

## Tonight

- Clustering/Mahout
- Storm (streaming engines)
- Spark (in-memory engines)

## **Jaccard Similarity**

- Ratio between
  - Intersection of two sets and
  - Union of two sets

- Intuitive tests
  - Similarity S ♠S is 1

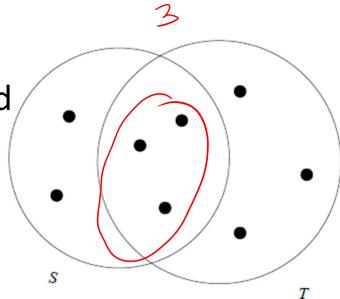


Figure 3.1: Two sets with Jaccard similarity 3/8

- Similarity S <> T with zero overlap is 0
- Otherwise similarity is between 0 and 1

#### Distance Measures

- Dist(x, y) >= 0
- Dist(x, y) = 0 iff x=y
- Dist(x,y) = Dist(y, x)
- Dist(x,y) <= Dist(x,z) + Dist(z, y) x</li>

Euclidean Distance

$$d([x_1, x_2, \dots, x_n], [y_1, y_2, \dots, y_n]) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

#### Jaccard Distance

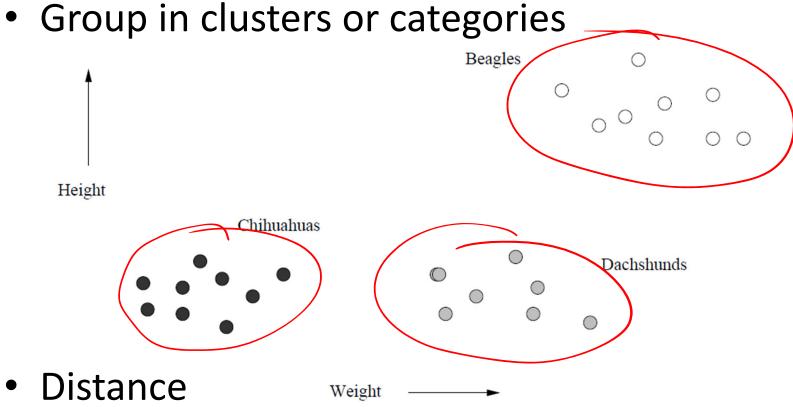
• (Dist(x, y) = 
$$1 - SIM(x, y)$$

- Cosine distance
  - Angle between vectors
- Edit distance
  - # of insertion/deletions for transforming x to y
- Hamming distance
  - # of elements that differ between x and y

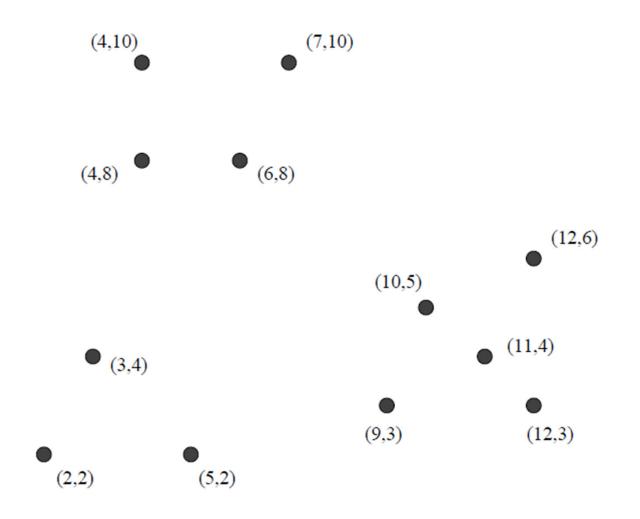
## Clustering

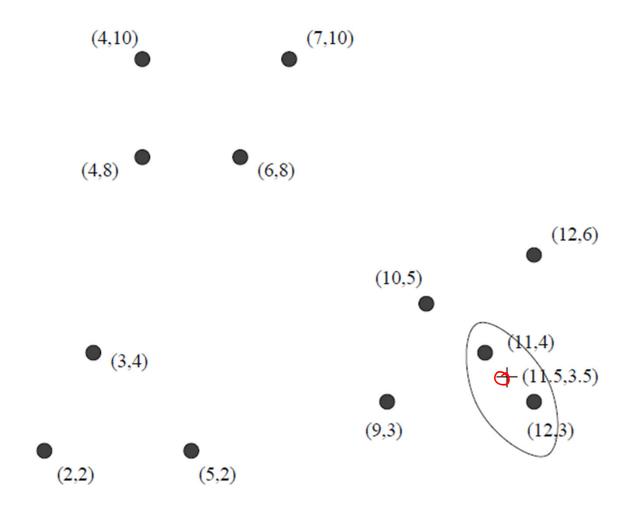
Set of elements in space

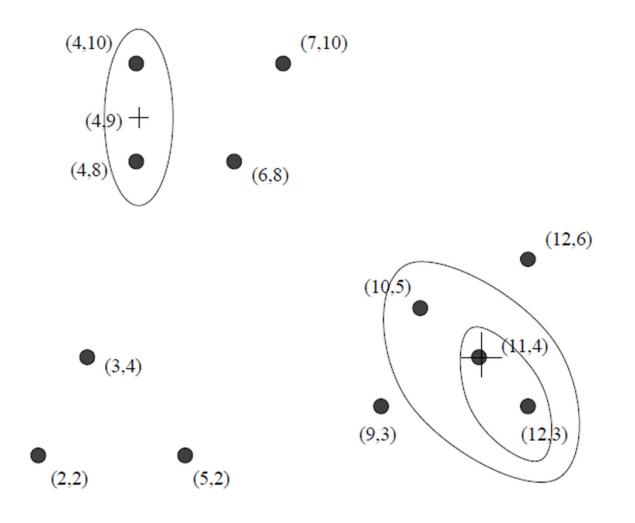
• Croup in clusters or categorie

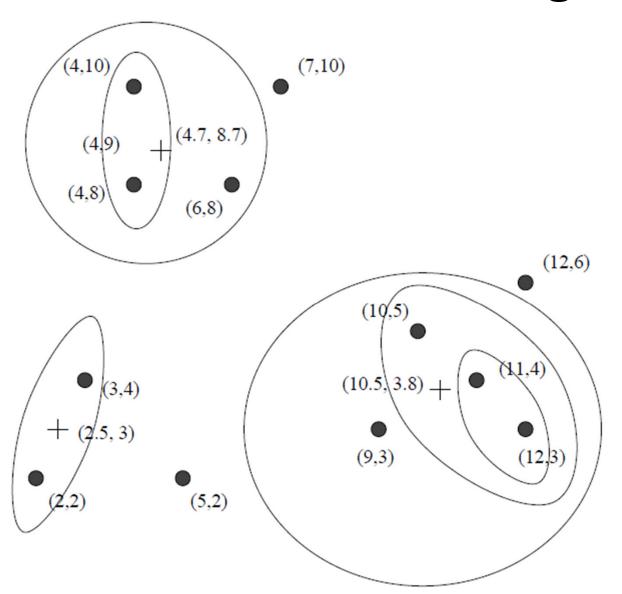


- Start from 1 point => 1 cluster
- Apply pairwise merging
  - Simpler distance measure requirements
- Stop condition
  - k # of clusters, cluster too large, etc.
- No back-tracking

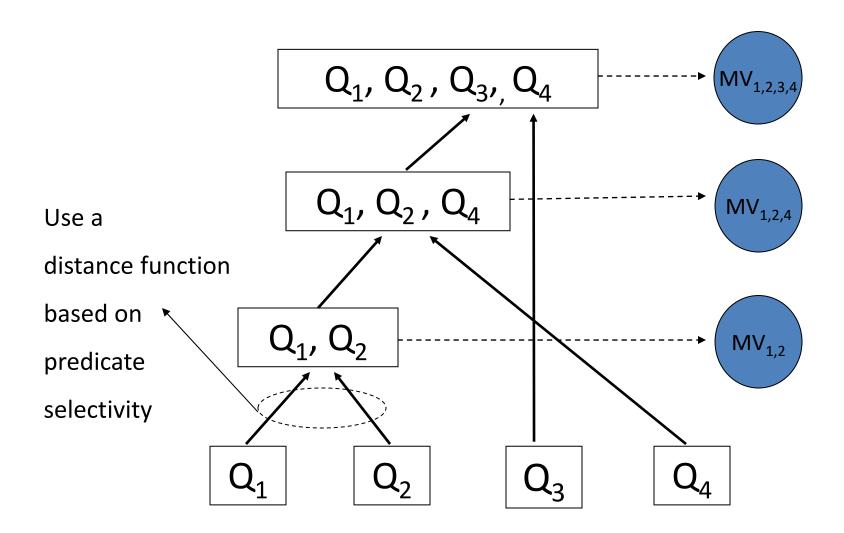








# Merging Materialized Views: Hierarchical Query Clustering



## Selectivity Vector for a SQL Query

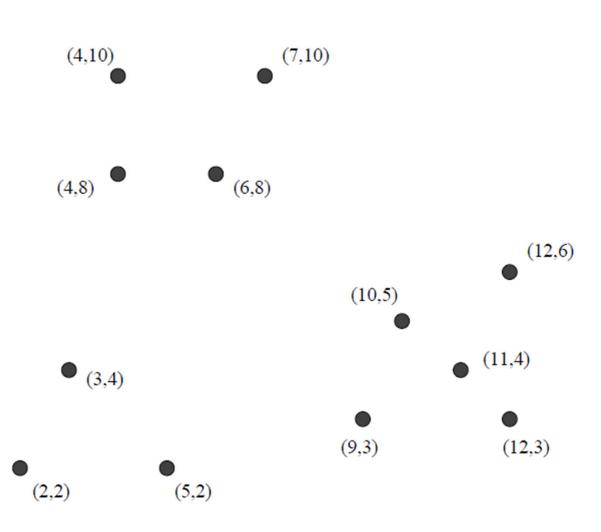
```
SELECT SUM(revenue)
FROM lineorder, dwdate, part, supplier
WHERE ... category = 'MFGR#12'
AND nation = 'CANADA' AND year = 1997;
(category: 0.27, nation: 0.20, year: 0.14)
                                            14%
(0.27, 1.0, 0.20, 0.14)
```

## K-Means Clustering

- Choose an initial set of centers
  - Well distributed (likely in different clusters)
  - Converges\*
- Assign every point to the closest centroid
- Find the new centroids
- Repeat
  - N steps
  - Or until converges

## K-Means Clustering

What is the best value for k?

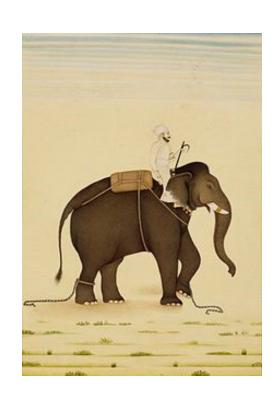


## K-Means Clustering

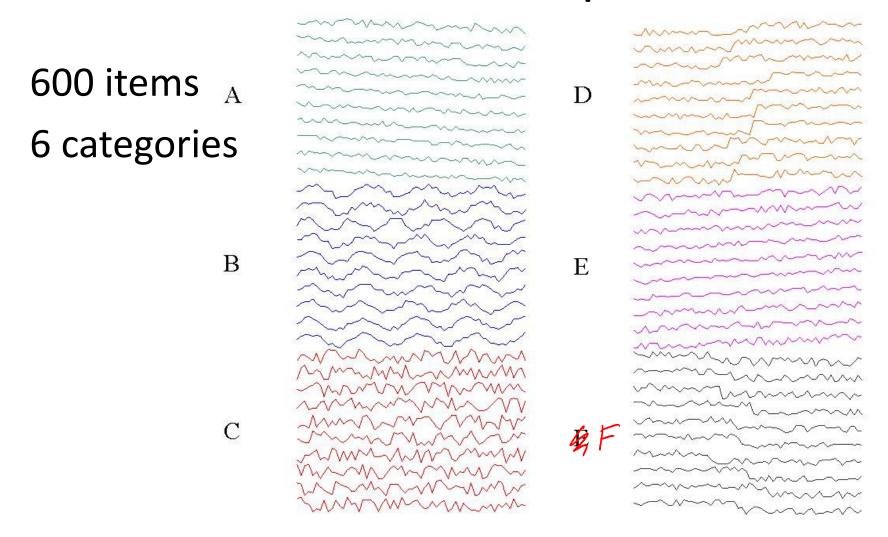
http://shabal.in/visuals/kmeans/4.html

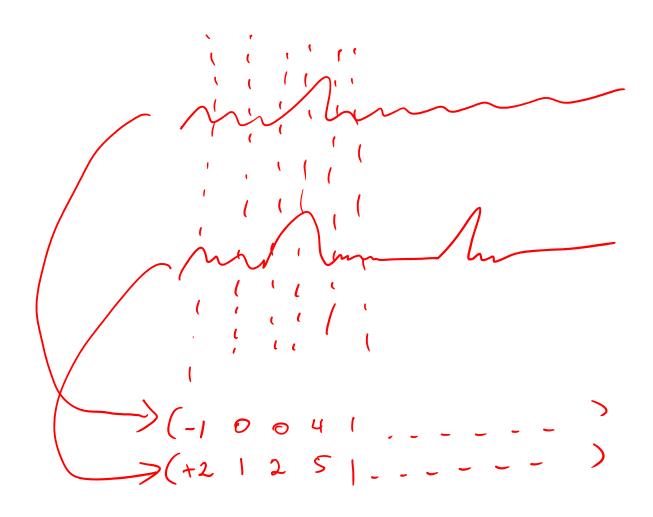
#### Mahout

- Machine learning library
- Community contribution
- Provides a number of algorithms
  - Clustering
    - K-Means, Fuzzy K-Means, ...
  - Recommenders
    - User and item based recommenders
  - Classifiers
- Runs on top of Hadoop/Spark/etc.



## Mahout Example

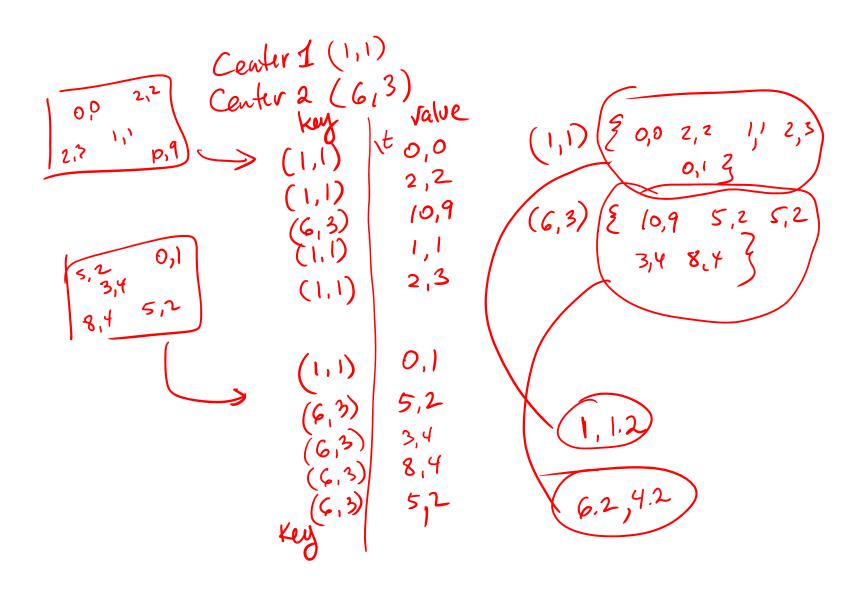




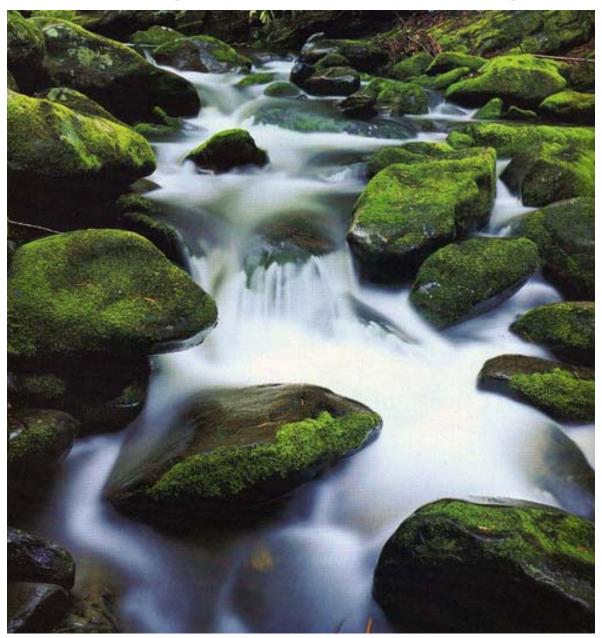
## Example KMeans Run

\$MAHOUT\_HOME/bin/mahout org.apache.mahout.clustering.syntheticcontrol.kmeans.Job

- Defaults testdata, Euclidian vector dist, 10 iterations
- KMeans Mapper
  - Assign the point to nearest cluster
- KMeansCombiner
  - Partially aggregate the points towards the new centroid
- KMeansReducer
  - Receive all keys and compute the new cluster centroid
- KMeansDriver iterate n times (or converge)



# **Streaming Database Engines**



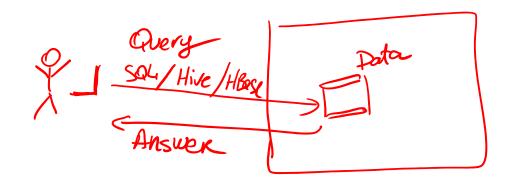
## Data-Stream Management

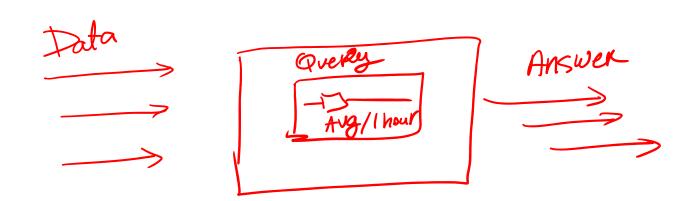
- Traditional DBMS
  - Data stored in finite, persistent data sets
- Data streams
  - Distributed, continuous, unbounded, time-varying
- Data-Stream Management
  - Network monitoring
  - Data security
  - Sensor network
  - Financial applications

#### **Streaming Data**

- Types of data (infinite)
  - Sensor data
  - Stock market quotes
  - Network monitoring data
- Types of queries
  - One pass (cannot rewind or store)
  - Report last day average
  - Report changes in sensor data
  - Collect information about single user

# HDFS/Oracle/Postgre SQ4





## Data-stream-management System

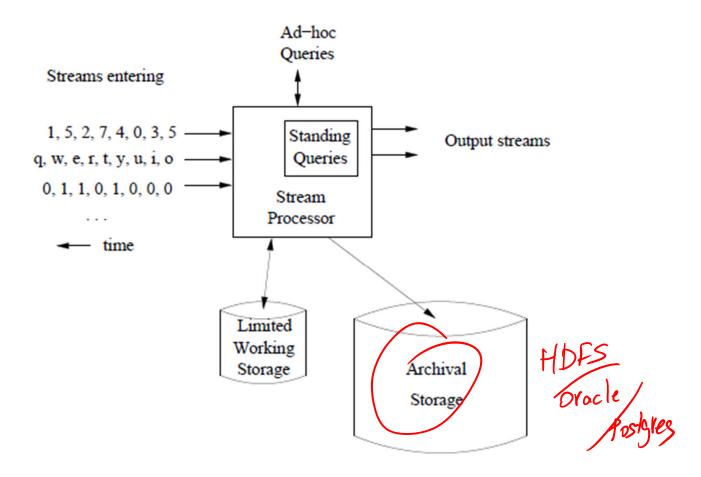
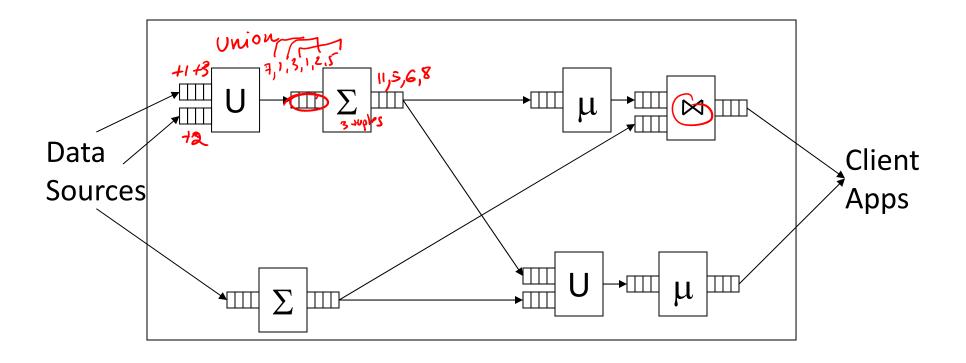


Figure 4.1: A data-stream-management system

## Stream Processing



- Data sources push tuples continuously
- Operators process windows of tuples

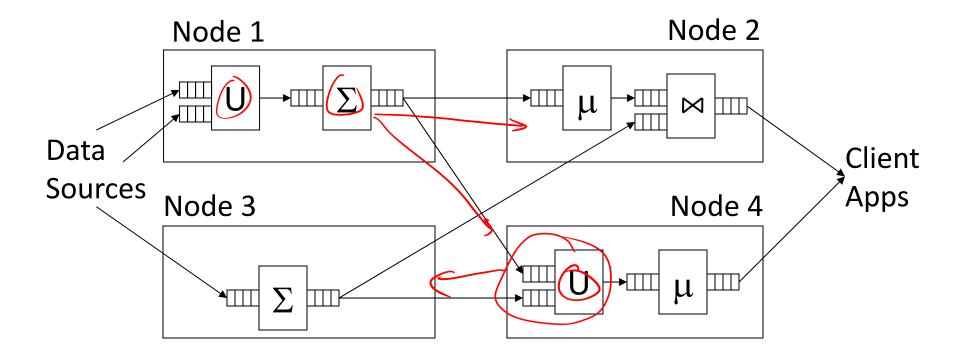
## Examples

- Ocean temperature sensors
  - GPS+temperature every 1/10 second (3.5MB)
  - More sensors (x1000, or x1000000)
- Image data
  - London has millions cameras
- Internet monitoring
  - Intruder alert
  - Any real-time activity spike, traffic

## **Operations**

- Filter the stream
  - E.g., select (temperature > 100)
- Aggregate the stream
  - # of cars entered the garage
    - Window (per-hour? Over last hour?)
- Join two streams
  - Match "buy" and "sell" offers

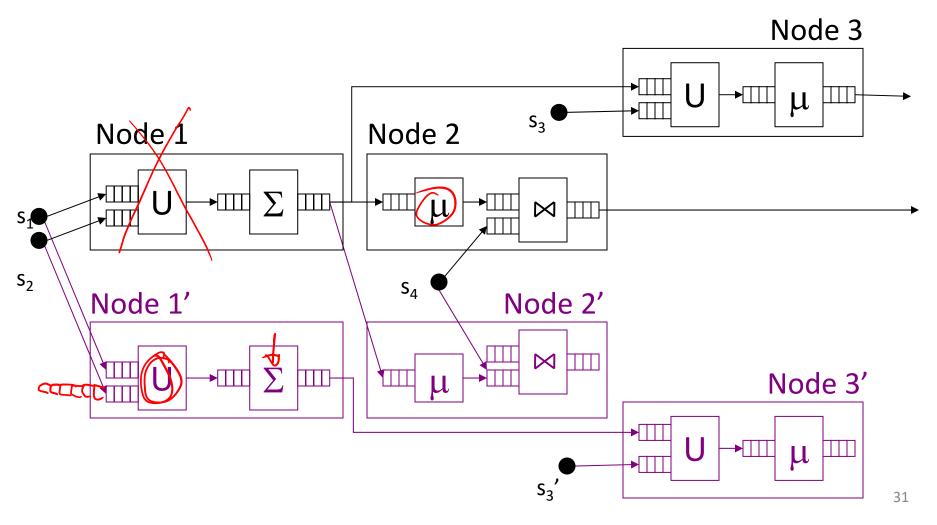
## Distributed Stream Processing

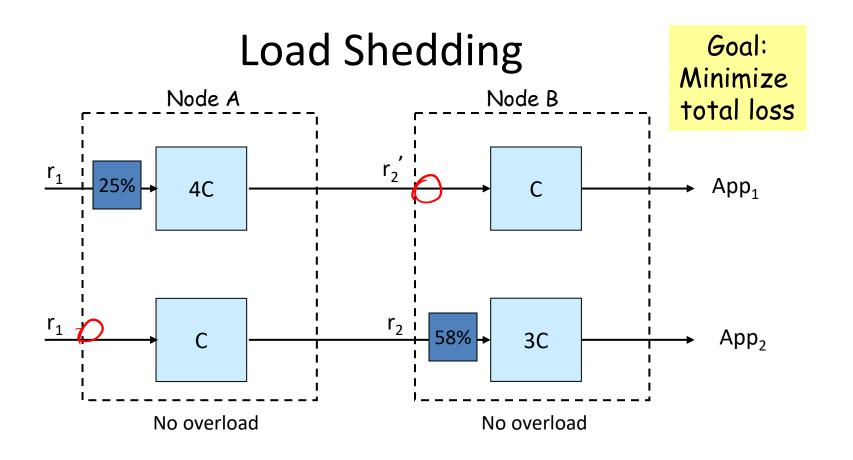


- Data sources push tuples continuously
- Operators process windows of tuples

### Fault-Tolerance through Replication

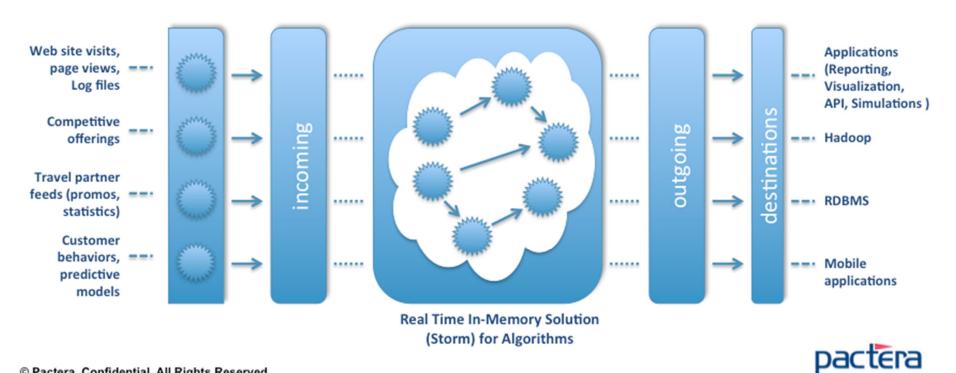
Goal: Tolerate node and network failures





Plan	Loss
Local	App <sub>1</sub> : 25% App <sub>2</sub> : 58%

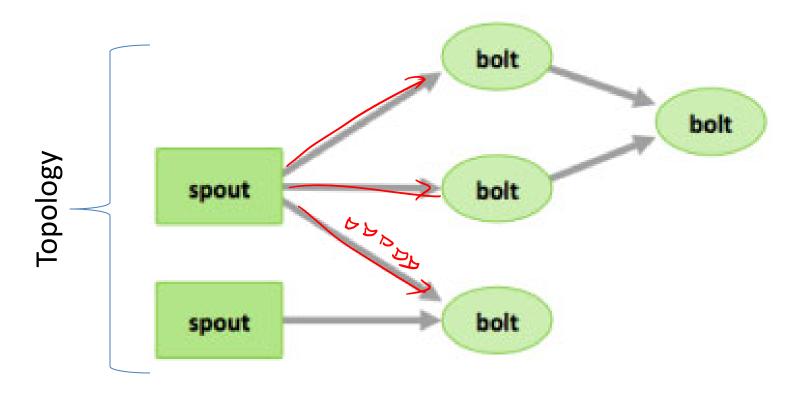
#### Storm



© Pactera. Confidential. All Rights Reserved.

## **Storm Components**

- Spout = source of tuples
- Bolt = processing operator

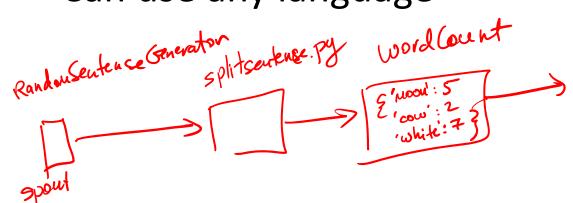


## **Storm Components**

- Nimbus node
  - JobTracker
- Supervisor Nodes
- ZooKeeper Nodes
  - Coordination
- Stream
  - Unbound sequence of tuples

## **Running Storm**

- Submit "topologies"
  - WordCountTopology.java
- Spouts / Bolts
- Can use any language



## **In-Memory Databases**

Main Memory database system (MMDB)

- Data resides permanently on main physical memory
- Backup copy on disk

Disk Resident database system (DRDB)

- Data resides on disk
- Data may be cached into memory for access
   Main difference is that in MMDB, the primary copy lives permanently in memory

#### MMDB?

- Is it reasonable to assume that the entire database fits in memory?
  - Yes, for some applications
  - -100 machines x 64GB RAM = 6.4TB RAM
- What is the difference between a MMDB and a regular database with a very large cache?
  - Even if all data fits in memory, the structures and algorithms are designed for disk access.

# Differences in properties of main memory and disk

- The access time for main memory is orders of magnitude less than for disk storage
- Main memory is normally volatile, while disk storage is not
- The layout of data on disk is much more critical than the layout of data in main memory

## Spark

- General engine for large-scale data processing
  - Batch/real-time/stream
  - Machine learning
  - Graph workloads
- In memory processing
  - Resilient Distribute Dataset (RDD)
  - HDFS or in-memory\*

## Spark

- bin/pyspark (shell)
  - textFile = sc.textFile("README.md")
  - textFile.count() # number of items/lines
  - textFile.first() # first item
  - linesWithSpark =
     textFile.filter(lambda line: "Spark" in line)

**Spark Principles** 

Madine 1
Machine3

- Distributed datasets
  - data = [1, 2, 3, 4, 5]
  - distData = sc.parallelize(data)
  - distData.reduce(lambda a, b: a + b)
- distFile = sc.textFile("data.txt")
  - hdfs://... , s3n://
  - lengths = distFile.map(lambda s: len(s))
  - totalLength = lengths.reduce(lambda a, b: a + b)

## Spark Setup

- Download Spark binary
  - http://spark.apache.org/downloads.html
- Mode
  - Using stand-alone mode
  - YARN
    - Spark pre-built with support, security

## Spark Setup

- sbin/start-all.sh
  - Workers file
- http://public-IP:8080/
  - Cluster view

## Spark Example

```
# Read file from HDFS
text_file = sc.textFile("hdfs://ec2-18-144-12-192.us-west-
1.compute.amazonaws.com/data/README.md")
lengths = text file.map(lambda s: len(s))
print text file.take(100)
lengths.foreach(myprint)
print lengths.take(100)
totalLength = lengths.reduce(lambda a, b: a + b)
```