions (count, collect, save)
- RDDs automatically rebuilt on machine failure



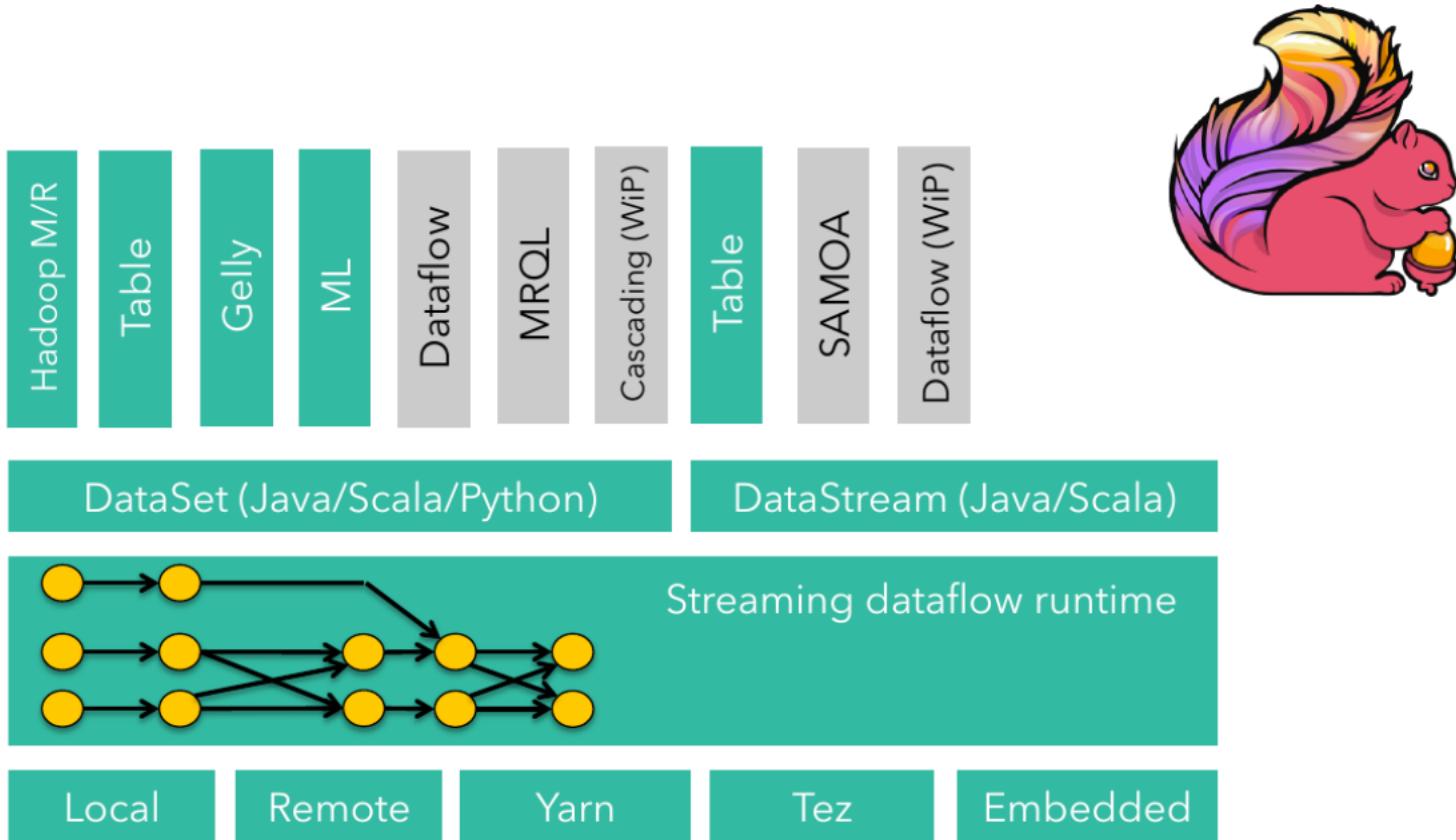
# Apache Spark

***Discretized Streams: Fault-tolerant streaming computation at scale***  
(Zaharia et al., 2013)

***Spark SQL: Relational Data Processing in Spark***  
(Armbrust et al., 2015)

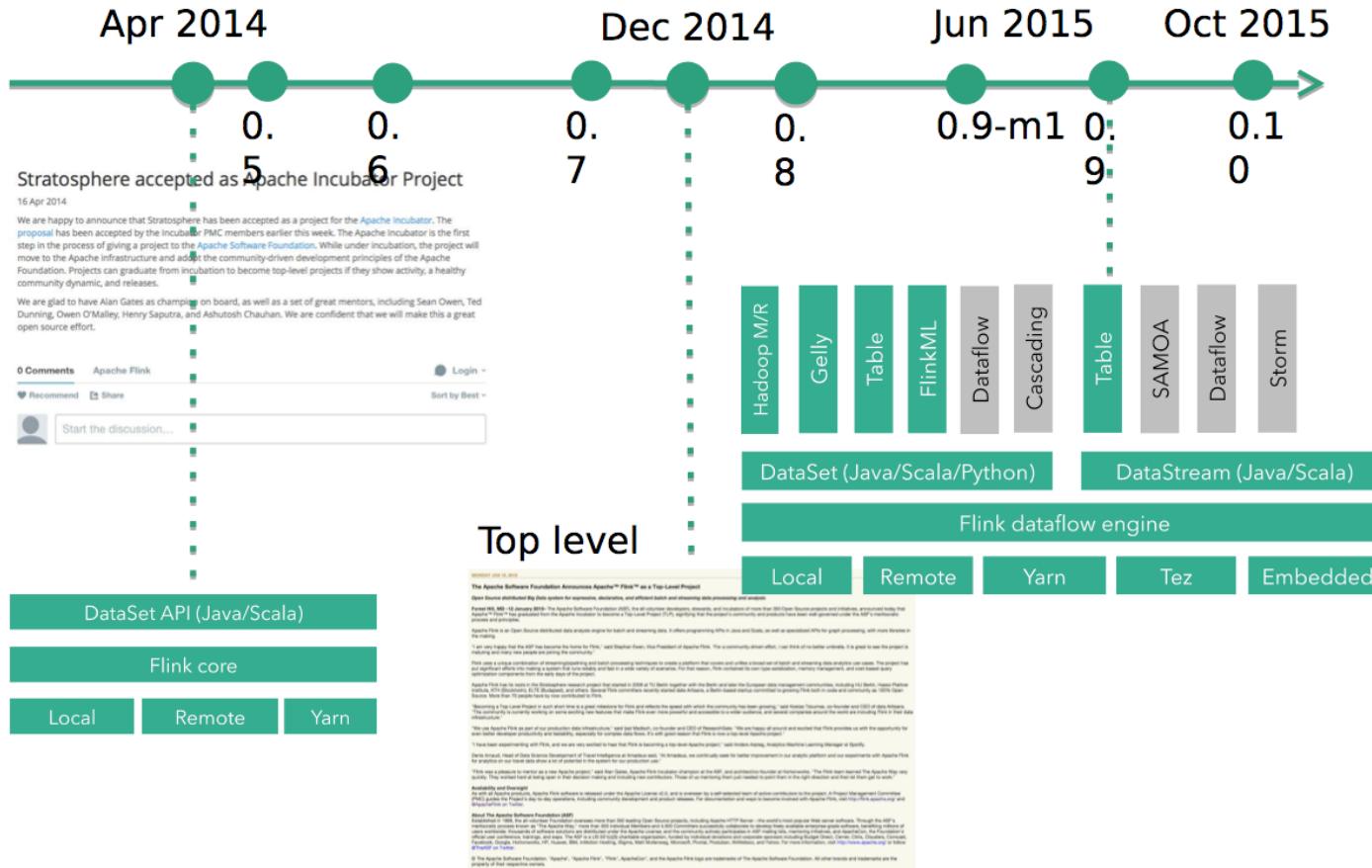


# Apache Flink



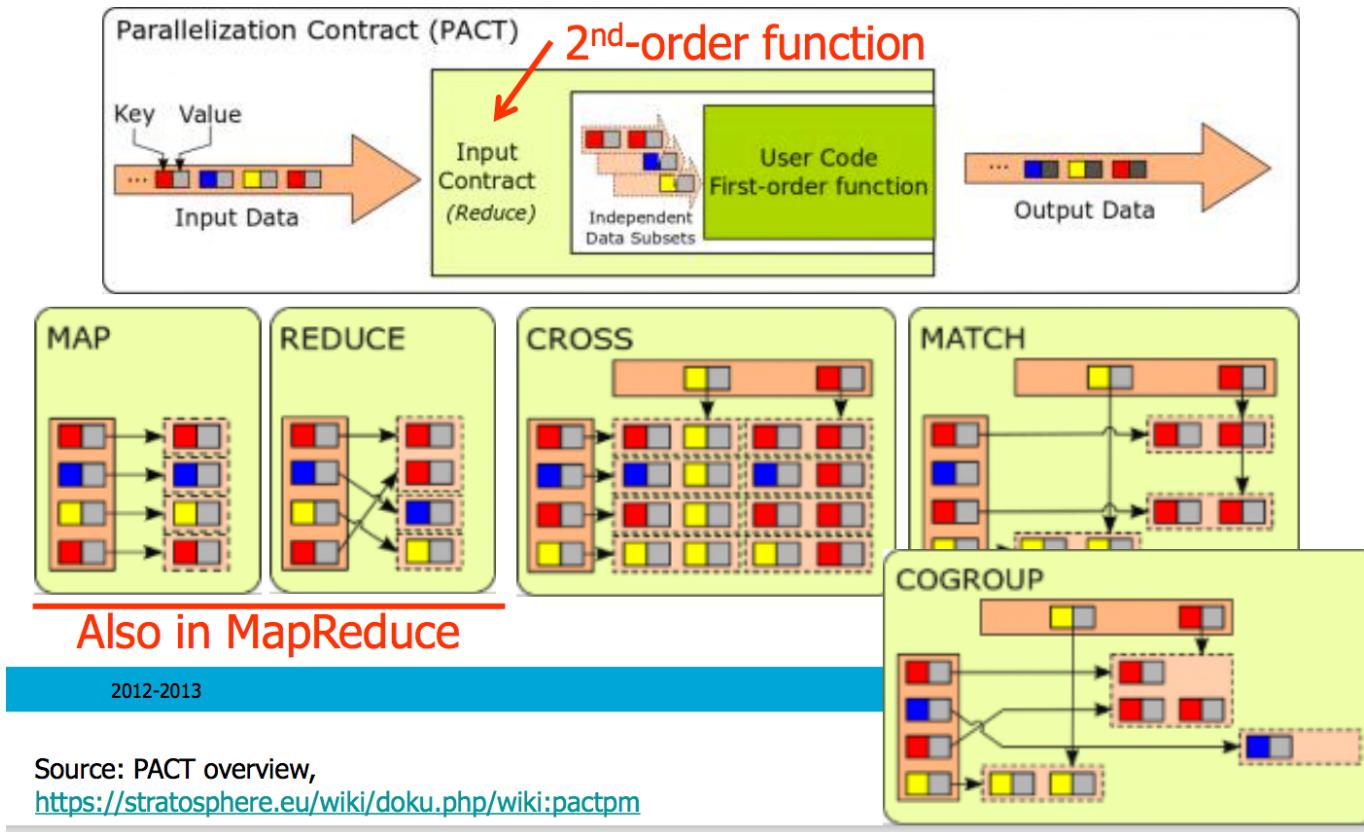
# Apache Flink

***The Stratosphere platform for big data analytics (Alexandrov et al., 2014)***



# Apache Flink

## Stratosphere Programming Contracts (PACTs) [1/2]



# Apache Flink

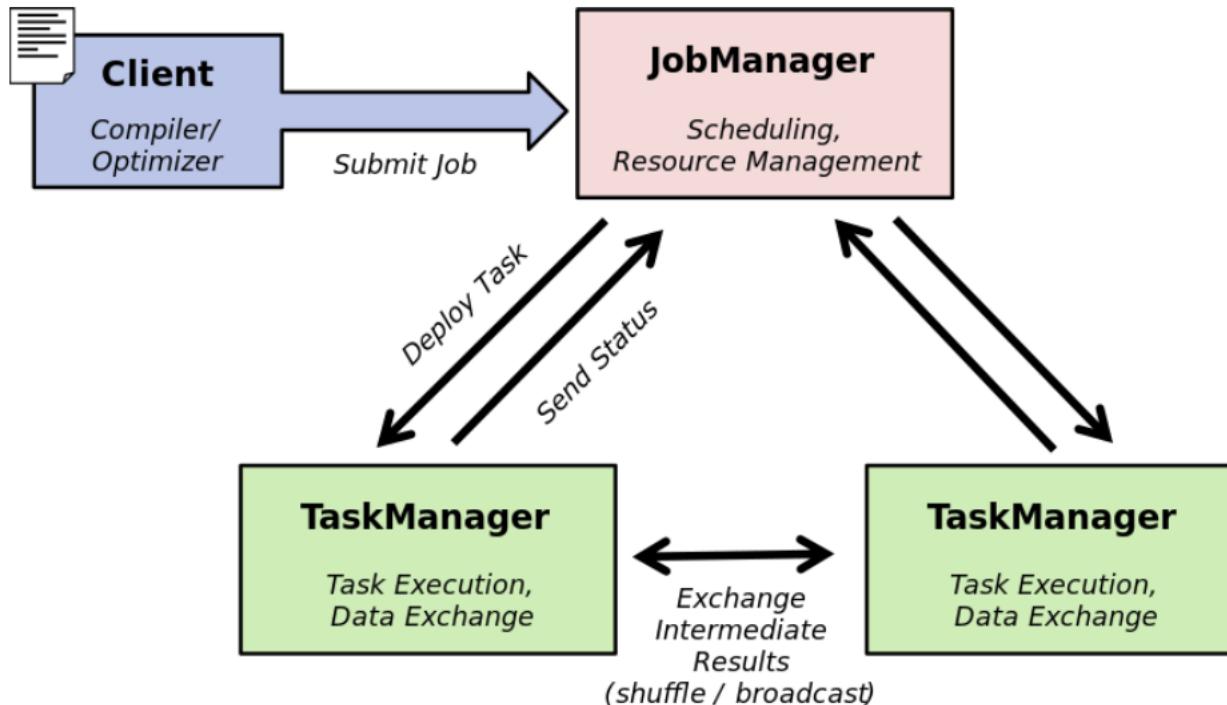
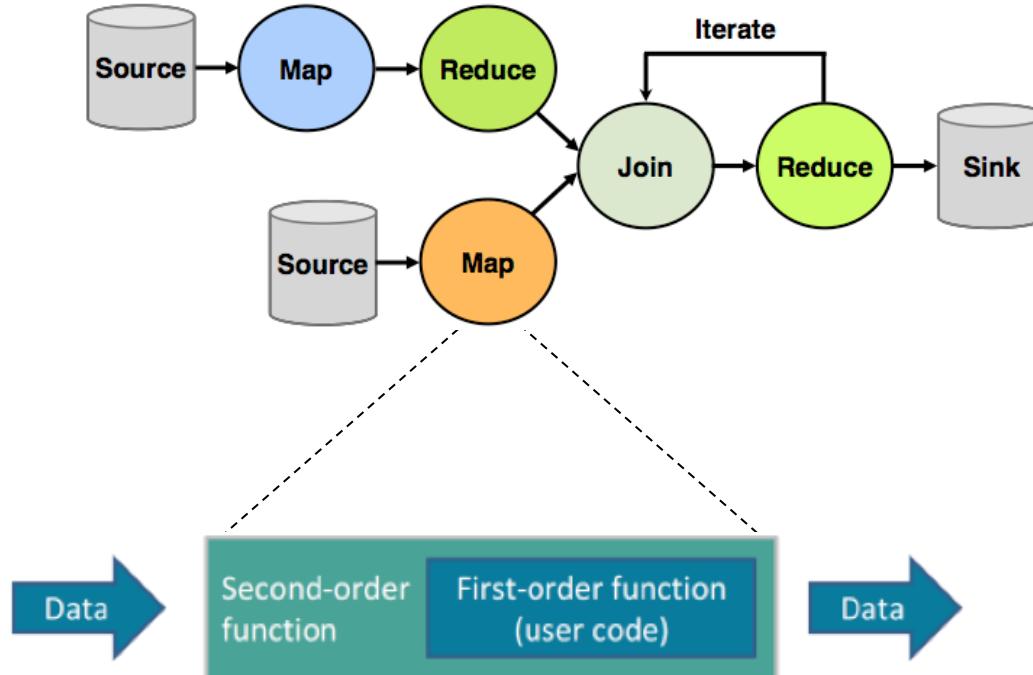


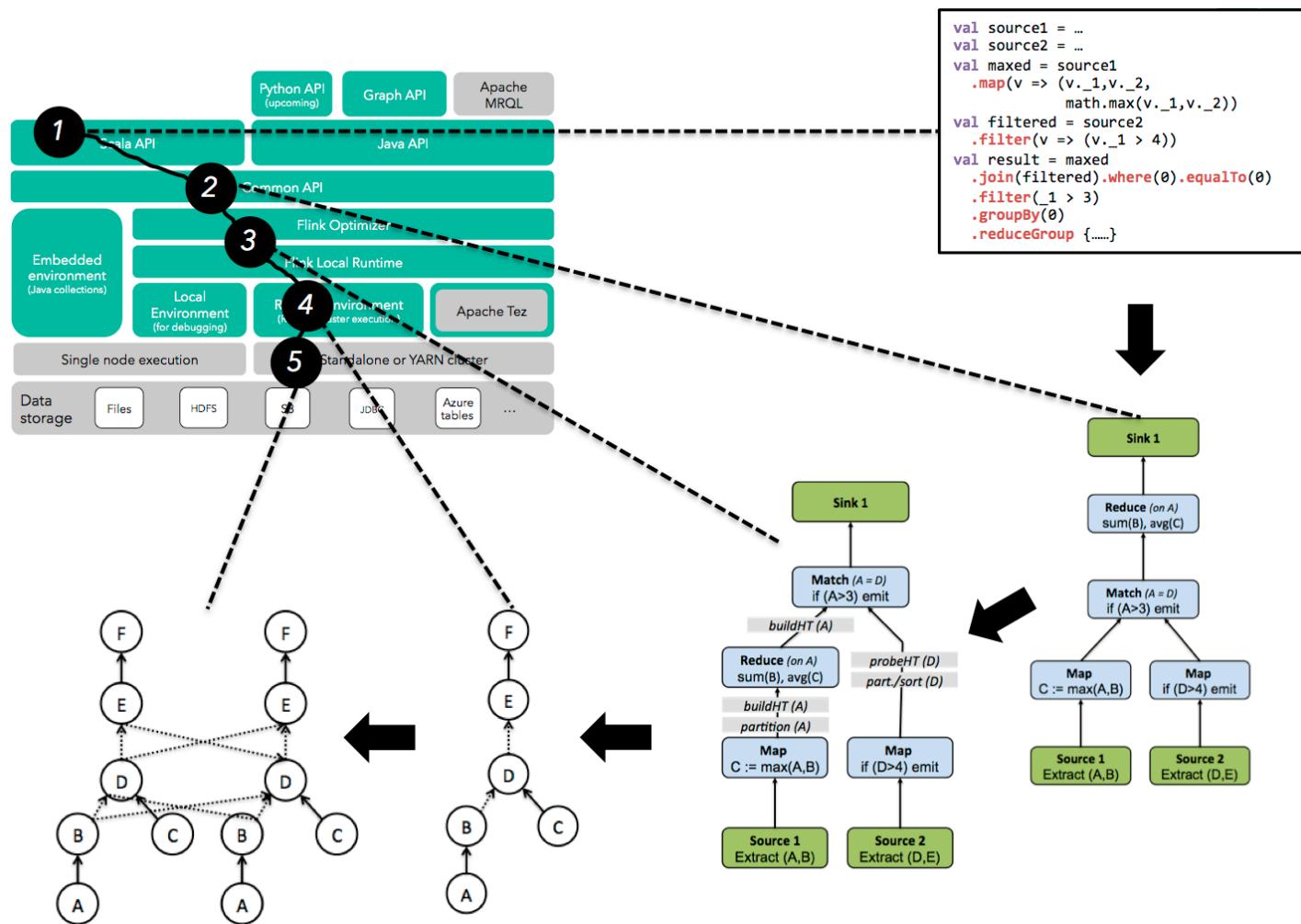
Figure: The JobManager is the coordinator of the Flink system  
TaskManagers are the workers that execute parts of the parallel programs.

# Apache Flink

A **program** is expressed as an arbitrary **data flow** consisting of **transformations**, **sources** and **sinks**.



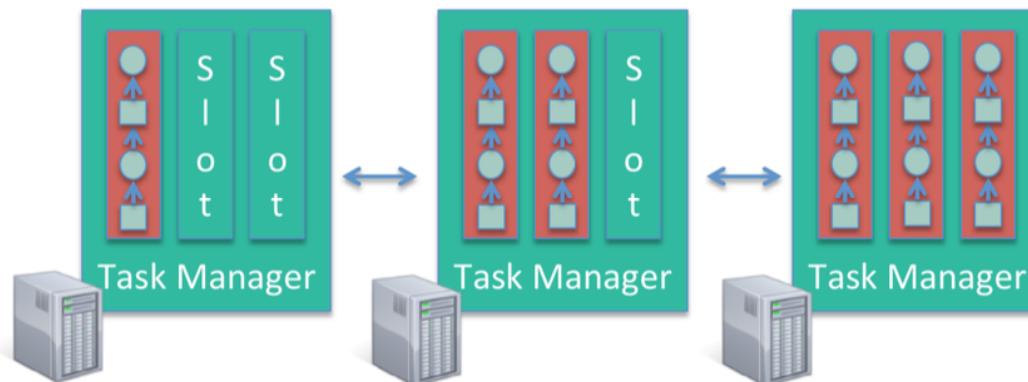
# Apache Flink



# Apache Flink

## Task Manager

- ① Operations are split up into **tasks** depending on the specified parallelism
- ② Each parallel instance of an operation runs in a separate **task slot**
- ③ The scheduler may run several tasks from different operators in one task slot



# Contents

---

1. RDF Stream Processing
2. Big Data Technologies
- 3. Current work**
  - 3.1 Queries on batch (static)
  - 3.2 Queries on streaming
4. To do

# Contents

---

1. RDF Stream Processing
2. Big Data Technologies
3. Current work
  - 3.1 Queries on batch (static)**
  - 3.2 Queries on streaming
4. To do

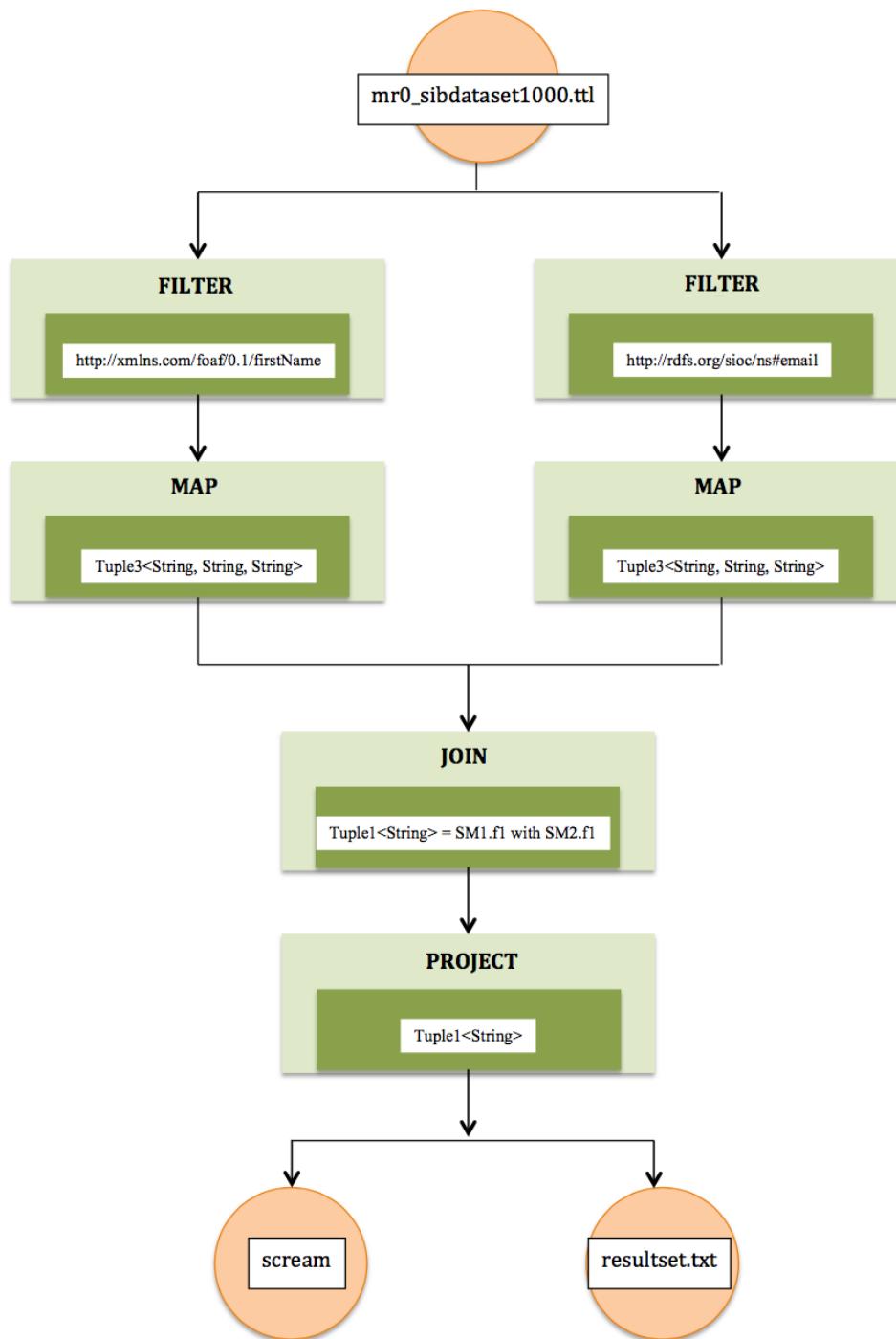
### 3. Current Work - batch

---

```
PREFIX sibp: <http://www.ins.cwi.nl/sib/person/>
PREFIX foaf: <http://xmlns.com/foaf/0.1>
PREFIX sioc: <http://rdfs.org/sioc/ns#>

SELECT ?name
FROM <mr0_sibdataset1000.ttl>
WHERE
{ ?person foaf:firstName ?name ;
            sioc:email      ?email
}
```

```
(project (?name)
  (bgp
    (triple ?person <http://xmlns.com/foaf/0.1firstName> ?name)
    (triple ?person <http://rdfs.org/sioc/ns#email> ?email)
  ))
```



# Contents

---

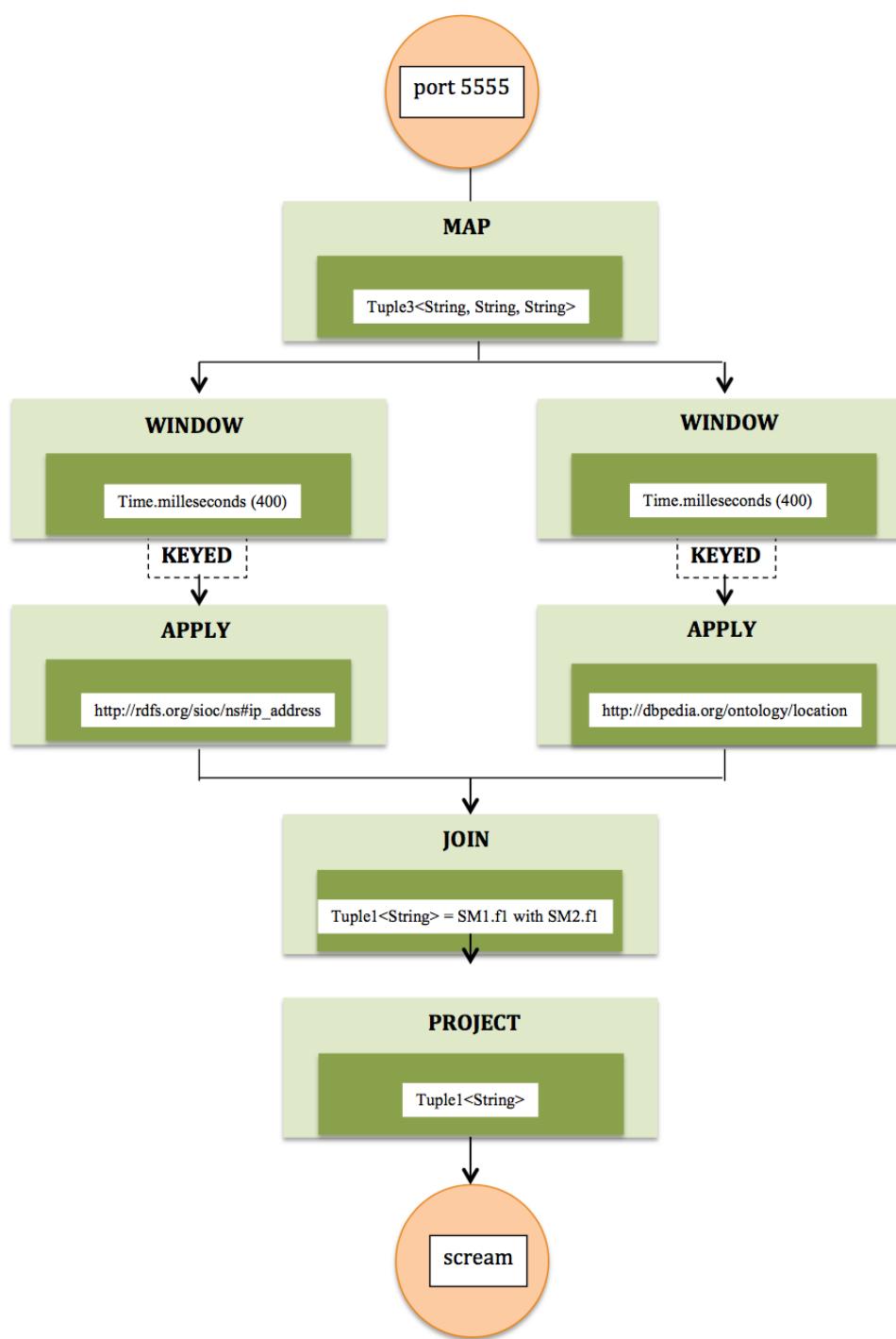
1. RDF Stream Processing
2. Big Data Technologies
3. Current work
  - 3.1 Queries on batch (static)
  - 3.2 Queries on streaming**
4. To do

### 3. Current Work - streaming

```
String cqels_query = "PREFIX sibp: <http://www.ins.cwi.nl/sib/person/>" +  
    "PREFIX foaf: <http://xmlns.com/foaf/0.1>" +  
    "PREFIX sioc: <http://rdfs.org/sioc/ns#>" +  
    "SELECT ?name " +  
    "WHERE {" +  
        " STREAM <http://univalle.org/streams/rfid> [RANGE 400ms] " +  
        " { ?person foaf:firstName ?name ." +  
        " ?person sioc:email ?email }" +  
    "};
```

```
PREFIX sibp: <http://www.ins.cwi.nl/sib/person/>  
PREFIX foaf: <http://xmlns.com/foaf/0.1>  
PREFIX sioc: <http://rdfs.org/sioc/ns#>  
  
SELECT ?name  
WHERE  
{ GRAPH <http://univalle.org/streams/rfid>  
{ ?person foaf:firstName ?name .  
?person sioc:email ?email  
}  
}
```

```
(project (?name)  
  (graph <http://univalle.org/streams/rfid>  
    (bgp  
      (triple ?person <http://xmlns.com/foaf/0.1firstName> ?name)  
      (triple ?person <http://rdfs.org/sioc/ns#email> ?email)  
    )))
```



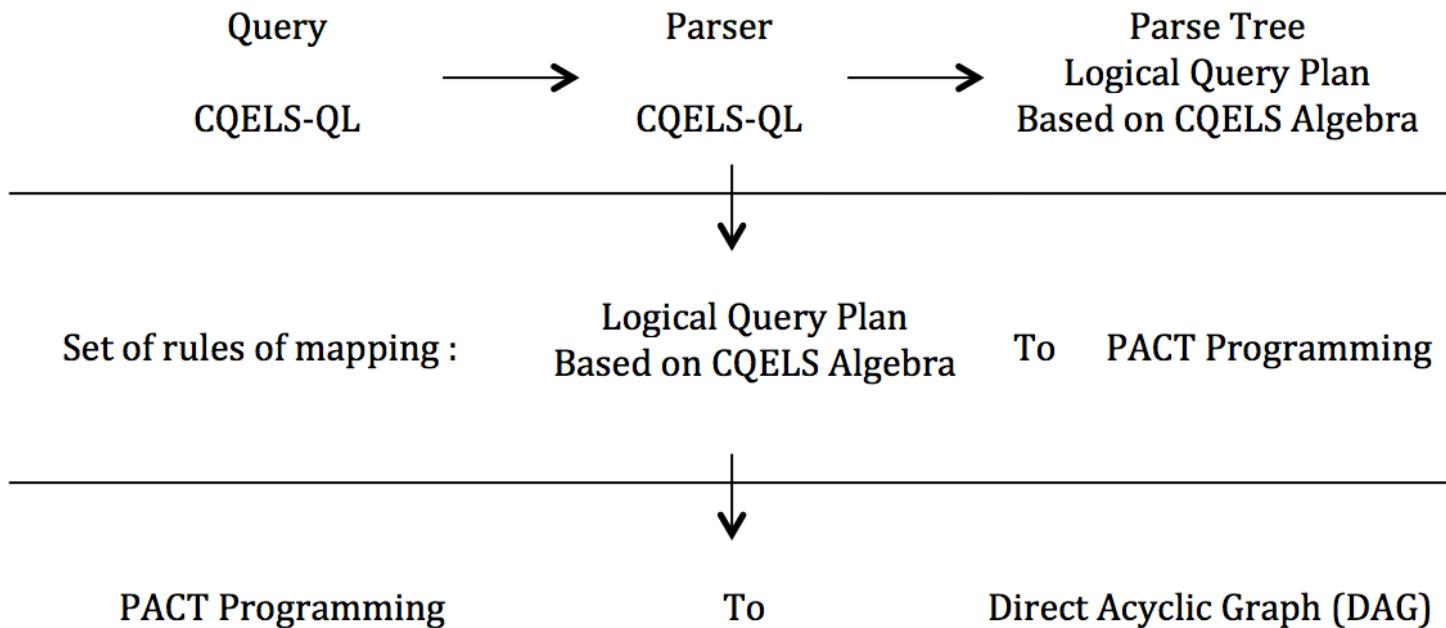
# Contents

---

1. RDF Stream Processing
2. Big Data Technologies
3. Current work
  - 3.1 Queries on batch (static)
  - 3.2 Queries on streaming
4. To do

## 4. To do

---



## 4. To do

---

```
String cqels_query="PREFIX lv: <http://deri.org/floorplan/>" +
  "PREFIX dc: <http://purl.org/dc/elements/1.1/>" +
  "PREFIX foaf: <http://xmlns.com/foaf/0.1/>" +
  "SELECT ?loc2 ?locName (count(distinct ?coAuth) as ?noCoAuths)" +
  "FROM NAMED <http://deri.org/floorplan/>" +
  "WHERE {" +
    "  GRAPH <http://deri.org/floorplan/>" +
    "    {?loc2 lv:name ?locName." +
    "    ?loc2 lv:connected ?loc1}" +
    "    STREAM <http://deri.org/streams/rfid> [TRIPLES 1]" +
    "      {?auth lv:detectedAt ?loc1}" +
    "    STREAM <http://deri.org/streams/rfid> [RANGE 30s]" +
    "      {?coAuth lv:detectedAt ?loc2}" +
    "      {?paper dc:creator ?auth.}" +
    "      ?paper dc:creator ?coAuth." +
    "      ?auth foaf:name \"AUTHORNAME\"^^<http://www.w3.org/2001/XMLSchema#string> }" +
    "    FILTER(?auth!=?coAuth)" +
  "}" +
  "GROUP BY ?loc2 ?locName";
```

## 4. To do

---

```
PREFIX lv: <http://deri.org/floorplan/>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX dc: <http://purl.org/dc/elements/1.1/>

SELECT ?loc2 ?locName (count(distinct ?coAuth) AS ?noCoAuths)
FROM NAMED lv:
WHERE
{ GRAPH lv:
  { ?loc2 lv:name ?locName .
    ?loc2 lv:connected ?loc1
  }
  GRAPH <http://deri.org/streams/rfid>
    { ?auth lv:detectedAt ?loc1 }
  GRAPH <http://deri.org/streams/rfid>
    { ?coAuth lv:detectedAt ?loc2 }
  { ?paper dc:creator ?auth .
    ?paper dc:creator ?coAuth .
    ?auth foaf:name "AUTHORNAME"^^<http://www.w3.org/2001/XMLSchema#string>
  }
  FILTER ( ?auth != ?coAuth )
}
GROUP BY ?loc2 ?locName
```

## 4. To do

---

```
(project (?loc2 ?locName ?noCoAuths)
  (extend ((?noCoAuths ?.0))
    (group (?loc2 ?locName) ((?.0 (count distinct ?coAuth)))
      (filter (!= ?auth ?coAuth)
        (join
          (join
            (join
              (graph <http://deri.org/floorplan>
                (bgp
                  (triple ?loc2 <http://deri.org/floorplan/name> ?locName)
                  (triple ?loc2 <http://deri.org/floorplan/connected> ?loc1)
                ))
              (graph <http://deri.orgstreams/rfid>
                (bgp (triple ?auth <http://deri.org/floorplan/detectedAt> ?loc1))))
            (graph <http://deri.orgstreams/rfid>
              (bgp (triple ?coAuth <http://deri.org/floorplan/detectedAt> ?loc2))))
        (bgp
          (triple ?paper <http://purl.org/dc/elements/1.1/creator> ?auth)
          (triple ?paper <http://purl.org/dc/elements/1.1/creator> ?coAuth)
          (triple ?auth <http://xmlns.com/foaf/0.1/name> "AUTHORNAME"^^<http://www.w3.org/2001/XMLSchema#string>)))
        )))))
```

# References

---

- [1] Vincenzo Gulisano, Ricardo Jimenez-Peris, Marta Patino-Martinez, Claudio Soriente, and Patrick Valduriez. Streamcloud: An elastic and scalable data streaming system. *IEEE Trans. Parallel Distrib. Syst.*, 23(12):2351–2365, December 2012.
  - [2] Jean-Paul Calbimonte, Oscar Corcho, and Alasdair J.G. Gray. Enabling ontology- based access to streaming data sources. In ISWC 2010, volume 6496, pages 96–111. Springer Berlin Heidelberg, 2010.
  - [3] Juan Sequeda and Oscar Corcho. Linked stream data: A position paper. In SSN09, pages 148–157, 2009.
  - [4] Danh Le-Phuoc, Minh Dao-Tran, Josiane Xavier Parreira, and Manfred Hauswirth. A native and adaptive approach for unified processing of linked streams and linked data. In Proceedings of the 10th International Conference on The Semantic Web - Volume Part I, ISWC’11, pages 370–388, Berlin, Heidelberg, 2011
  - [5] Danh Le-Phuoc, Hoan Nguyen Mau Quoc, Chan Le Van, and Manfred Hauswirth. Elastic and scalable processing of linked stream data in the cloud. In The Semantic Web – ISWC 2013, volume 8218 of Lecture Notes in Computer Science, pages 280–297. Springer Berlin Heidelberg, 2013.
  - [6] IBM, Paul Zikopoulos, and Chris Eaton. Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data. McGraw-Hill Osborne Media, 1st edition, 2011.
-