
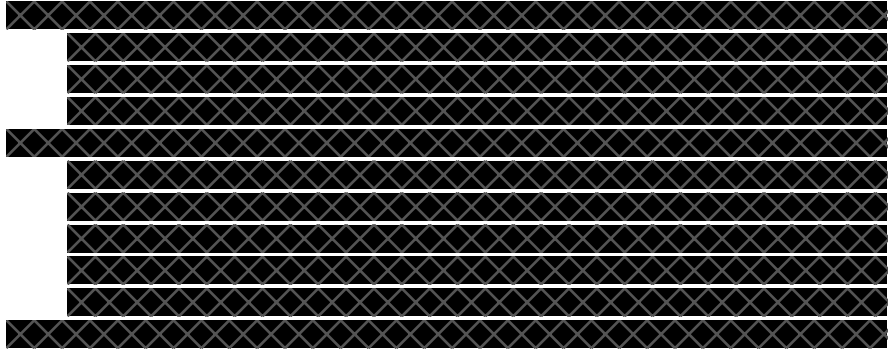


Disclaimer: these notes are meant to be helpful and to expand the core course materials. They are not necessarily complete and representing the whole content of the course. They are built on a mix of course sources and students' notes. Use them just as additional source of information.

## Contents

	4
<b>2 Introduction</b>	<b>5</b>
2.1 Data Driven Decision making . . . . .	5
2.2 The Definition of Big Data . . . . .	6
2.3 The Roles in the Data World . . . . .	6
<b>3 Unstructured Data Models</b>	<b>7</b>
3.1 NOSQL General Concepts . . . . .	7
3.1.1 Differences with the traditional data model . . . . .	7
3.1.2 Schema Less Approach . . . . .	7
3.1.3 Data Lakes . . . . .	8
3.1.4 Scalability . . . . .	8
3.1.5 Transactional Properties in NoSQL . . . . .	9
3.2 Graph Stores . . . . .	12
3.2.1 Graph Theory . . . . .	12
3.2.2 Graph Databases . . . . .	14
3.2.3 Neo4J . . . . .	15
3.3 Key-Value Stores . . . . .	19
3.3.1 Introduction . . . . .	19
3.3.2 Redis . . . . .	19
3.4 Columnar Databases . . . . .	23
3.4.1 Introduction . . . . .	23
3.4.2 Cassandra . . . . .	25
3.4.3 Cassandra's Architecture . . . . .	28
3.5 Document Databases . . . . .	32
3.5.1 Intoduction . . . . .	32
3.5.2 MongoDB . . . . .	33
<b>4 Streaming Data Engineering</b>	<b>35</b>
4.1 Introduction . . . . .	35
4.2 Streaming . . . . .	39
	
4.6 Spark . . . . .	63

4.6.1	Introduction . . . . .	63
4.6.2	Spark APIs . . . . .	64
<b>5</b>	<b>Data Pipelines</b>	<b>78</b>
5.1	Data Ingestion . . . . .	78
5.1.1	Introduction . . . . .	78
5.1.2	Web APIs . . . . .	78
5.1.3	Web Scraping . . . . .	81
5.2	Data Wrangling . . . . .	82
5.2.1	The need for clean data . . . . .	82
5.2.2	Data cleansing . . . . .	83
5.3	Crowdsourcing . . . . .	87
5.3.1	Human Computation . . . . .	87
5.3.2	Gamification . . . . .	89

## 2 Introduction

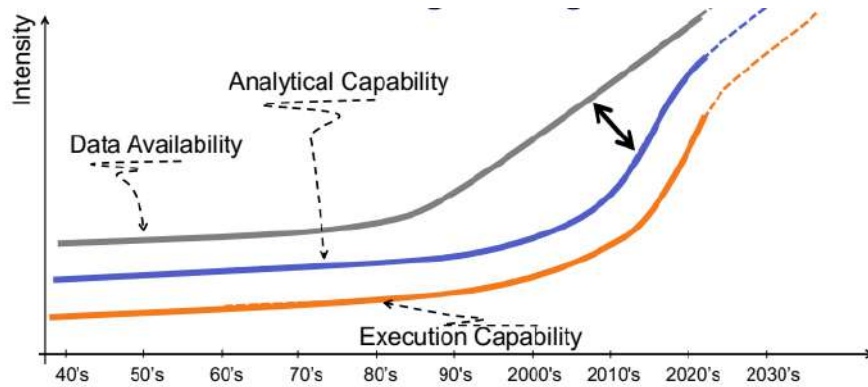
### 2.1 Data Driven Decision making

Nowadays companies need a reliable source of data for performing decision making, and big enterprises have huge amounts of data from which they need to extrapolate valuable information.

With the amount of data derived from IT systems companies can now structure their organization in a data-driven way, as opposed to in the past when they had to take decisions based on instinct.

Data-driven decisions are more effective, predictable and profitable. Such an approach is only possible if we have a great amount of data to analyze and if data is organized and stored in a good way.

The available data is an upper bound for the analytical capability of a company, which in turn is an upper bound to the execution capability.



## 2.2 The Definition of Big Data

So in order to be able to have a data-driven approach we need to be able to handle the so called "Big Data", a term used to describe data with the following characteristics ("The 4 V's"):

- Volume: data that ranges in size from Terabytes to Petabytes (in Italy a company with tens of Terabytes of analytical data is considered pretty big)
- Variety: data that can assume many different forms, from structured (like relational tables), to semi-structured (like JSON and XML) and even completely unstructured (text and multimedia files)
- Velocity: information that flows continually creating streams, that often needs to be analyzed in a time frame of a few seconds
- Veracity: big data can often be imprecise and unpredictable, and this unreliability needs to be addressed and managed

## 2.3 The Roles in the Data World

The data from the operational systems of a company is a raw material that needs to be refined using Data Science techniques in order to extract valuable information that can fuel the decisional process of a company and so the duties of a Data Scientist are:

- Obtaining predictable, actionable insights from data
- Creating data product with immediate impact
- Communicating relevant business stories from data
- Building confidence in decisions that drive business value

Instead Data Engineers are the ones who build the underlying system that allows Data Scientist to analyze, develop and deploy working and ever-evolving data models.

In other words a Data Engineer handles the lifecycle of data in a processing pipeline, ensuring that it is always available and usable by whoever needs it.

To summarize, these are the typical responsibilities of a DE:

- Design and build data processing systems
- Operationalize machine learning models
- Ensure solution quality

## 3 Unstructured Data Models

### 3.1 NOSQL General Concepts

#### 3.1.1 Differences with the traditional data model

The Big Data world has a lot of differences from the traditional data managing approach. First we need to move from the relational world to new kinds of unstructured data storage options.

Also there is a shift in the hardware technologies, going from the typical vertically scalable, on premise mainframes to cloud based, horizontally scalable commodity machines

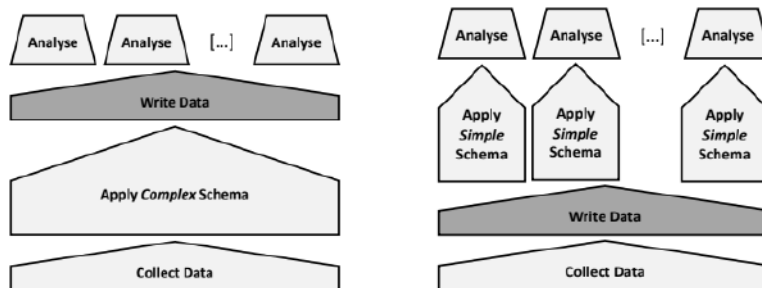
Lastly there is a change in data processing culture, from a rear view reporting approach based on relational algebra to a smarter model of continuous analytics with streaming and machine learning techniques.

#### 3.1.2 Schema Less Approach

This new types of data model require lot of flexibility, which is usually implemented by getting rid of the traditional fixed schema of the relational model, using a so called "schema-less" approach where data doesn't need to strictly respect a predefined structure in order to be stored in the database.

This way the definition of the schema is postponed from the moment in which the data is written to the moment when the data is read. This approach is called Schema-On-Read as opposed to the traditional Schema-On-Write. This new philosophy has many advantages over the old one:

- Since we don't have to check the correctness of the schema the database writes are much faster
- We have a smaller IO throughput since when reading data we can only fetch the properties we need for that particular task
- The same data can assume a different schema according to the needs of the analytical jobs that is reading it, optimizing performance.



### 3.1.3 Data Lakes

The better define this new kind of data storage systems,the concept of Data Lake was invented: a reservoir of data, which is just a rough container, used to support the needs of the company.

A Data Lake has incoming flows of both structured and unstructured data and outgoing flows consumed by the analytical jobs.

This way organizations have only one centralized data source for all its analytical needs,as opposed to silos based approaches.

Building such a huge and centralized data storage system has many challenges: it's important to organize and index well all the incoming data,or we could end up with a huge Data Swamp where it's impossible for everyone to find what they're looking for.

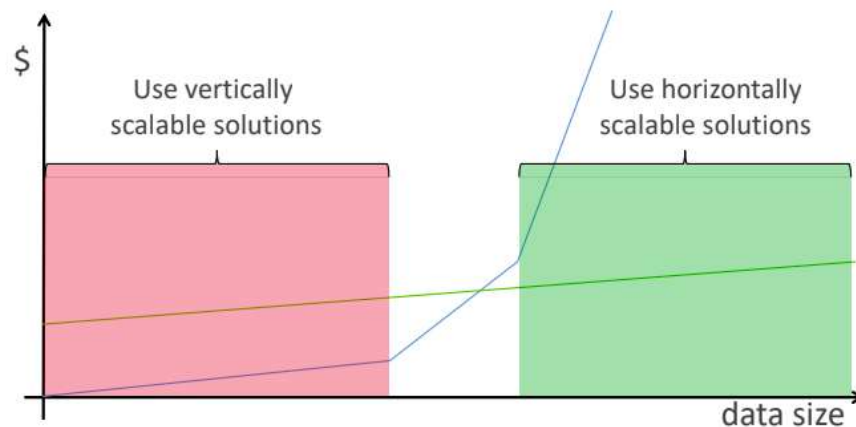
And so usually in a Data Lake the incoming raw data gets incrementally refined and enriched in order to improve accessibility,by adding metadata and indexing tags and by performing quality checks.

### 3.1.4 Scalability

As we previously mentioned,traditional data systems scale vertically,meaning that when the machine doesn't have enough ram or disk to store all the data it just gets replaced with a more powerful and bigger one.

This approach only scales up to a certain point and after that it just isn't efficient anymore,so when dealing with Big Data we need to scale horizontally instead,meaning that the computing system is composed of many commodity machines and when there's a need for more resources more machines are added without replacing the old ones.

## Vertical vs. Horizontal scalability



As can be seen from the graph,in vertical scalability cost scales exponentially with respect to data size,while in horizontal scalability cost scales linearly.

### 3.1.5 Transactional Properties in NoSQL

In the relational world we are used to having the concept of a transaction, an elementary unit of work encapsulated by begin and commit commands characterized by ACID properties:

- Atomicity: the operations contained in a transaction either all fail or all succeed
- Consistency: the state of the db respects the integrity constraints before and after the transaction is executed
- Isolation: every transaction doesn't affect and isn't affected by other concurrent transactions
- Durability: a transaction produces durable changes in a db

These properties are very important and often even fundamental in a traditional SQL based OLTP application

Big Data systems who have an architecture based on horizontal scalability and distributed systems unfortunately can't give complete ACID transactions guarantees.

This is what essentially stated by the CAP theorem.

CAP is a word composed by the initials of 3 important features in distributed systems:

- Consistency: all nodes see the same data at the same time
- Availability: Node failures do not prevent other survivors from continuing to operate (a guarantee that every request receives a response about whether it succeeded or failed)
- Partition Tolerance: the system continues to operate despite arbitrary partitioning due to network failures (e.g., message loss)

A distributed system can satisfy any two of these guarantees at the same time but not all three. In big data system who rely heavily on network communication Partition Tolerance is essential, so designing a distributed system means having to handle a tradeoff between Availability and Consistency. The extreme choice would be to abandon completely one between these 2 features, but in real system the data engineer needs to finetune the level of Consistency and Availability to the need of the service.



- AP: A partitioned node returns a correct value, if in a consistent state; a timeout error or an error, otherwise e.g., DynamoDB, CouchDB, and Cassandra (typical use case: banking applications)
- CP: A partitioned node returns the most recent version of the data, which could be stale. e.g., MongoDB, Redis, AppFabric Caching, and MemcacheDB (typical use cases: media streaming, statistics, social media, news websites)

The traditional relational systems, being monolithic don't need Partition Tolerance and so they are both Available and Consistent (AC). To address the need for distributed big data system new data model and technologies have been invented.



As we mentioned previously these new kind of technologies (which are called NoSQL, that can mean both not SQL and not only SQL) cannot give the same ACID guarantees of SQL databases, so a new weaker and generic set of characteristics called BASE has been invented to describe features of big data systems.

- Basic Availability: the system can always fulfill requests, but the answer could be partially consistent
- Soft State: The solid state of relational system is abandoned
- Eventual Consistency : At some point in the future data will converge to a consistent state (consistent state is not granted immediately but eventually)

These features are voluntarily vague and generic since they can be finetuned to the needs of the application as stated by the CAP theorem.

## 3.2 Graph Stores

The first kind of databases to have success after the monopoly of SQL are the graph storage solutions.

We'll now have an introduction about graph as a data structure in general, and we will then go deeper into the theory behind the graph based databases and finally talk about a real world graph DB called Neo4J.

### 3.2.1 Graph Theory

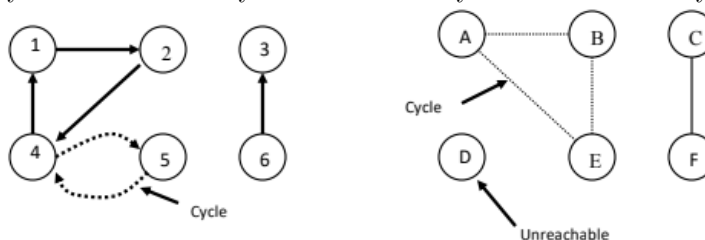
The Graph is a data structure that was first used in 1736 to represent a city and its water canals and have been used ever since for a huge number of optimizations of problems that can be modeled has a set of entities connected by a relation of some kind. From a mathematical point of view, a graph is an ordered triple  $G:(V,E,f)$  where

- $V$  is a set of nodes/vertex.
- $E$  is a set of arcs/edges.
- $f$  is a function that maps every edge of  $E$  to an unordered pair of nodes of  $V$

A path is a sequence of consecutive vertex, and a node  $A$  is reachable from another one  $B$  if exists a path that goes from  $A$  to  $B$ .

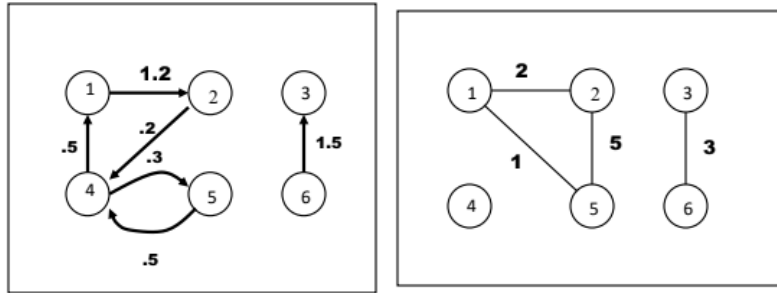
If a path starts from a node and returns to the same node its called a cycle, a graph is called acyclic if there it doesn't have any cycle, cyclic otherwise.

A graph is connected if every node is reachable from any other node and is strongly connected if every node it's directly connected to every other one.



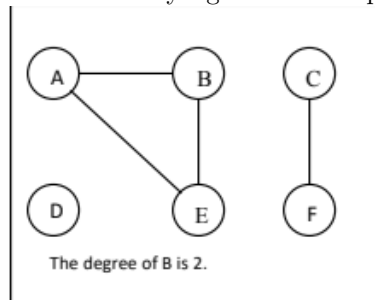
If a graph has  $n$  nodes then it can have up to  $v = n(n-1)$  edges, if  $v \approx n(n-1)$  then the graph is called dense, otherwise it's called sparse.

Nodes and vertex can be enriched by putting additional data on them, the latter case for example is commonly used to represent some kind of cost of the connection (in this case we have a weighted graph).



If the edges can have a direction, meaning that given the two nodes A and B  $(A,B) \neq (B,A)$  can be true, then the graph is called directed.

The number of incident edges in a node it's called the degree of a node; this is a very simple metric but it can be very significant in representing the importance



of a node in a graph.

One last important concept is the subgraph, which is a graph whose vertex and edges are a subset of a supergraph. The physical implementation of a graph in a computer can be done in many ways, the most famous of which are the adjacency matrix (where  $a_{ij} = 1$  if  $(i,j)$  belongs to E and  $a_{ij} = 0$  otherwise) which is efficient for dense graphs and the edge list, more efficient for sparse ones instead.

### 3.2.2 Graph Databases

So why did the need for a graph based database arise? The answer is, somewhat surprisingly, that relational databases are not good at managing relationships! The table based structure of relational databases makes it hard to represent relationships between rows in the same table, and moreover whenever someone needs to find a relationship between records of different tables the db has to perform a JOIN operation, which is usually very expensive.

So for the use cases in which relationships are the most important feature of our data (e.g. social network friendships) it would be best to go with a technology who can implement relationships in a native and efficient way, and that's where graph dbs come in.

In these kind of storage systems data are represented as entities connected by information rich relations, just like in a real graph.

- Sailor [Reserves] Boat



- (:Sailor) -[:reserves]-> (:Boat)

The typical query based approach for managing data in relational databases isn't a good fit for graph dbs, in these cases it's preferred a different approach called Graph Matching, or more in general Pattern Matching. In these technique the user specifies a particular shape or structure he wants to find in the graph, and the system searches for and returns all the subgraphs that match that request.

### 3.2.3 Neo4J

The most popular graph database is Neo4J, implemented in Java by Neo Technologies, and it has the following characteristics:

- Operational db
- ACID guarantees
- Not efficient for large scale graph analysis
- Nodes and edges are a native feature
- Each node has an identifier tag and can have many properties
- Node can have types, called labels
- Schemaless (nodes with new types can be created at any time)
- traditional CA architecture

Queries are expressed with a custom made declarative language called Cypher. With Cypher is easy to express queries based on relationships (who are the main focus in a graph db), and it's much more efficient than SQL in doing operations equivalent to the JOINS in relational databases

Let's start now making some examples of Cypher Commands:

Create a node of type Crew and label Neo, with attribute name of value 'Neo':

```
CREATE (Neo:Crew {name:'Neo'})
```

Add an edge of type "Knows" from the node Neo to the node Morpheus

```
(Neo) - [:KNOWS] -> (Morpheus)
```

Create and index on a Node attribute (used as a starting point for a query)

```
CREATE INDEX ON :Customer(customerID)
```

Create a constraint enforcing uniqueness of attribute customerID on nodes with type customer:

```
CREATE CONSTRAINT ON (c:Customer)  
ASSERT c.customerID IS UNIQUE;
```

Import Data from CSV file:

```
LOAD CSV WITH HEADERS FROM "file:customers.csv" AS row  
CREATE (:Customer {companyName: row.CompanyName,  
customerID: row.CustomerID, phone: row.Phone});
```

Create Node with multiple labels:

```
CREATE (n:Actor:Director {name:'Clint Eastwood'})
```

Merge operator, creates new nodes only if it doesn't exist another node with the same label:

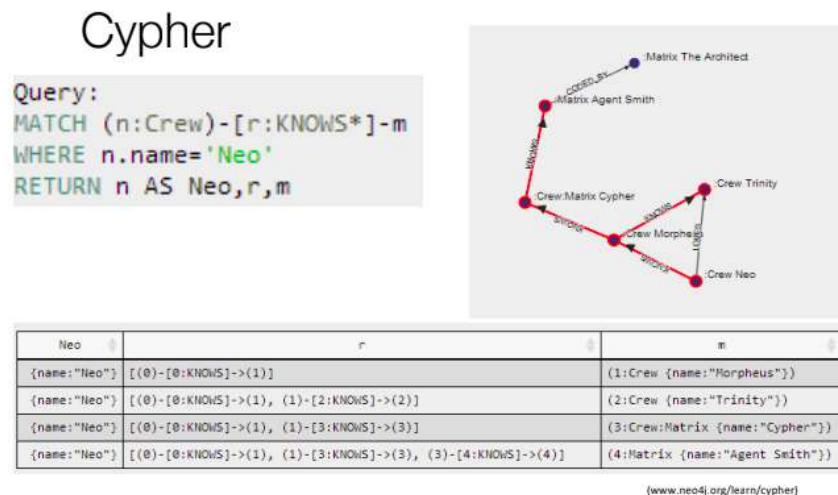
```
LOAD CSV WITH HEADERS FROM "file:///transfers.csv" AS Row
MERGE (player:Player {id:row.playerUri})
ON CREATE SET player.name = row.name,
              player.position = row.playerPosition
```

Let's now take a look at Cypher's query system, which is as previously mentioned based on pattern matching and has the following structure:

```
START
MATCH Pattern Matching
WHERE Expressions, Predicates
RETURN Output
```

Match describes the shape of the subgraph we are looking for (similar to SQL's FROM but much more flexible as it isn't limited to tables), Return describes the shape of the query output (similar to SQL's SELECT) and WHERE sets some conditions on the attributes of the nodes.

Example: search for all nodes who have a direct or indirect relation "Knows" to node of type Crew with attribute name "Neo"; return the base node, the relation chain and the end node:



This is a more general query structure with additional features such as aggregations, skip and limit:

- MATCH (user)-[:FRIEND]-(friend)
- WITH user, count(friend) AS friends
- ORDER BY friends DESC
- SKIP 1 LIMIT 3
- RETURN user

There are many possible pattern to be searched for, here's a list with the most popular ones:

<code>(n:Person)</code>	Node with Person label.
<code>(n:Person:Swedish)</code>	Node with both Person and Swedish labels.
<code>(n:Person {name: \$value})</code>	Node with the declared properties.
<code>()-[r {name: \$value}]-()</code>	Matches relationships with the declared properties.
<code>(n)--&gt;(m)</code>	Relationship from n to m.
<code>(n)--(m)</code>	Relationship in any direction between n and m.
<code>(n:Person)--&gt;(m)</code>	Node n labeled Person with relationship to m.
<code>(m)-[:KNOWS]-(n)</code>	Relationship of type KNOWS from n to m.
<code>(n)-[:KNOWS :LOVES]-&gt;(m)</code>	Relationship of type KNOWS or of type LOVES from n to m.
<code>(n)-[r]-&gt;(m)</code>	Bind the relationship to variable r.
<code>(n)-[*1..5]-&gt;(m)</code>	Variable length path from 1 to 5 rels. from n to m.
<code>(n)-[*]-&gt;(m)</code>	Variable length path of any number of rels. from n to m
<code>(n)-[:KNOWS]-&gt;(m {property: \$value})</code>	A relationship of type KNOWS from a node n to a node m with the declared property.

One additional feature is the possibility of adding custom pattern written in Java to Cypher by loading the JAR.



## 3.3 Key-Value Stores

### 3.3.1 Introduction

In many applications performance is an essential priority, and often a small delay in response time could mean a big financial loss.

For these kind of needs key-value stores were built, a type of db which guarantees the best possible lookup time of any storage system.

Key-value stores are built upon the assumption that any entity can be seen as a value pointed by a key, and so when you are looking for a certain value you need to search for the key associated to it.

In these contexts keys are just used to retrieve a specific value, unlike relational databases where primary and foreign keys are means to do searches and mapping between tables; this obviously grants a significant speed boost during reads.

Key-Value stores can be used to improve performance at many different levels:

- Database
- Caching Layers
- Message Brokers

### 3.3.2 Redis

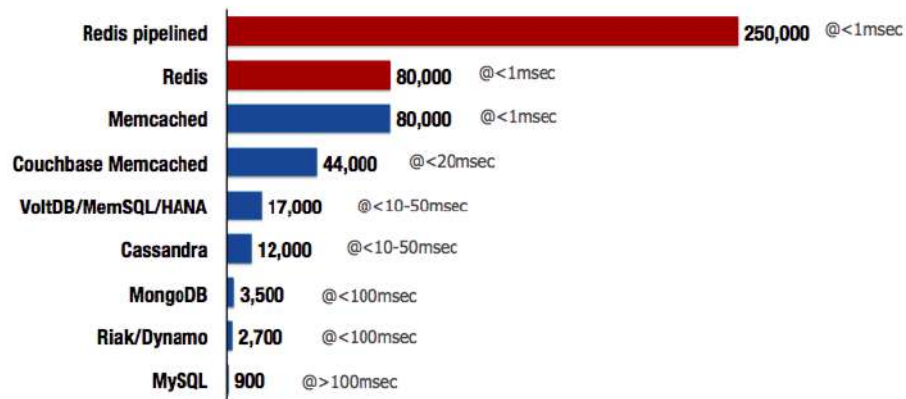
Redis is the most popular Key-Value store. It was built by a single developer, Salvatore Sanfilippo, its performance and ease of use were made possible by simplifying its architecture as much as possible:

- Written in C
- Single Threaded
- Atomic operations only
- No official support for Windows
- Executes most commands with constant  $O(1)$  complexity
- Based on 2 fundamental commands: SET(k,v), which sets the value v to the key, k and GET(k), which returns the value associated to the key k
- In Memory in DB

These characteristics make Redis very fast on reads, but not very efficient for general purpose long term storage system (since being in-memory means that a single shutdown will cost you all your data).

Redis is the best technology for the use cases where:

- Speed is critical
- You need fast access to complex data structures
- Dataset can fit in memory
- Data is not critical



Other than the basic GET and SET commands there are a few other utility functions, but they only build upon the 2 original ones. There are commands to delete and check existence of a key and get the type of a value, for example.

Another important feature is the ability to set an expiration time to a key, after which the key is deleted from the db. Since Redis is often used as a caching layer, this feature is very important.

Upon the basic dictionary-like behaviour of Redis richer and more complex features are built.

One of these it's the possibility of handling the values in the key as more complex and abstract data types that can be associated to a specific key-value pair, such as:

- String
- Hash
- List
- Set
- Sorted Set
- Geospatial Set
- HyperLog

This way we can retrieve values by searching its keys and once we have those values treat them as more complex data structures.

## Redis Commands - Lists & Hashes

Lists	
<b>Push on either end</b> <code>redis&gt; RPUSH jobs "foo"</code> <code>(integer) 1</code> <code>redis&gt; LPUSH jobs "bar"</code> <code>(integer) 1</code>	RPUSH/LPUSH [key value] $O(1)$
<b>Pop from either end</b> <code>redis&gt; RPOP jobs</code> <code>"foo"</code> <code>redis&gt; LPOP jobs</code> <code>"bar"</code>	RPOP/LPOP [key] $O(1)$
<b>Blocking Pop</b> <code>redis&gt; BRPOP jobs</code> <code>redis&gt; BRPOP jobs</code>	BRPOP/BLPOP [key] $O(1)$
<b>Pop and Push to another list</b> <code>redis&gt; RPOPLPUSH src dst</code>	RPOPLPUSH [src dst] $O(1)$
<b>Get an element by index</b> <code>redis&gt; LINDEX jobs 1</code> <code>"foo"</code>	LINDEX [key index] $O(N)$
<b>Get a range of elements</b> <code>redis&gt; LRANGE jobs 0 -1</code> <code>1. "bar"</code> <code>2. "foo"</code>	LRANGE [key start stop] $O(N)$

Hashes	
<b>Set a hashed value</b> <code>redis&gt; HSET user:1 name John</code> <code>(integer) 1</code>	HSET [key field value] $O(1)$
<b>Set multiple fields</b> <code>redis&gt; HMSET user:1 lastname Smith visits 1</code> <code>OK</code>	HMSET [key field value ...] $O(1)$
<b>Get a hashed value</b> <code>redis&gt; HGET user:1 name</code> <code>"John"</code>	HGET [key field] $O(1)$
<b>Get all the values in a hash</b> <code>redis&gt; HGETALL user:1</code> <code>1) "name"</code> <code>2) "John"</code> <code>3) "lastname"</code> <code>4) "Smith"</code> <code>5) "visits" 6) "1"</code>	HGETALL [key] $O(N)$ : $N$ =size of hash.
<b>Increment a hashed value</b> <code>redis&gt; HINCRBY user:1 visits 1</code> <code>(integer) 2</code>	HINCRBY [key field incr] $O(1)$

## Redis Commands - Sets & Sorted sets

Sets	Sorted sets
<b>Add member to a set</b> <code>redis&gt; SADD admins "Peter"</code> <code>(integer) 1</code> <code>redis&gt; SADD users "John" "Peter"</code> <code>(integer) 2</code>	<b>Add member to a sorted set</b> <code>redis&gt; ZADD scores 100 "John"</code> <code>(integer) 1</code> <code>redis&gt; ZADD scores 50 "Peter" 200 "Charles" 1000 "Mary"</code> <code>(integer) 3</code>
<b>Pop a random element</b> <code>redis&gt; SPOP users</code> <code>"John"</code>	<b>Get the rank of a member</b> <code>redis&gt; ZRANK scores "Mary" (integer) 3</code>
<b>Get all elements</b> <code>redis&gt; SMEMBERS users</code> <code>1) "Peter"</code> <code>2) "John"</code>	<b>Get elements by score range</b> <code>redis&gt; ZRANGEBYSCORE scores 200 +inf WITHSCORES</code> <code>1) "Charles" 2) 200</code> <code>3) "Mary" 4) 1000</code>
<b>Union multiple sets</b> <code>redis&gt; SUNION users admins</code> <code>1) "Peter"</code> <code>2) "John"</code>	<b>Increment score of member</b> <code>redis&gt; ZINCRBY scores 10 "Mary" "1010"</code>
<b>Diff. multiple sets</b> <code>redis&gt; SDIFF users admins</code> <code>1) "John"</code>	<b>Remove range by score</b> <code>redis&gt; ZREMRANGEBYSCORE scores 0 100</code> <code>(integer) 2</code>

Let's now talk about some of the ways to make Redis more scalable and fault tolerant, in order to be able to integrate it safely in Big Data applications:

- **Persistence:** while all main operations are done in main memory, Redis provides a couple of features to implement disk persistence. The first one allows you to take a snapshot of the db and save it to a file (RDB), while the second one works by keeping track of changes in append only log file (AOF)
- A Redis instance known as the master, ensures that one or more instances known as the slaves, become exact copies of the master. Clients can connect to the master or to the slaves. Slaves are read only by default.
- **Partitioning:** Breaking up data and distributing it across different hosts in a cluster; can be implemented in different layers:  
 Client: Partitioning on client-side code.  
 Proxy: An extra layer that proxies all redis queries and performs partitioning  
 Query Router: instances will make sure to forward the query to the right node. (i.e Redis Cluster).
- **Failover:** it can be either Manual, Automatic with Redis Sentinel (for master-slave topology) or Automatic with Redis Cluster (for cluster topology)

## 3.4 Columnar Databases

### 3.4.1 Introduction

In the recent years there has been a ever growing need for technologies capable of handling large scala data analysis Thas need was born because dataset of "Big Data" size have often very different characteristics than the ones usually stored in realtional databases:

- Data Large and unstructured
- Lots of random reads and writes
- Less use of foreing keys in favour of denormalized data structures

What this new kind of data needs instead is:

- Incremental scalability
- Speed
- No single point of failure
- Low cost and administration
- Horizontal scalability

To address this new needs a new set of databases was invented:columnnar databases,who owe their name to the fact that they store data by column,as opposed to traditional dbs that store data by row.

Row based databases work well in an OLTP application because it is easy to insert and update records and there isn't much need for bulk reads;however in OLAP environments where analytical jobs often read a huge number of records at once it's fundamental to be able to load only the necessary columns,otherwise the overhead of having to load the entire rows in main memory would be unsustainable.

Of course these aren't the only differences between the two approaches,here are some more pros and cons of choosing a columnar database:

Pros:

- Data compression
- Improved bandwidth utilization
- Improved code pipelining
- Improved cache locality

Cons:

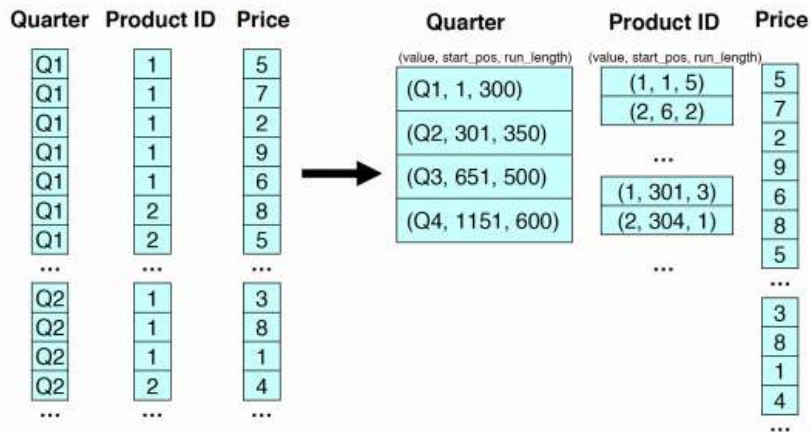
- Increased disk seek times
- Increased cost of Inserts

- Increased tuple reconstruction costs

What we can gather from this is that columnar dbs are suitable for read-mostly,read-intensive,large data repositories like the ones used by OLAP systems.

Let's now focus on one of the biggest advantages of columnar databases,which is lossless data compression: since it happens very often that a column is composed of a few different values repeated many times, what we can do is to group all the equal values together and instead of storing that same values many times just store that datum once alongside the number of times that value is repeated. For example if we have one column called "Quarter" who only has the values "Q1","Q2","Q3" and "Q4" instead of repeating the full string for every record we can just sort the db by that column and save tuples of this sort ("Q1",200),meaning that the value "Q1" appears 200 times (this approach is called Run-Length Encoding).

## Run-Length Encoding



This approach exploits the so called spacial-locality,which means storing similar data values next to each other.

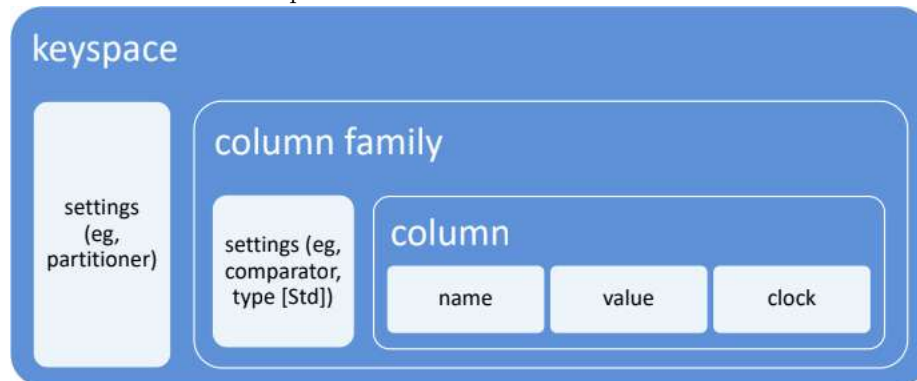
### 3.4.2 Cassandra

Cassandra is a columnar database with some key-value features. It was originally designed at Facebook and is now maintained by the Apache Foundation as an open-source project.

Its data model is based around the concept of column family, which is a group of columns. Column families are organized by rows identified by a unique key (hence the key value approach). While this approach may seem similar to relational tables at first glance, there are many significant differences that make it a NoSQL db:

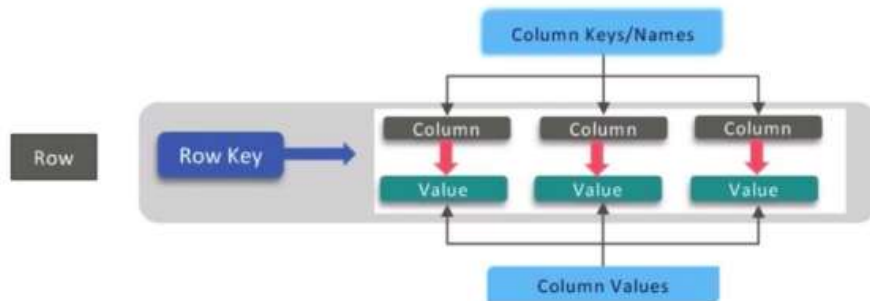
- Each cell in a table could be missing or different from the others
- Schemaless
- Increased tuple reconstruction costs
- Support for key getters and setters

Let's now take a deeper look at Cassandra's data model:



Each column is composed by a name, a value and a clock (which is the timestamp of the last update on that column). The value of a column can contain other columns inside him and in that case that column is called a super-column. Columns are grouped in families, and many group families make up a keyspace (which is the equivalent of a database in SQL)

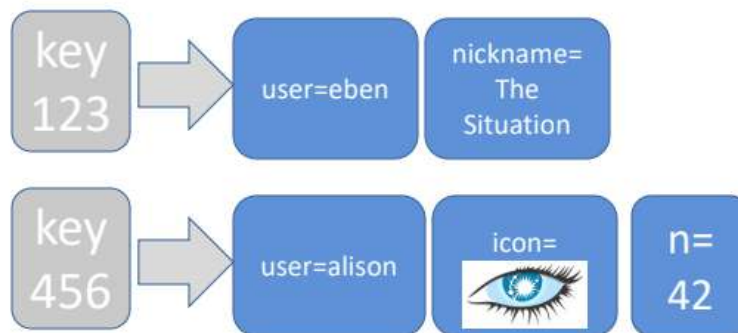
We can see Cassandra also as a key-value store, where the key is the row identifier and the values are the schemaless column families.



So a keyspace is a group of column families sharing the same application configurations, and the column family is a group of column with values similar or related to each other.

It's important to specify similar and not equal values because, unlike relational dbs, Cassandra can be seen as row-oriented, where every row is made up of many columns with values of similar kind but not necessarily equal to the ones in the other rows of the same cf.

Example of 2 rows in a User Column Family:



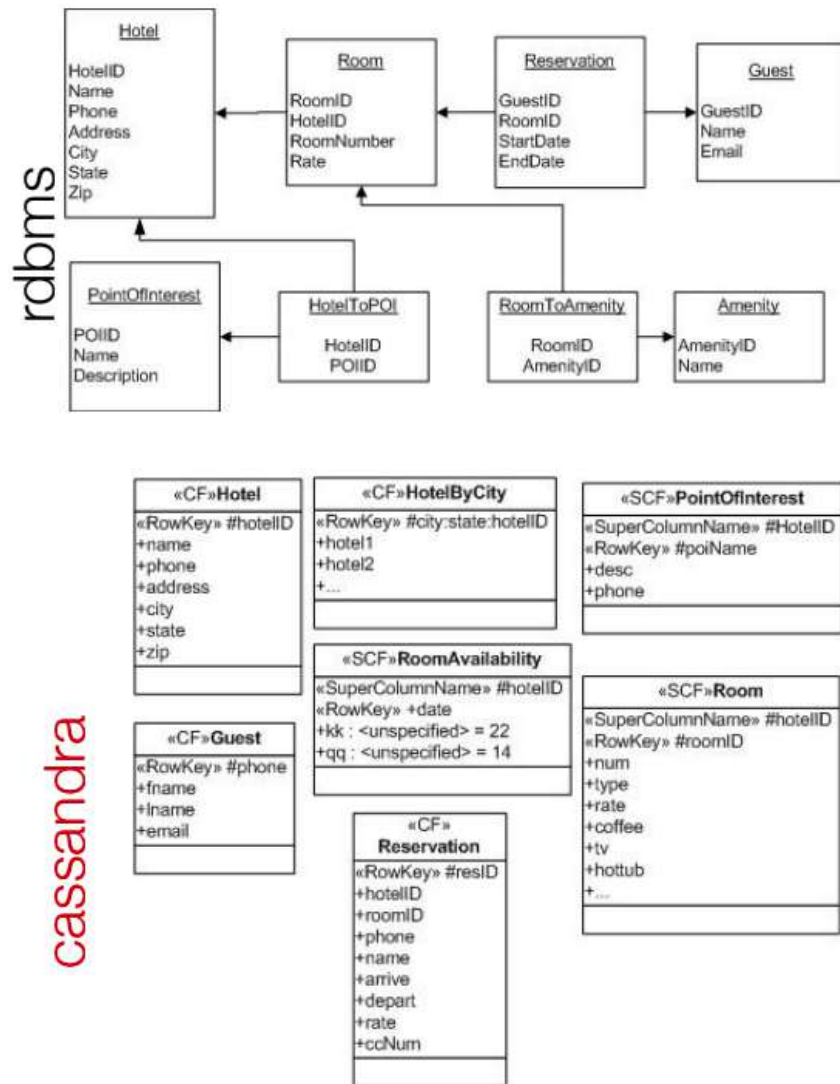
This flexible structure of the columns implies what is maybe the most peculiar characteristic of Cassandra, which is that you usually can't write queries on column values.

Cassandra in fact performs a complete paradigm shift from relational databases: while in traditional databases you design the db schema first and then you write queries over it (Domain-based-approach), in Cassandra it's used a Query First Approach, where designers first think about the requests the system will need to be able to answer to, and then they build the db structure around them.

The Domain-based and Query-based philosophies can be seen as the corresponding implementations of the schema-on-write and schema-on-read approaches we have seen in chapter 1



This is what a rdmbms schema looks like compared to a Cassandra database:



Cassandra can be seen as an "index factory", meaning that for each query there is a table optimized to handle it. For example if we have a User that we need to search by id or by city, we'll have to add a User cf indexed by id and another User cf with the same data but indexed by city, so that each query runs on an optimized table.

So the UserByCity cf would be a list of rows with the city as a key and a column for each id of the users living in that city. This way a row single row can have

a huge number of different columns and that is why columnar db are also often called Big Column stores.

If we need to search an Entity by more than 1 property at the same time then we'll need to create a cf with as key the aggregate value of all those properties.

```
<<cf>>USER
Key: UserID
Cols: username, email, birth date, city, state

How to support this query?

SELECT * FROM User WHERE city = 'Scottsdale'

• Use an aggregate key
state:city: { user1, user2}

Create a new CF called UserCity:

• Get rows between A2: & A2:
for all Arizona users

<<cf>>USERCITY
Key: city
Cols: IDs of the users in that city.
Also uses the Valueless Column pattern

• Get rows between A2:Scottsdale & A2:Scottsdale1
for all Scottsdale users
```

This big change of perspective while gives great advantages on reads implies data duplication, and that is one of the reasons why Cassandra is almost exclusively used in OLAP systems with multiple Terabytes sized data.

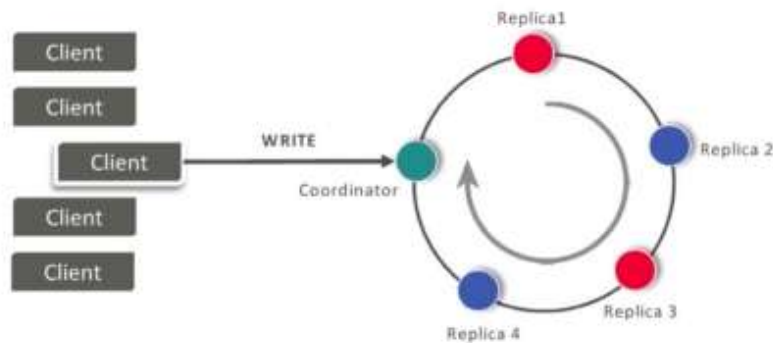
### 3.4.3 Cassandra's Architecture

One Cassandra's main difference from RDBMS it's that it relies on data duplication and denormalization and following it's query-first philosophy it doesn't support JOINS (instead of joining 2 different tables you just build an unique column family).

Here's a list of it's main characteristics:

- Tuneably consistent (we can't have full consistency as stated by the CAP theorem)
- Fast writes
- Highly available
- Fault tolerant
- Linear,elastic scalability
- Decentralized
- around 12 client languages
- O(1) dht

Cassandra is a distributed system implemented with a Ring Architecture: this means that the cluster is masterless and made up of equivalent nodes, and each of them as a predecessor and a successor. Each node can only send data to its successor and receive data from its predecessor. When a client wants to write data to the database any one of the nodes can answer to that request and if it does it becomes the coordinator for that operation, and after it has finished writing all data to disk he passes the data to its successor so that it can be replicated in other nodes.



Cassandra was designed to enable very fast writes, in order to do so it had to abandon many transactional guarantees. In order to avoid locks on resources a WORM (Write Once Read Many) methodology, where data is written only one time and then never updated again.

The write process works in an asynchronous way: the client sends the data to the coordinator which sends it to all the replica nodes responsible for that key via partitioning function. This method grants atomicity for a given key, since in the case of a fault in a replica during the write, the coordinator can just keep sending the data to the working nodes which in turn will send the data again to the broken replica once it's up and running again.

If all replicas are down, the coordinator can buffer the data stream for up to an hour.

To speed up even more the write process all records are initially stored in main memory and then flushed to disk later on (in case of a fault the node can restart interrupted operations since every action is written to a changelog).

When dealing with reads and deletes we need to handle the problem of consistency (since we have many replicas of the data stored in different nodes there's always the risk of having inconsistent data between them).

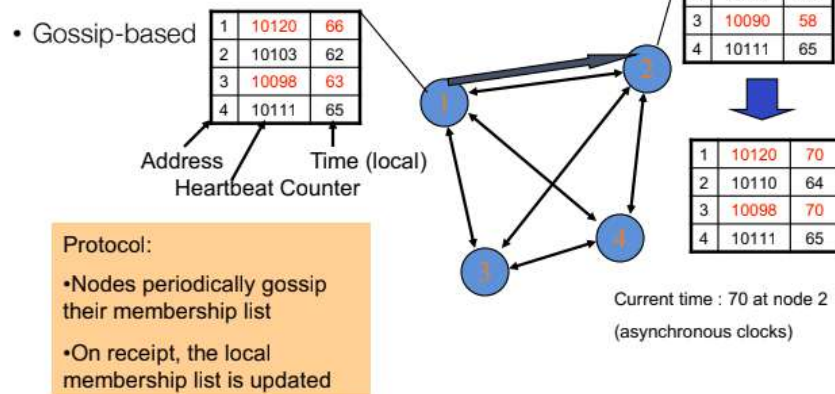
A Delete operation works by adding a tombstone to the item in the changelog and later on during the compaction that item will be effectively deleted from all replicas. The Read operation is similar to the Write, but with some differences: the coordinator fetches the same data from multiple replicas and if there is any inconsistency the right data is determined by a quorum mechanism and then a read repair operation is performed. This makes the reads a bit slower than the writes but they're still very fast

As we mentioned before consistency in Cassandra can be tuned in order to maximize or minimize performance; this is done by changing the quorum policy

- ANY: any node (may not be replica)
- ONE: at least one replica
- QUORUM: quorum across all replicas in all datacenters
- LOCAL\_QUORUM: in coordinator's DC
- EACH\_QUORUM: quorum in every DC
- ALL: all replicas all DCs

Heartbeat communication between nodes is performed by using the gossip protocol: since if each node were to ping every other node there would be an huge and unnecessary stress on the network, each machine chooses a small number of its peers to communicate with and in turn each of this peers will communicate with an equal number of different nodes. This way status updates can spread quickly in the cluster while keeping communications limited.

## Cluster Membership



Every node keeps a membership list of the others machines in the cluster alongside the timestamp of the last heartbeat message that was received from them,after a certain period of time without has passed receiving any message from one of its peers that node is considered down.

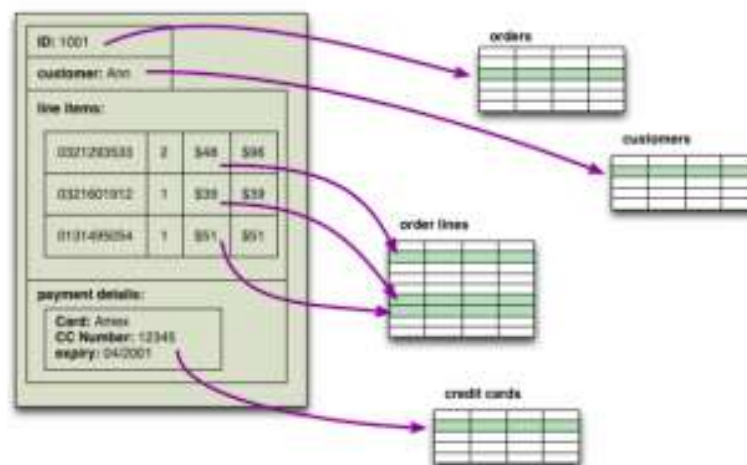
## 3.5 Document Databases

### 3.5.1 Introduction

Document databases are probably the most widespread kind of NoSQL technologies. They are called in this way because they deviate from the entity-based denormalized data model of relational dbs and prefer an approach based on denormalized and aggregated data typical of business document.

A document is a complex object made up of what would be many different SQL tables, and so it can be composed of different pieces of data with different structures.

Some of the advantages of document stores are:



- Flexible structure that handles well schema changes
- Solves mismatches impedance problems
- full JSON compatibility, which guarantees very ease of integration with web technologies

### 3.5.2 MongoDB

MongoDB is the most popular document database, it's built using a CP architecture.

Its data model is based on documents structured just like a JSON file, many documents form a collection. Every document must have a unique id and can have any number of nested documents inside him.

**Embedded documents:**



This data model is convenient for many applications (especially web based ones) since unlike relational databases you don't have to reconstruct the complex business objects from the normalized tables with expensive joins or aggregations techniques.

This approach gives you the possibility of structuring your data using the granularity which fits best your application needs.

It's however possible to establish references between documents, but it doesn't make much sense to overuse this feature since that would be just like recreating a relational database. Queries and CRUD operations between documents are made using built-in functions that can filter and operate on the objects.

## Read – mapping to SQL

SQL Statement	MongoDB commands
SELECT * FROM table	db.collection.find()
SELECT * FROM table WHERE artist = 'Nirvana'	db.collection.find({Artist:"Nirvana"})
SELECT* FROM table ORDER BY Title	db.collection.find().sort(Title:1)
DISTINCT	.distinct()
GROUP BY	.group()
>=, <	\$gte, \$lt

## Comparison Operators

Name	Description
\$eq	Matches value that are equal to a specified value
\$gt, \$gte	Matches values that are greater than (or equal to a specified value
\$lt, \$lte	Matches values less than or ( equal to ) a specified value
\$ne	Matches values that are not equal to a specified value
\$in	Matches any of the values specified in an array
\$nin	Matches none of the values specified in an array
\$or	Joins query clauses with a logical OR returns all
\$and	Join query clauses with a logical AND
\$not	Inverts the effect of a query expression
\$nor	Join query clauses with a logical NOR
\$exists	Matches documents that have a specified field



## 4 Streaming Data Engineering

### 4.1 Introduction

When a Data Engineer wants to design a scalable system, he needs to consider the dimensions by which processing systems are usually described by: Throughput, Latency and Message Size.

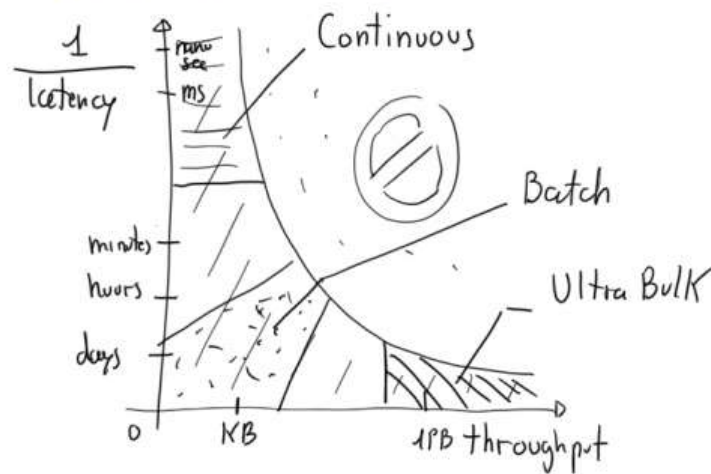
To understand these properties, let's make a comparison with the trucks of a transportation company: the time need for a truck to go through a certain road would be the latency (which means how much does it take to our server to answer to a request) and the amount of tracks that can travel at the same time would be the throughput (the amount of data our system can process in a given time frame).

Let's also say that before the company can send out a new track it needs to wait some time, the amount of which depends by how the speed and the load of the last sent truck (and in our analogy this would be the message size).

If we want to increase our performances there are many ways to do so: we can increase the speed of the trucks (vertical scaling), add lanes to the road (horizontal scaling) or decrease the time frame between trucks (reduce message size).

It's easy to understand that there is a tradeoff between throughput and latency: after a certain point, if you want to increase one you need to decrease the other.

### The throughput / latency trade-off



The graph above represents the solution space created by this tradeoff, the area below the curve represent the possible characteristic of a computing systems, the one above it's impossible to obtain.

The solution space can be divided into 3 main zones:

- The area with very high throughput but bad latency is called Ultra-Bulk
- The sector with extreme low latency obtained by sacrificing throughput is the Continuous zone
- The middle ground is called Batch

An example of Ultra Bulk is when companies decide to move their in-premise data to the cloud not by using the internet but by transporting the hard drives by car directly to their cloud provider's data centers. This solution may be counterintuitive but when the data size is very big its the fastest way to perform the migration (for example if the company had to move 1 PB of data it would take 37 days even with a 2.5 Gbps network,while with the transfer appliance the migration time doesn't depend on the size of the data but only on the distance from the datacenter).

Let's make now an example of a continuous case instead:let's say that a multinational operating all over the world has both retail stores and e-commerce websites,and this company wants to keep track of all sales and purchases in order to create a real time inventory

The sales could generate something like 200 MB/s,which isn't much for a network to handle but it could be difficult for the hard disks to have this kind of throughput if we have a single node architecture.In fact,in order to solve this problem we would need a continuous distributed system.

Let's now cover the last and more common case,the batch one: a company needs a data warehouse to host its analytical data coming from many sources (e-commerce,catalog,retail shops...);the data get periodically extracted from the source databases by using the data change capture technique (which means basically by reading the commits from the log file) and saved to a could blob storage in a raw text format (like CSV,JSON or XML),then a listener program loads all the data uploaded to the blob storage to the data warehouse.

The refresh period it's usually quite high and can vary from once every few minutes to even once a day.

Lastly Business Intelligence tools connect to the data warehouse and generate a report.

## 4.2 Streaming

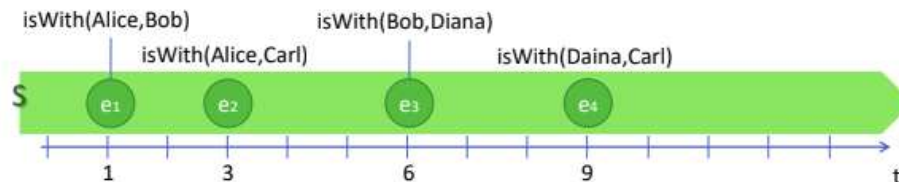
We will now dive deeper into the continuous case, also called Data Streaming. This approach performs a complete paradigm shift from traditional data processing: if in traditional systems data gets stored first and then processed later on, in streaming data gets transformed and analyzed continuously as it is generated.

Before going forward we need to talk about the Time Models which are the many possible ways in which the relation between the passing of time and information being generated can be interpreted.

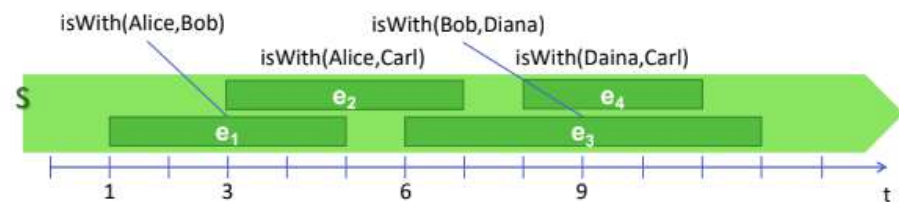
- Causal Time Model: the events happen one after the other, it's possible to determine an order of the events but we can't know how much time has passed in between but we can still make meaningful queries (e.g. Does Alice meet Bob before Carl? Who does Carl meet first?)



- Absolute Time Model: like causal time model, but the system keeps track of the time passed between each event, and this gives additional expressive power to our queries (e.g. How many people has Alice met in the last 5 minutes? Does Diana meet Bob and then Carl within 5 minutes?)



- Interval-based Time Model: like the absolute time model, but we also know how long an event lasts, and this gives us the chance to write even more queries (e.g. What meetings last less than 5 minutes? Which meetings are overlapping?)



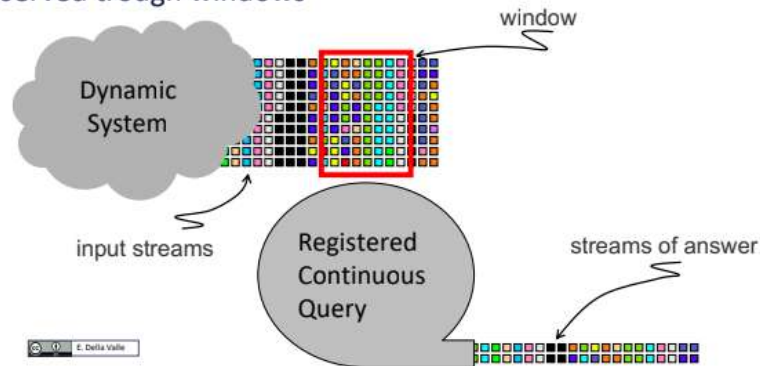
Different communities in the computer science world have implemented the streaming concept in different ways, let's take a look at them:

- Database Community: Standard DBMSs are passive, meaning that the storage system is static and the queries performed by the users are active and dynamic. Streaming systems work in an active way instead: the queries are always the same and they run on a continuous stream of ever-changing data, the so-called Data Streams, which are unbounded sequences of time-varying data.

This concept is implemented in practice by having analytical jobs working on limited windows of the data stream, as it would be impossible to operate on an unbounded source.

The database community uses heavily the time models to implement their streaming data systems.

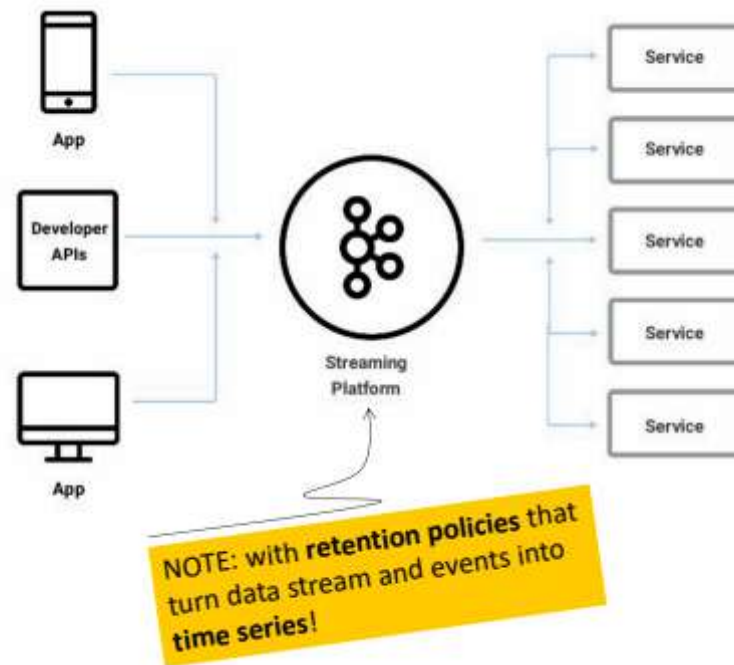
- Continuous queries registered over streams that are observed through windows



- Event-based Community: this model is based on event processing systems, where each component publishes data about a certain topic and subscribes to the topics he's interested in via a middleware. This kind of communication is asynchronous, anonymous, message based, multicast and implicit.

The Complex Event Processing (CEP) variant adds the ability to deploy rules that describe how composite events can be generated from primitive (or composite) ones

- Service Oriented Community: in this architecture component can communicate only by calling services, so there's a standardized protocol of communication in the system and no component can directly manipulate another one. This approach however is rigid and creates tight coupling between services, so a new version called Event-Driven Architecture (EDA) was developed in order to have a more flexible and extensible model. In EDA components communicate with services via data streams that act as a middleware, this way the messages don't get lost if the service goes down



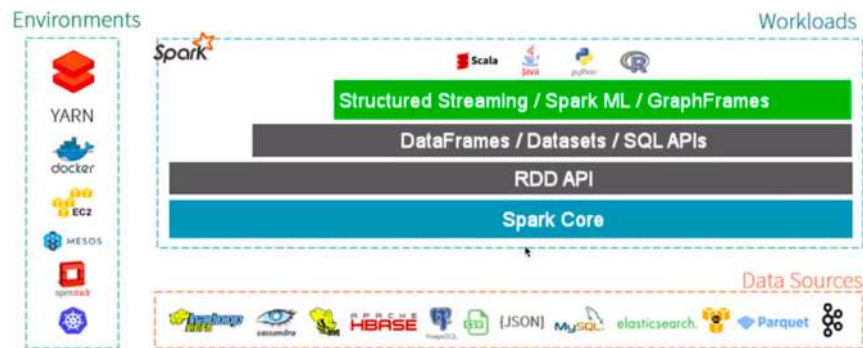
## 4.6 Spark

### 4.6.1 Introduction

Apache Spark is a unified analytics engine for big data processing, with built-in modules for streaming, SQL, machine learning and graph processing

It's similar to Hadoop but it has the purpose of being many times faster, and it achieves that goal by performing all computation in memory rather than on disk. While in Hadoop each partial output of a data processing job is saved in disk, Spark keeps everything it can in RAM and only stores in disk the final output of its jobs.

Spark has a layered architecture in which each layer has APIs that are used by the upper one. The higher the layer, the more abstract are the data structures. Spark has connectors for almost any existing data storage solution and it can run on a variety of environments (Yarn, Docker, Kubernetes, Mesos...).



These are its main data structures, from the most barebone to the most abstract:

- **RDD:** a distributed collection of JVM objects with functional operators
- **Dataframe:** relational db like collection of rows that support expression based operators and UDFs. Operations are very efficient thanks to logical plans and optimizers, that take advantage of its well defined internal representations.
- **Dataset:** abstraction that allows you to treat Dataframes like JVM objects. Operations aren't as fast as on pure Dataframes but they are still quite fast, typesafe and integrate well into traditional OOP programming. However they aren't as popular as DFs because they are not as convenient for data analysis with is what Spark is mainly used for.

## 4.6.2 Spark APIs

Let's go through each Spark API in detail:

RDD stands for: Resilient (Fault Tolerant), Distributed (split across many nodes and computed in parallel), Dataset. RDDs are immutable, that means that once they are initialized they can't change in any way. This allows to take advantage of many caching and locality strategies that greatly improve performance.

Since RDDs are distributed it's necessary to keep track of how each partition it's connected to the others and where each piece of data resides and which job generated it.

Here's an example of a wordCount job with RDDs written in Python:

```
# Open textFile for Spark Context RDD
text_file = sc.textFile("hdfs://...")
# Declare word count
res = (text_file.flatMap(lambda line: line.split())
      .map(lambda word: (word, 1))
      .reduceByKey(lambda a, b: a+b))

# Execute word count
res.collect()
```

On a RDDs it's possible to execute transformations (such as map, filter and group) or actions (such as count, collect and save). Transformations are lazily evaluated which means that they won't actually be executed even after their declaration until they are required by an action (which are eager instead).

Lazy execution has many advantages:

- Not forced to load all data at 1st step, technically impossible with large datasets
- Easier to parallelize operations: N different transformations can be processed on a single data element, on a single thread, on a single machine
- Allows optimizations

It's important to note that since RDDs are immutable each transformation doesn't modify the original but creates a new object.

Transformations can be either narrow if they can be executed within a single partition, or wide if they require data distribution between nodes (this operation is called shuffling). A Job is orchestrated by a Driver node and all the processing work on the partitions is done by Worker nodes.

Actions are operations that either return a value or write something to disk. Another huge contributor to Spark's speed is caching, which basically works by keeping in main memory RDD blocks that can be reused in the same job instead of loading them from disk each time they're needed. Example:

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
```

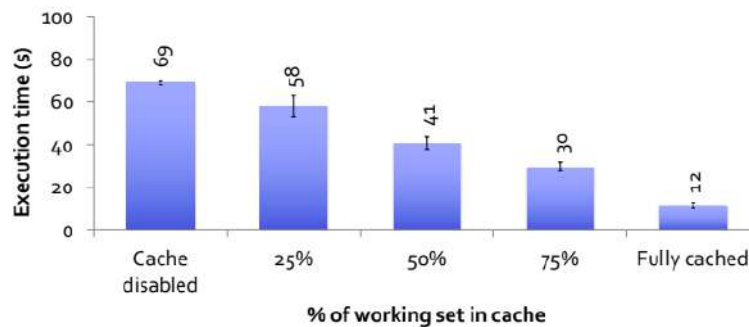
```
messages = errors.map(lambda s: s.split("\\t")[2])
```

```
//with this action the message RDD is kept in memory and reused in the next 2 lines
messages.cache()
```

```
messages.filter(lambda s: \"mysql\" in s).count()
messages.filter(lambda s: \"php\" in s).count()
```

Thanks to all this optimizations with Spark is now possible to make with few resources and in a small time tasks that would have taken huge amounts of time and money just a few years ago (like for example a full text search on Wikipedia).

## Caching and performance

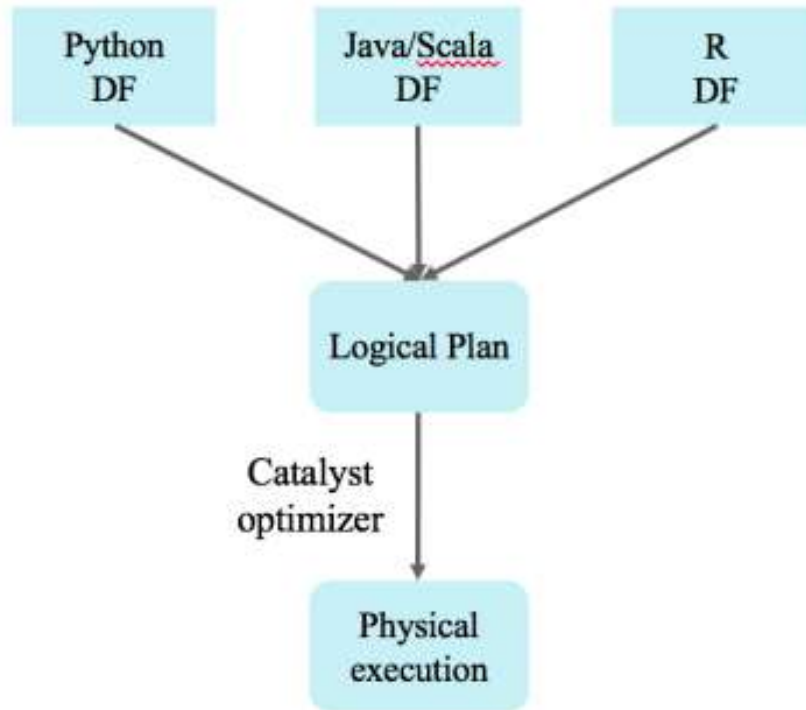


While RDDs are the building blocks of Spark, they are relatively low level abstractions and the APIs that Data Engineers usually use are the Dataframes. Dataframes are a collection of immutable data with named columns built on top of RDDs. They offer a user-friendly, cross-language API and support many performance optimizations. Here's an example of Dataframe Python code which utilizes SQL to query and join data coming from different sources.

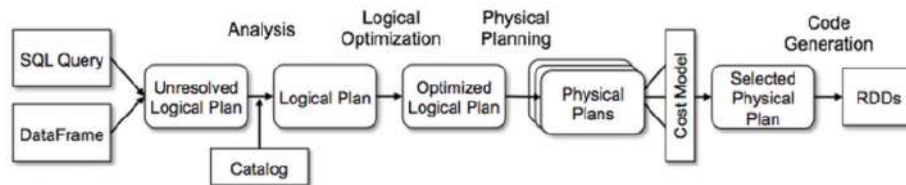
```
# Read a CSV
userDF = spark.read.csv(" ... /userData.csv")
# ... or Use DataFrame APIs and register a temp view
middleageSmokers = userDF.filter(col("smoker")== "N").filter(col("age")>40)
middleageSmokers.createOrReplaceTempView("middleageSmokers")
# ... read a CSV and register a temp view
spark.read.json(" ... /part-00000.json.gz").createOrReplaceTempView("iot_stream")
# ... execute SQL query or ...
spark.sql("SELECT avg(calories_burnt) FROM iot_stream JOIN middleageSmokers ON ...")
```



One very important feature of dataframes is that every operations has the same performance across every programming language. This is possible because the language APIs are just declarative interfaces calling a low level process that formulates the best logical plan for that operation and execute in by generating RDD code.



Here's the query manager architecture for Dataframes:

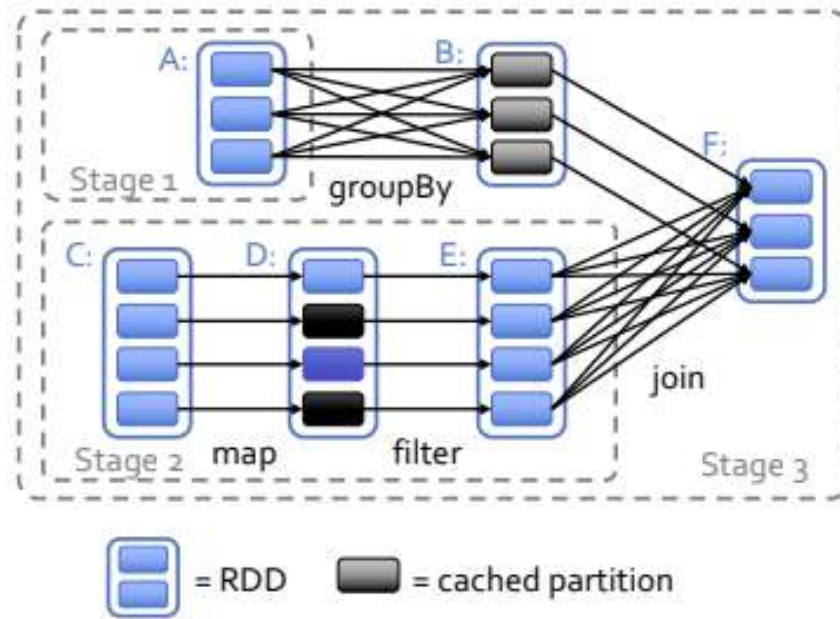


The last Spark API is the Dataset, which offers strong static typing over Dataframes. This feature (which is only useful in languages with static types like Java and Scala) allows for more compile time error checks compared to the loosely typed Dataframes.

### 4.6.3 Spark Scheduler

The DAG Scheduler is the scheduling layer of Spark, it transforms a logical execution plan into a physical one (a Job) considering data locality and partitioning (to avoid shuffles).

Each Job is a DAG of Stages which are broken down into Tasks



### 4.6.4 Spark Structured Streaming

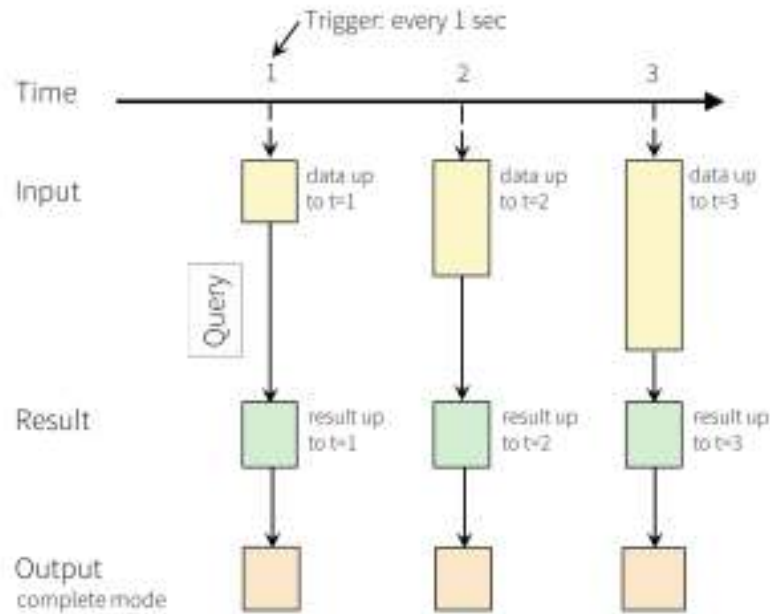
Spark Structured Streaming (SSS) provides fast, scalable, fault-tolerant, end-to-end exactly-once stream processing without the user having to reason about streaming. This means that Spark allows you to operate on streams by writing Dataframe and the developer could even forget he is working on a stream.

There are 2 possible operational modes in SSS

- Micro-batch processing engine: the default mode, it offers end-to-end latencies as low as 100 milliseconds and exactly-once fault-tolerance guarantees
- Continuous Processing: Added from Spark 2.3, it offers end-to-end latencies as low as 1 millisecond and at-least-once guarantees

The key idea behind this technology is treating a stream like a table that is being continuously appended, this way the user can write static queries over a Dataframe and Spark runs it as an incremental query over the unbounded table.

The programming model works by executing a new computation periodically, and the state calculated in the current computation is passed on to the next one so each computation works with all existing data by integrating past state and the newly received records.

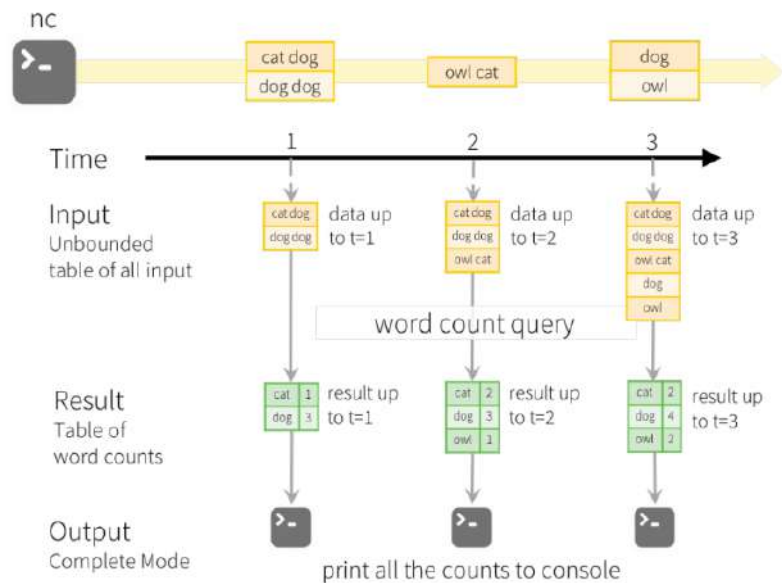


There are 3 possible output modes for stream processing jobs:

- Complete Mode: The entire updated result is written
- Append Mode: Only the new rows appended in the result are written, applicable only when existing rows in the result are not expected to change
- Update Mode: Only outputs the rows that have changed since the last trigger. If the query doesn't contain aggregations, it will be equivalent to Append mode.

Each mode is compatible only with some queries and data sinks and so the Data Engineer needs to evaluate the best mode for each use case.

Here's an example of a word count process in a streaming environment that shows how Spark computes new state from old state and current record (old records are discarded since only the state derived from them is necessary to the following computations).



### Model of the Quick Example

The SSS Dataframes API works by calling the `SparkSession.readStream()` function which returns a stream based Dataframe from a specified source, which could be a file (supported formats include CSV, JSON, ORC, Parquet), a Kafka Topic or even a Socket (even though this way you don't get Fault Tolerance).

Python Examples:

```
spark = SparkSession. ...
# Read text from socket
socketDF = spark
    .readStream
    .format("socket")
    .option("host", "localhost")
    .option("port", 9999)
    .load()
socketDF.isStreaming() # Returns True for DataFrames that have streaming sources
socketDF.printSchema()
```

```

# Read all the csv files written atomically in a directory
userSchema = StructType().add("name", "string").add("age", "integer")
csvDF = spark
    .readStream
    .option("sep", ";")
    .schema(userSchema)
    .csv("/path/to/directory")

# Input: streaming DataFrame with IOT device data with schema
{ device: string, deviceType: string, signal: double, time: DateType }

df = ...
# Select the devices which have signal more than 10
df.select("device").where("signal > 10")
# Running count of the number of updates for each device type
df.groupBy("deviceType").count()

# Alternatively, register a streaming DataFrame/Dataset as a temporary
view and apply SQL on it

df.createOrReplaceTempView("updates")
spark.sql("select count(*) from updates")

```

Windows operations are implemented through an additional timestamp column, and are viewed as grouping operations based on that column.  
Example

```

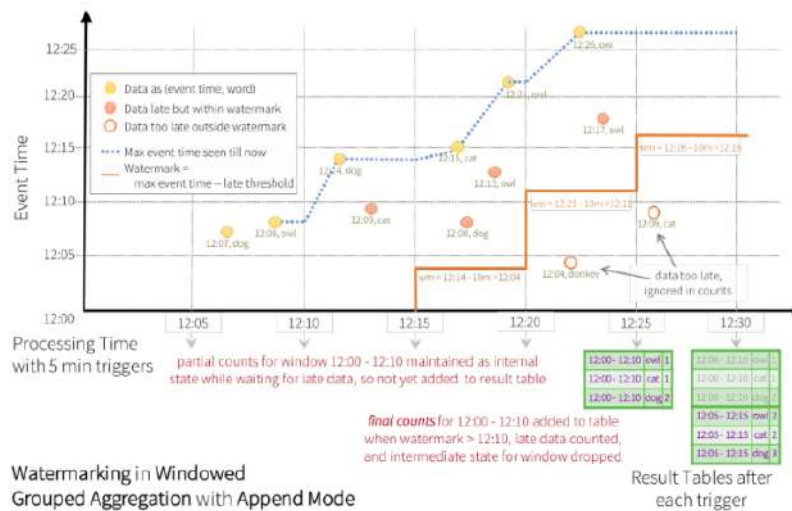
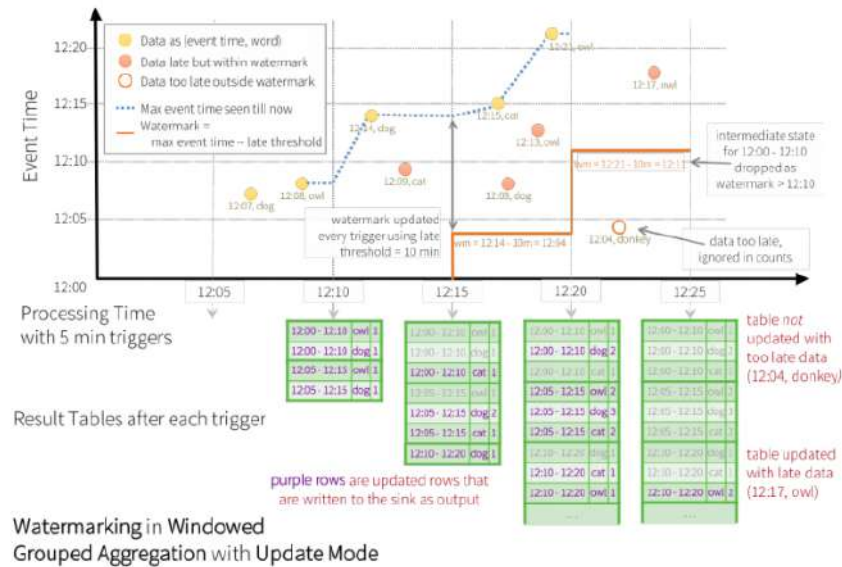
words = ... # streaming DataFrame of schema { timestamp: Timestamp, word: String }
# Group the data by window and word and compute the count of each group
windowedCounts = words.groupBy(
    window(words.timestamp, "10 minutes", "5 minutes"),
    words.word
).count()

```

An important topic in SSS is how late records (which are record with a different generation time and ingestion time) are managed.

The solution is to use watermarking, which means that in update mode the dataframes can be retroactively updated with late records if the delay is within a certain time limit, while in append mode the processing output is delayed for some time in order to wait for possible late entries.

This way you at least have the guarantee that data with a delay inferior to a certain time frame won't be ignored. Data with greater delays could either be ignored or considered, but the greater the delay the more likely it is that the record won't be taken into account.



Spark Structured Streaming supports both stream-to-stream and stream-to-table joins.

Of course native stream joins isn't possible, as it isn't possible to compare each row of 2 unbounded tables, so Spark performs this kind of joins by using watermarks that tell you how long should you keep on buffering before performing a static join with all the buffered data.

Example:

```
from pyspark.sql.functions import expr
impressions = spark.readStream. ...
clicks = spark.readStream. ...
# Apply watermarks on event-time columns
impressionsWithWatermark = impressions.withWatermark("impressionTime", "2 hours")
clicksWithWatermark = clicks.withWatermark("clickTime", "3 hours")
# Join with event-time constraints
impressionsWithWatermark.join(
    clicksWithWatermark,
    expr("""
        clickAdId = impressionAdId AND
        clickTime >= impressionTime AND
        clickTime <= impressionTime + interval 1 hour
        """)
)
```

Here's the join compatibility list:

- Inner: supported, optionally specify watermark on both sides + time constraints for state cleanup
- Right/Left Outer: conditionally supported, must specify watermarks on left/right + time constraints for correct results, optionally specify watermark on left for all state cleanup
- Full Outer: not supported.

The only topic left to discuss is the stream configuration. When creating a query on a stream in fact you need to specify the following:

- Output sink: Data format, location, etc.
- Output mode: what gets written
- Query name: Optionally, a unique name of the query
- Trigger interval: Optionally, the trigger interval
- Checkpoint location: for the end-to-end fault-tolerance

## 5 Data Pipelines

### 5.1 Data Ingestion

#### 5.1.1 Introduction

Data ingestion is the first and fundamental step of any Data Analysis Pipeline. The focus of this section is on how is it possible to collect data from publicly available sources over the web.

It is in fact a common practice nowadays to integrate the proprietary data coming from OLTP databases with data coming from public web resources in order to also have a data source which isn't coming from the Company domain and (and could as such be biased in many ways) . We'll now dive deep into the two main ways of collecting data from the web: APIs and Scraping.

#### 5.1.2 Web APIs

An API (Application Program Interface) is a set of routines, protocols and tools for building software applications. A Web API is a API which is based on the HTTP protocol, and it's usually used to give programmatic access to data and platforms.

The main advantage of Web APIs are:

- Separation between model and presentation
- Regulate access to the data (traceable accounts, enable access to only a portion of the available data)
- Avoid direct access to the platform website
- Request throttling
- Avoid server congestion
- Avoid DOS/DDOS
- Provide paid access to full (or higher volume) data

To access a resource via Web API, you need to specify a method and an URL:

#### URL format (RFC 2396):

scheme:[//[user:password@]host[:port]][/]path[?query][#fragment]

The diagram shows a URL: `abc://user:pass@example.com:123/path/data?key=value&k2=v2#fragid1`. Brackets below the URL identify its components: `abc` is the **scheme**; `//user:pass` are the **credentials**; `example.com` is the **domain**; `:123` is the **port**; `/path/data` is the **path**; `?key=value&k2=v2` is the **querystring**; and `#fragid1` is the **fragment**.



There are many commands (GET,POST,PUT,DELETE...) each with its own purpose. A request can be tweaked by adding parameters to end of the url;Example:

```
GET https://api.twitter.com/.../search/tweets.json
?q=%23expo2015milano
&lang=it
&result_type=recent
&count=100
&token=1235abcd
```

To give a more standardized structure to the APIs,most developers follow the REST paradigm.

The RESTful approach is:

- Client-Server: a uniform interface separates clients from servers.
- Stateless: the client-server communication is constrained by no client context being stored on the server.
- Cacheable: clients and intermediaries can cache responses.
- Layered system: clients cannot tell whether are connected directly to the server.
- Uniform interface: simplifies and decouples the architecture.
- Resource centric

A REST API in comparison to a generic one has the advantage of having a standard/predictable URL format:

```
METHOD http://domain/collection/item
GET http://example.com/user/1234
```

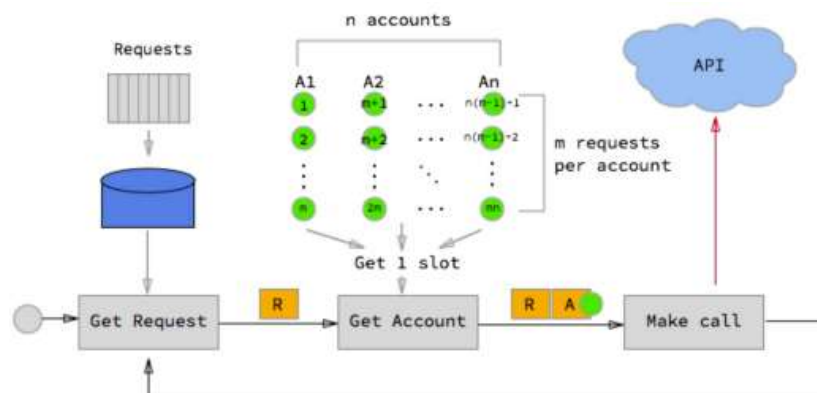
REST also standardizes the behaviour of a Method call based on the context (collection/item)

VERB	Collection	Item
POST	Create a new item.	Not used
GET	Get list of elements.	Get the selected item.
PUT	Not used	Update the selected item.
DELETE	Not used	Delete the selected item.

Almost all the APIs require a kind of user authentication.The user must register to the developer of the provider to obtain the keys to access the API. Most of the main API providers use the OAuth protocol to authenticate a user.

The user registers to the developer portal to obtain a key and a secret. Platform authenticates the user via key/secret pair and supply a token that can be

In the following example a Round Robin policy is used in order to assign requests from a stack to the nodes of a crawler:



### 5.1.3 Web Scraping

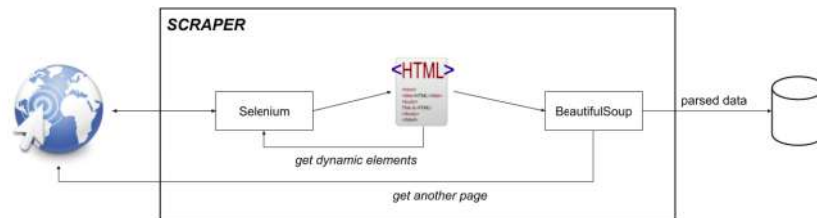
If there's no API available to get data from a website, the next best thing is web scraping.

Web scraping, web data mining, web harvesting are synonyms to one single concept: "Build an agent which can extract, parse, download and organize useful information directly from HTML source code". To get meaningful data with web scraping you need to download the HTML code of the page you want to scrape and then analyze its structure and extract the information you need.

To navigate the HTML structure we can use a query language called XPath, which combined with regular expressions allows you to extract content from web pages in a very precise way.

If we need to scrape a dynamic page the process becomes more complex, as we may need to overcome authorization forms or perform some more complicated action in order to get the content we need.

Python is a very good language for building web scraper since it offers many good libraries and tools to help you. One very popular toolset is the combination of Selenium (a headless browser with support for dynamic behaviour) and BeautifulSoup (a HTML navigator)



It's important to keep in mind that scraping should be used just as a last resort when there's no API available, since scrapers are very difficult to maintain (if the website changes structure the scraper doesn't work anymore) and many website providers don't allow scraping, so you could run into legal problems if you do.

## 5.2 Data Wrangling

### 5.2.1 The need for clean data

The step after Data Acquisition and before Analysis in the data management flow is Data Wrangling, also called Data Cleaning or Data Preparation.

In this step data is cleaned, tested and prepared to be the best input possible for the analysis process. In other words, we need to ensure that our data has high quality, which is a property defined by these metrics:

- Accuracy: The data was recorded correctly.
- Completeness: All relevant data was recorded.
- Uniqueness: Entities are recorded once.
- Timeliness: The data is kept up to date (and time consistency is granted).
- Consistency: The data agrees with itself.

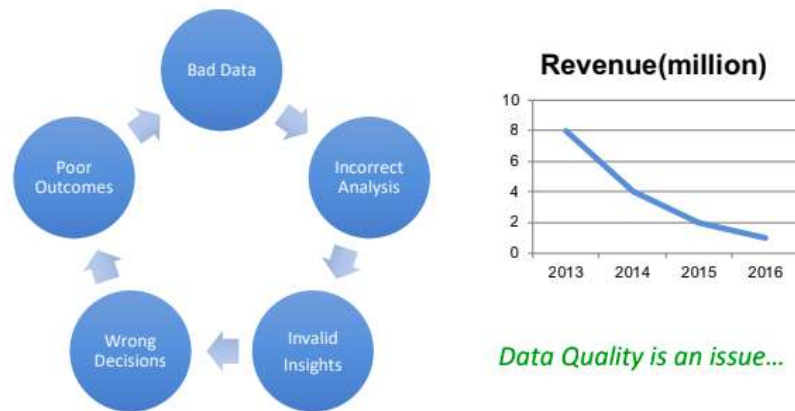
Unfortunately it isn't easy to measure and define these metrics (e.g. how do you know that you're missing data? How do you know that the data is accurate). In fact, these metrics could be:

- Unmeasurable: Accuracy and completeness are extremely difficult, perhaps impossible to measure.
- Context independent: No accounting for what is important. E.g., if you are computing aggregates, you can tolerate a lot of inaccuracy.
- Incomplete: What about interpretability, accessibility, metadata, analysis, etc.
- Vague: The conventional definitions provide no guidance towards practical improvements of the data.

For these reasons data wrangling is often the most crucial, difficult and predominant (statistics say that in average 80% of the time of a data scientist is spent cleaning data) task of a data scientist/engineer.

If the data wrangling step is done poorly, it may lead to analysis being run on bad data which could in turn lead to bad decision making in company and bad business.

Research says that this is a very common issue even in big companies, who often can't even trust their own internally generated data.

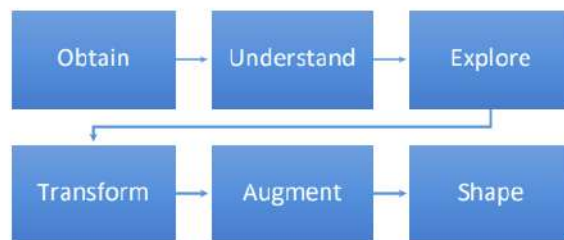


### 5.2.2 Data cleansing

So how can we transform raw unreliable data into refined high quality material? The common solution that has been spreading in the last few years is to improve the traditional ETL(Extract-Transform-Load) approach by adding more steps that test and improve the quality raw data.

#### Iterative process

- Understand
- Explore
- Transform
- Augment
- Visualize



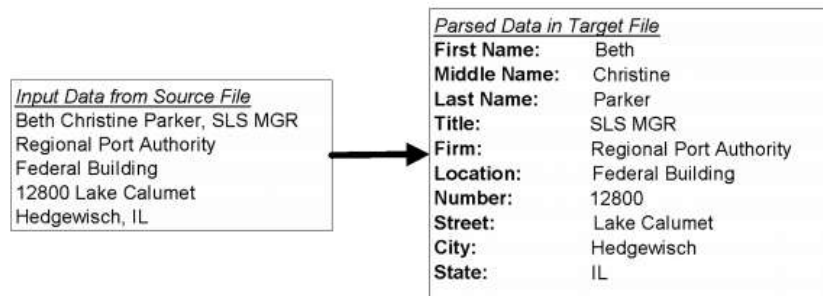
The first of this steps is the so-called Data cleansing or Data Scrubbing, which is the act of detecting and correcting (or removing) corrupt or inaccurate records from a data set.

The term refers to identifying incomplete, incorrect, inaccurate, partial or irrelevant parts of the data and then replacing, modifying, filling in or deleting this dirty data (when it's just too noisy to be fixed).

But what does it mean to have dirty data? to give you an idea, here are some examples:

- Dummy Values
- Absence of Data
- Multipurpose Fields
- Cryptic Data
- Contradicting Data
- Shared Field Usage
- Inappropriate Use of Fields
- Violation of Business Rules
- Reused Primary Keys
- Non-Unique Identifiers
- Data Integration

In practice, data cleansing consists of a few actions that are common to perform on raw data: for example, data may need to be parsed, which means to be converted from some kind of raw format (often plain text) to a more formal structure.  
Example:



Even after being parsed, data may still have missing fields or wrong values, which need to be corrected by using sophisticated data algorithms and secondary data sources.

Then, data may need to be standardized, by applying conversion routines to transform data into its preferred (and consistent) format using both standard and custom business rules, as well as coherent measurement units,

The next step is analyzing and identifying relationships between matched records and consolidating/merging them into ONE representation, which is usually done by pattern matching on historical data sources. If after the pattern matching two pieces of data appear to be referring to the same record, one of them is discarded.

Only after all these steps data can be considered reliable.

Sometimes data can come in a completely unstructured shape (like PDFs or free text), and complex algorithms need to be applied in order to extract some kind of structure from it.

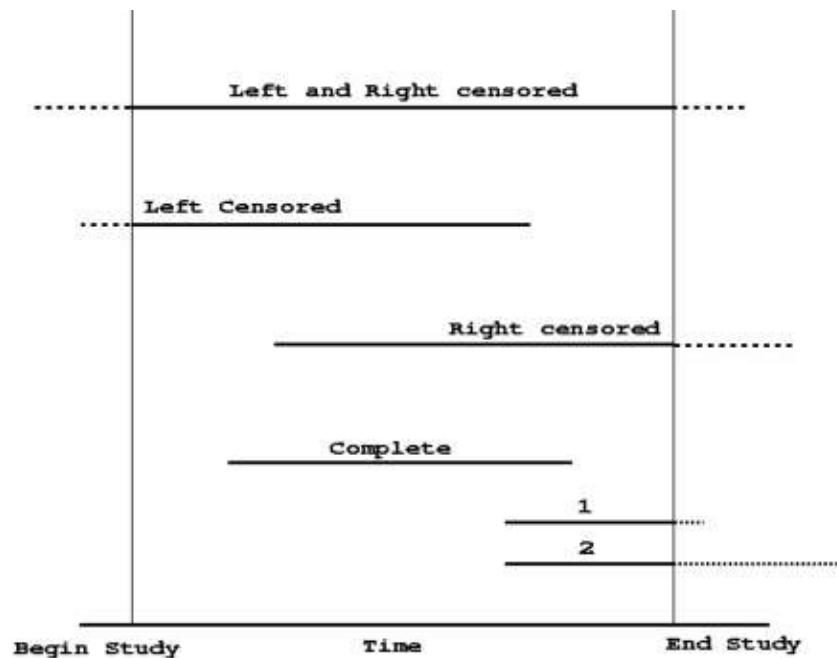
Sometimes it may be necessary to perform an operation called Data Munging, which are potentially lossy transformations applied to a piece of data or a file, or vague data transformation steps that are not yet completely clear (E.g.,

removing punctuation or html tags, data parsing, filtering, and transformation) Let's focus now a bit on maybe the most problematic issue when dealing with raw data: missing data. We may miss entire files or just some fields, or the data could be present but damaged.

There can be many reason for missing data: source system faults, ingestion errors (e.g. wrong schema), data who is partial by nature (like empty forms).

We can deal with these kind of situations in a number of ways:

- When there's a field missing, delete the entire record. This solution is very risky as can lead to losing big quantities of data and should therefore not be used if there's an alternative.
- Fill missing fields with estimators like the average. This method is convenient and simple but it assumes that the data distribution between records is uniform and ignores possible bias which could influence the result.
- Fill missing fields with more accurate algorithms based on attribute relationships.
- Use more complex methods like Markov Chain Monte Carlo who calculate the field by analyzing pattern in other records.

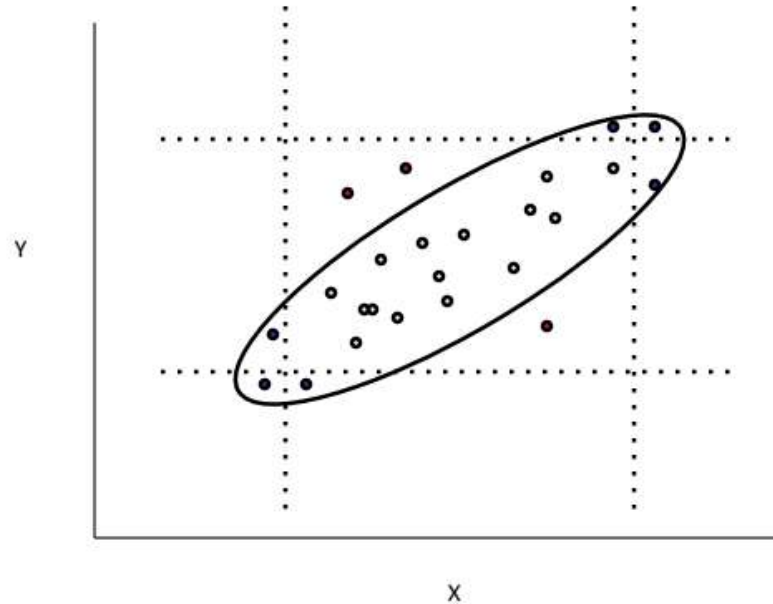


Censored time intervals

This picture shows another significant cause for missing data: interval based

data ingestion. In these cases it's very likely that the values ingested at the start and end of the process will be truncated.

Another possible cause for anomalies in data are outliers, records who have values completely outside of the data distribution which need to be eliminated. There are many techniques for detecting outliers, the simplest of which is plotting with control charts.



In the previous picture the points outside the black circle are outliers.

Another significant issue which needs to be taken into account is the "matching key problem", which can occur when we need to identify the same resource coming from different sources.

There are many tools which help dealing with this issues, ranging from programming libraries (e.g. Numpy, Pandas and SkLearn for Python) to full fledged applications like Trifacta Wrangler.