

# Database Systems

## 12 Distributed Analytics

**Matthias Boehm**

Graz University of Technology, Austria  
Computer Science and Biomedical Engineering  
Institute of Interactive Systems and Data Science  
BMVIT endowed chair for Data Management

# Hadoop History and Architecture

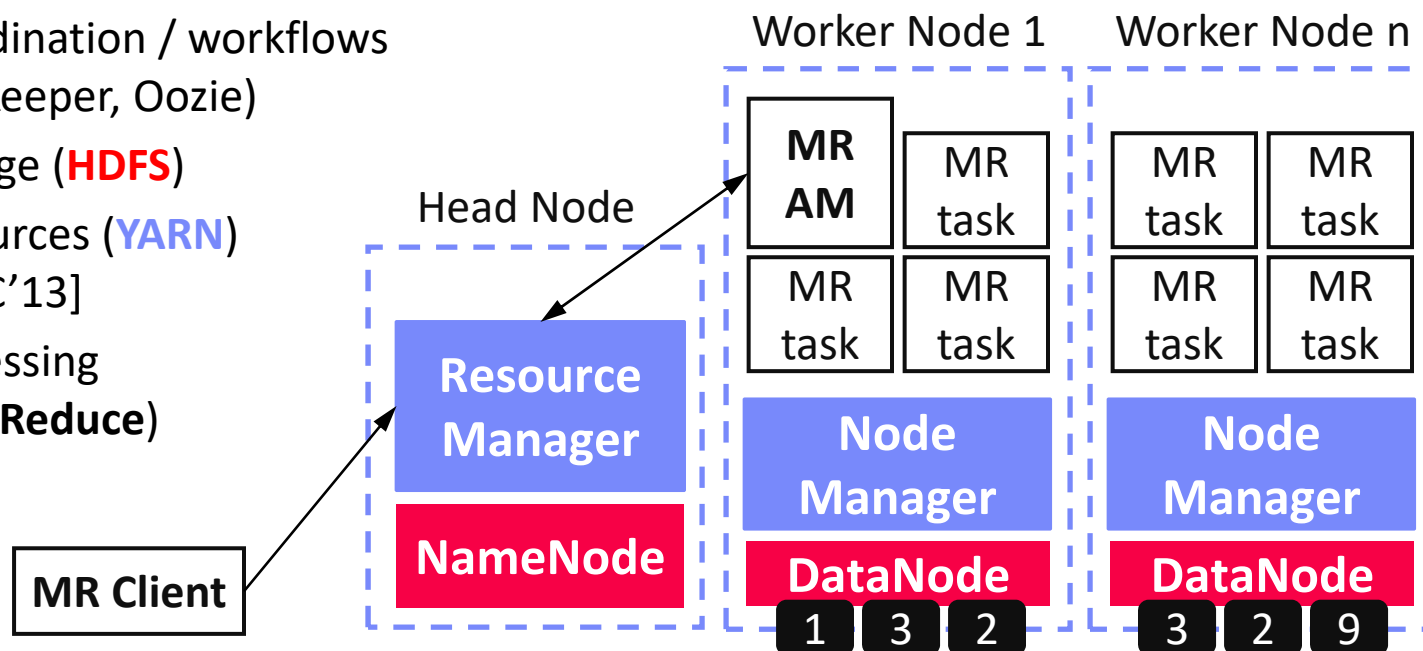


- Recap: Brief History

- Google's GFS [SOSP'03] + MapReduce [ODSI'04] → Apache Hadoop (2006)
- Apache Hive (SQL), Pig (ETL), Mahout (ML), Giraph (Graph)

- Hadoop Architecture / Eco System

- Management (Ambari)
- Coordination / workflows (Zookeeper, Oozie)
- Storage (**HDFS**)
- Resources (**YARN**) [SoCC'13]
- Processing (**MapReduce**)



# MapReduce – Programming Model

## Overview Programming Model

- Inspired by functional programming languages
- Implicit parallelism** (abstracts distributed storage and processing)
- Map** function: key/value pair  $\rightarrow$  set of intermediate key/value pairs
- Reduce** function: merge all intermediate values by key

## Example `SELECT Dep, count(*) FROM csv_files GROUP BY Dep`

Name	Dep
X	CS
Y	CS
A	EE
Z	CS

Collection of  
key/value pairs

```
map(Long pos, String line) {
  parts ← line.split(",")
  emit(parts[1], 1)
}
```

CS	1
CS	1
EE	1
CS	1

```
reduce(String dep,
  Iterator<Long> iter) {
  total ← iter.sum();
  emit(dep, total)
}
```

CS	3
EE	1

# MapReduce – Execution Model

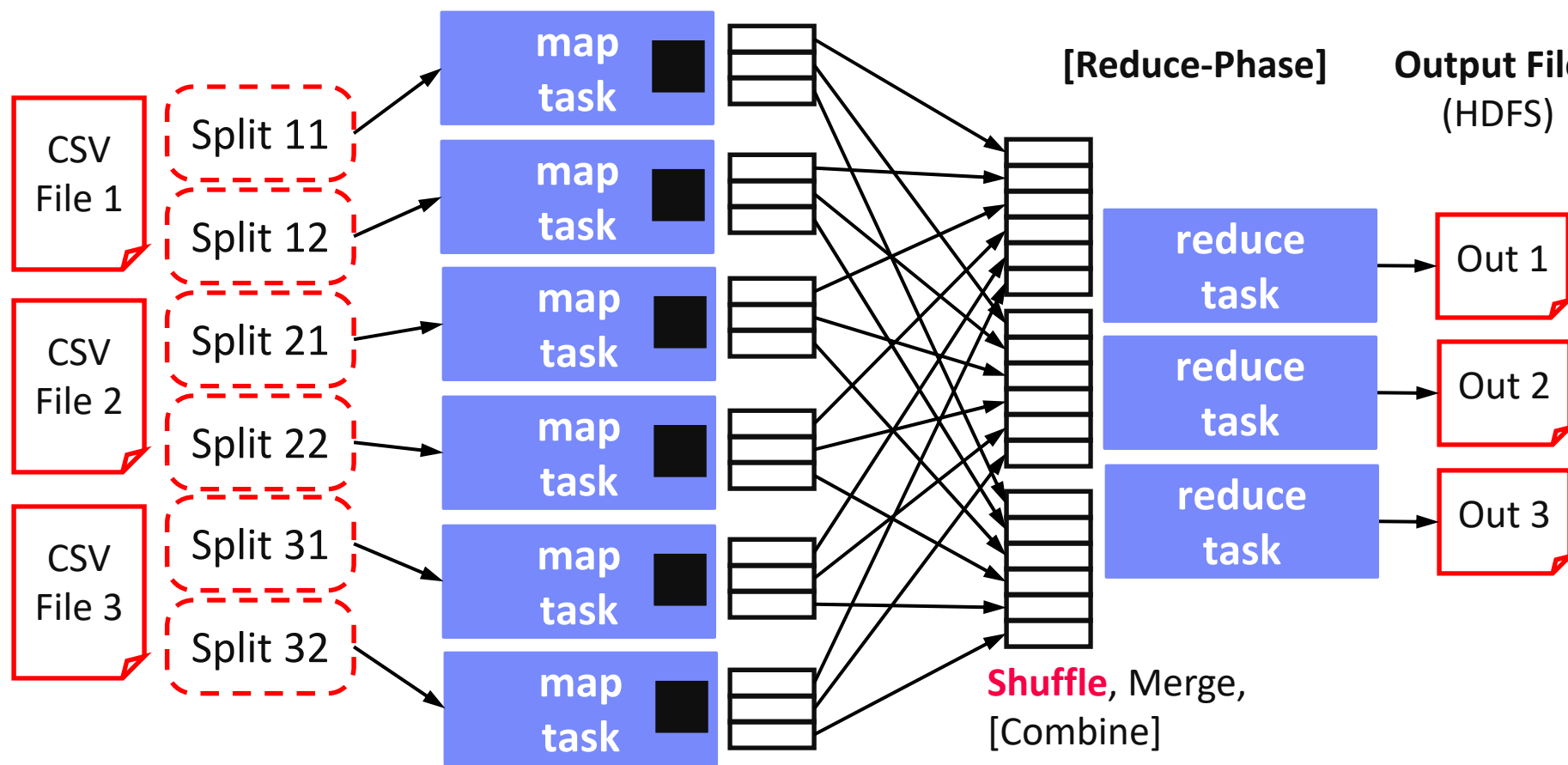
- #1 Data Locality (delay sched., write affinity)
- #2 Reduced shuffle (combine)
- #3 Fault tolerance (replication, attempts)

Input CSV files  
(stored in HDFS)

Map-Phase

[Reduce-Phase]

Output Files  
(HDFS)



Sort, [Combine], [Compress]


w/ #reducers = 3

# Spark History and Architecture

## ■ Summary MapReduce

- Large-scale & fault-tolerant processing w/ UDFs and files → **Flexibility**
- Restricted functional APIs → **Implicit parallelism and fault tolerance**
- **Criticism: #1 Performance, #2 Low-level APIs, #3 Many different systems**

## ■ Evolution to Spark (and Flink)

- Spark [HotCloud'10] + RDDs [NSDI'12] → **Apache Spark** (2014) 
- **Design: standing executors with in-memory storage**, lazy evaluation, and fault-tolerance via RDD lineage
- **Performance:** In-memory storage and fast job scheduling (100ms vs 10s)
- **APIs:** Richer functional APIs and general computation DAGs, high-level APIs (e.g., DataFrame/Dataset), unified platform

## ➔ But many shared concepts/infrastructure

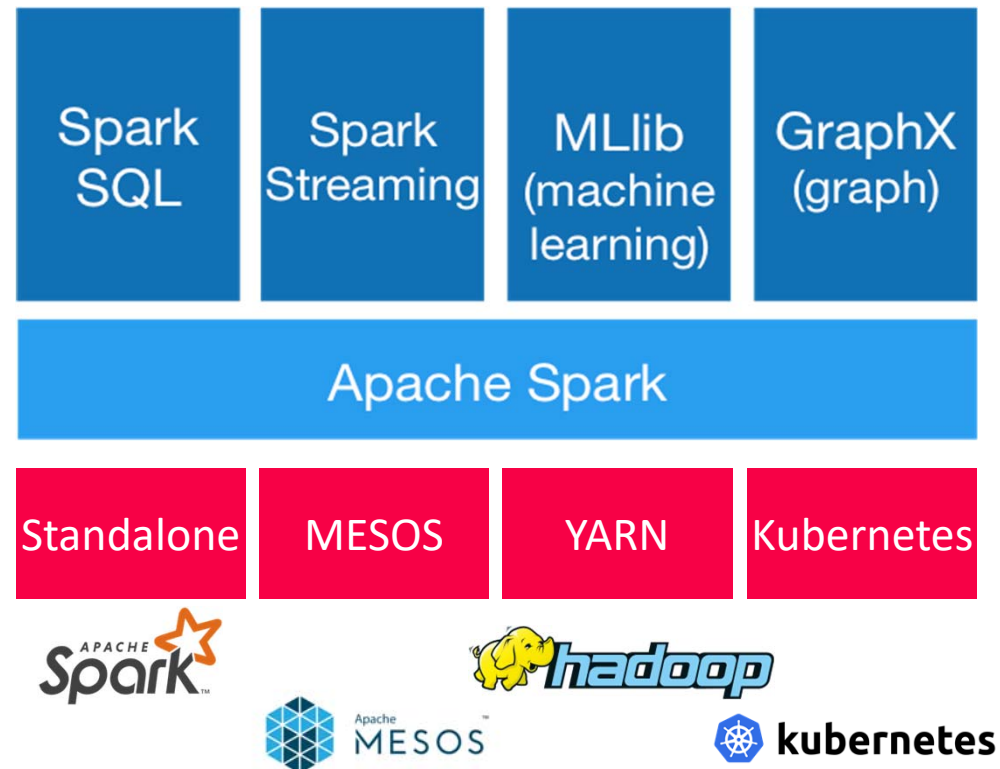
- **Implicit parallelism through dist. collections** (data access, fault tolerance)
- Resource negotiators (YARN, Mesos, Kubernetes)
- HDFS and object store connectors (e.g., Swift, S3)

# Spark History and Architecture, cont.

## High-Level Architecture

- **Different language bindings:**  
Scala, Java, Python, R
- **Different libraries:**  
SQL, ML, Stream, Graph
- Spark core (incl RDDs)
- **Different cluster managers:**  
Standalone, Mesos, **Yarn**, **Kubernetes**
- Different file systems/  
formats, and data sources:  
**HDFS**, **S3**, SWIFT, **DBs**, **NoSQL**

[<https://spark.apache.org/>]



- Focus on a **unified** platform for data-parallel computation

# Resilient Distributed Datasets (RDDs)

## ■ RDD Abstraction

- **Immutable**, partitioned  
collections of key-value pairs
- **Coarse-grained** deterministic operations (transformations/actions)
- Fault tolerance via lineage-based re-computation

JavaPairRDD

<MatrixIndexes,MatrixBlock>

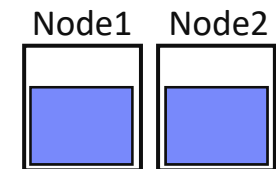
## ■ Operations

- Transformations:  
define new RDDs
- Actions: return  
result to driver

Type	Examples
Transformation ( <b>lazy</b> )	<code>map</code> , <code>hadoopFile</code> , <code>textFile</code> , <code>flatMap</code> , <code>filter</code> , <code>sample</code> , <code>join</code> , <code>groupByKey</code> , <code>cogroup</code> , <code>reduceByKey</code> , <code>cross</code> , <code>sortByKey</code> , <code>mapValues</code>
Action	<code>reduce</code> , <code>save</code> , <code>collect</code> , <code>count</code> , <code>lookupKey</code>

## ■ Distributed Caching

- Use fraction of worker **memory for caching**
- Eviction at granularity of individual partitions
- **Different storage levels** (e.g., mem/disk x serialization x compression)



# Partitions and Implicit/Explicit Partitioning

- **Spark Partitions**

- Logical key-value collections are split into **physical partitions**
- Partitions are granularity of **tasks, I/O** (HDFS blocks/files), **shuffling, evictions**

- **Partitioning via Partitioners**

- Implicitly on every data shuffling
- Explicitly via `R.repartition(n)`

**Example Hash Partitioning:**

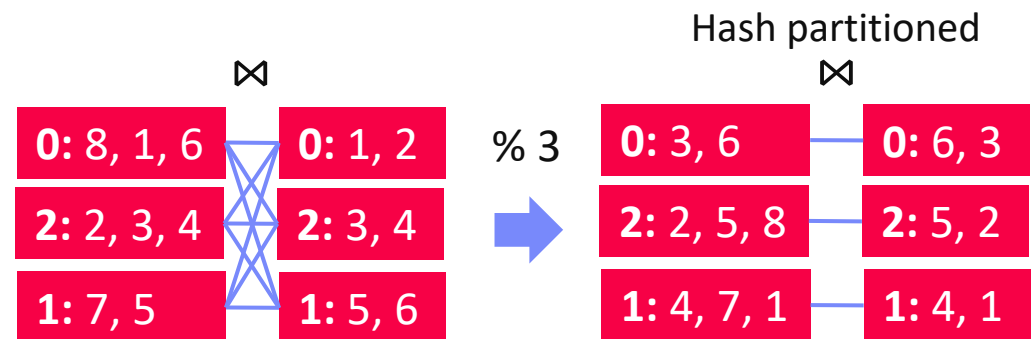
For all (k,v) of R:  
 $pid = hash(k) \% n$

- **Partitioning-Preserving**

- All operations that are guaranteed to keep keys unchanged (e.g. `mapValues()`, `mapPartitions()` w/ `preservesPart` flag)

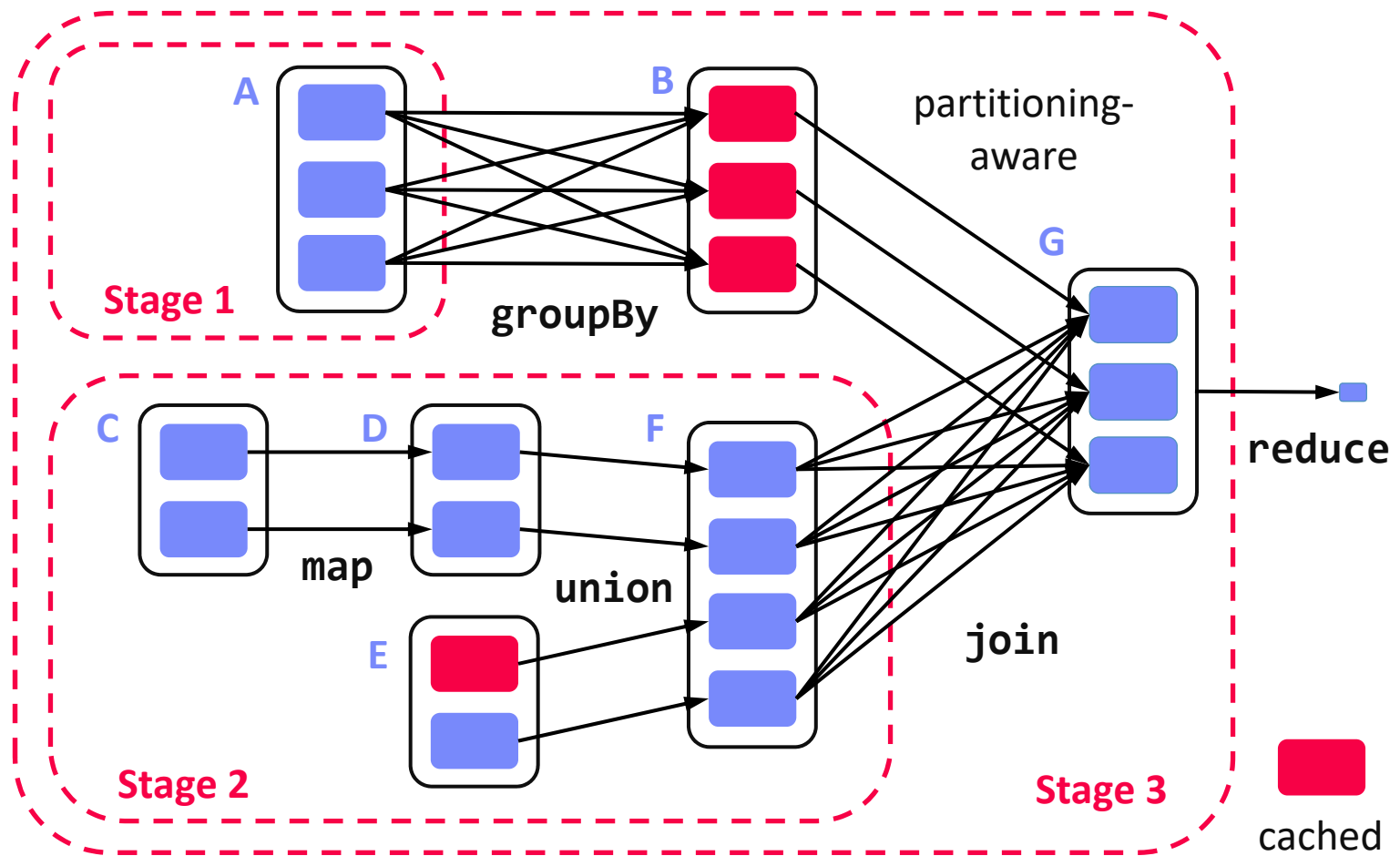
- **Partitioning-Exploiting**

- Join: `R3 = R1.join(R2)`
- Lookups: `v = C.lookup(k)`





# Lazy Evaluation, Caching, and Lineage



[Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauly, Michael J. Franklin, Scott Shenker, Ion Stoica: Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing. **NSDI 2012**]

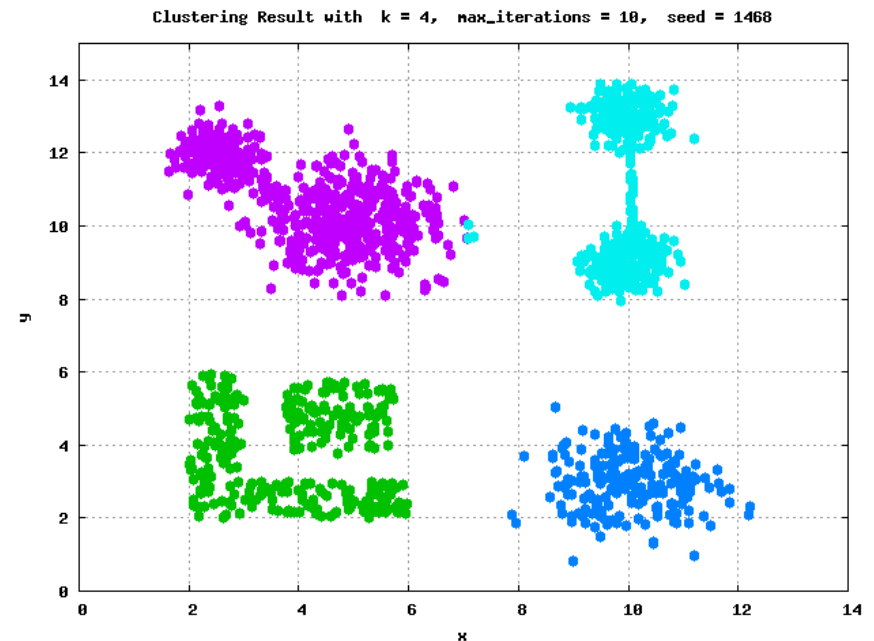
# Example: k-Means Clustering

## ■ k-Means Algorithm

- Given dataset  $D$  and number of clusters  $k$ , find cluster centroids (“mean” of assigned points) that minimize within-cluster variance
- Euclidean distance:  $\text{sqrt}(\text{sum}((\mathbf{a}-\mathbf{b})^2))$

## ■ Pseudo Code

```
function Kmeans(D, k, maxiter) {
  C' = randCentroids(D, k);
  C = {};
  i = 0; //until convergence
  while( C' != C & i<=maxiter ) {
    C = C';
    i = i + 1;
    A = getAssignments(D, C);
    C' = getCentroids(D, A, k);
  }
  return C'
}
```



## Example: K-Means Clustering in Spark

```
// create spark context (allocate configured executors)
JavaSparkContext sc = new JavaSparkContext();

// read and cache data, initialize centroids
JavaRDD<Row> D = sc.textFile("hdfs://user/mboehm/data/D.csv")
    .map(new ParseRow()).cache(); // cache data in spark executors
Map<Integer,Mean> C = asCentroidMap(D.takeSample(false, k));

// until convergence
while( !equals(C, C2) & i<=maxiter ) {
    C2 = C; i++;
    // assign points to closest centroid, recompute centroid
    Broadcast<Map<Integer,Row>> bC = sc.broadcast(C)
    C = D.mapToPair(new NearestAssignment(bC))
        .foldByKey(new Mean(0), new IncComputeCentroids())
        .collectAsMap();
}

return C;
```

Note: Existing library algorithm

[\[https://github.com/apache/spark/blob/master/mllib/src/main/scala/org/apache/spark/mllib/clustering/KMeans.scala\]](https://github.com/apache/spark/blob/master/mllib/src/main/scala/org/apache/spark/mllib/clustering/KMeans.scala)

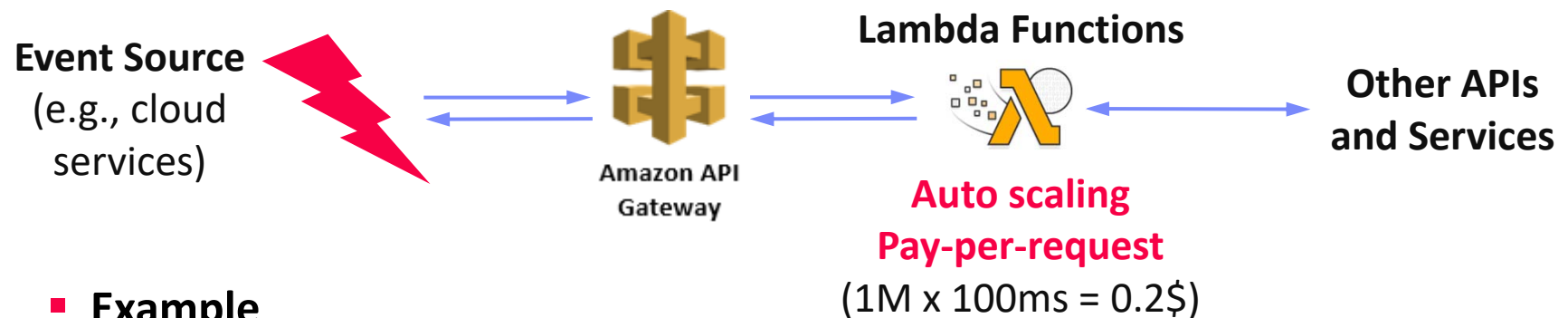
# Serverless Computing

[Joseph M. Hellerstein et al: Serverless Computing: **One Step Forward, Two Steps Back**. CIDR 2019]



## Definition Serverless

- **FaaS**: functions-as-a-service (event-driven, stateless input-output mapping)
- Infrastructure for deployment and auto-scaling of APIs/functions
- Examples: **Amazon Lambda**, **Microsoft Azure Functions**, etc



## Example

```
import com.amazonaws.services.lambda.runtime.Context;
import com.amazonaws.services.lambda.runtime.RequestHandler;

public class MyHandler implements RequestHandler<Tuple, MyResponse> {
    @Override
    public MyResponse handleRequest(Tuple input, Context context) {
        return expensiveStatelessComputation(input);
    }
}
```

# Conclusions and Q&A

- **Summary 11/12 Distributed Storage/Data Analysis**
  - Cloud Computing Overview
  - Distributed Storage
  - Distributed Data Analytics
  
- **Next Lectures (Part B: Modern Data Management)**
  - **13 Data stream processing systems [Jun 03]**
    - Including introduction of **Exercise 4**
  - **Jun 17: Q&A and exam preparation**