



Analysis of Community Interactions using Spark GraphX

Faculty: Michael Enudi

About Me.



- › Lives and works in Johannesburg, South Africa
- › Senior Software engineer with over 10years of working experience writing enterprise java applications, architecting data solutions.
- › Cloudera Certified Spark and Hadoop Dev.
- › Oracle Certified SQL Expert
- › Oracle Certified Java Master
- › Sun Certified Java Business Component Dev.
- › Sun Certified Java Programmer
- › Big data enthusiast

LinkedIn → <https://www.linkedin.com/in/michaelenudi>

What are Graphs??

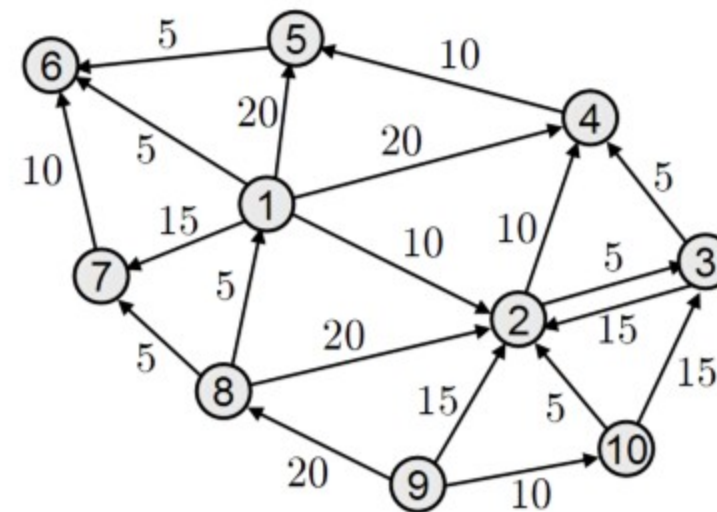
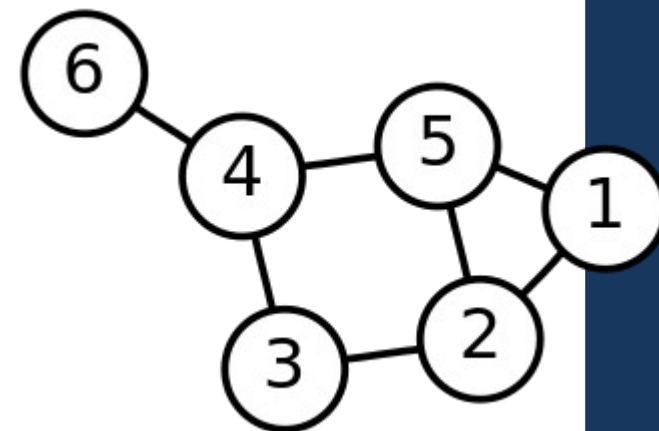
Graphs are mathematical structures used in computing to model entities (of same or different types) and their relationship.

It is based on the concept or subject of mathematics called graph theory.

The objects or entities are called Vertex/Vertices or Node/Nodes while the relationship are called Edge/Edges.

From the [documentation page...](#)

GraphX unifies ETL, exploratory analysis, and iterative graph computation within a single system. You can view the same data as both graphs and collections, transform and join graphs with RDDs efficiently, and write custom iterative graph algorithms using the Pregel API.



Types of Graphs

- Simple graphs
- Bipartite graphs
- Directed and Undirected graphs
- Multigraphs
- Property graphs
- [See more ...](#)

lines or curves for the edges. Graphs are one of the objects of study in [discrete mathematics](#).

The edges may be directed or undirected. For example, if the vertices represent people at a party, then person A shakes hands with a person B only if B also shakes hands with A . In contrast, if any edge from a person A to a person B represents a directed edge, then the former type of graph is called an *undirected graph* and the edges are called *undirected edges*.

Graphs are the basic subject studied by [graph theory](#). The word "graph" was first used in the 18th century.

Contents [\[hide\]](#)

1 Definitions

1.1 Graph

1.2 Adjacency relation

2 Types of graphs

2.1 Distinction in terms of the main definition

2.1.1 Undirected graph

2.1.2 Directed graph

2.1.3 Oriented graph

2.1.4 Mixed graph

2.1.5 Multigraph

2.1.6 Simple graph

2.1.7 Quiver

2.1.8 Weighted graph

2.1.9 Half-edges, loose edges

2.2 Important classes of graph

2.2.1 Regular graph

2.2.2 Complete graph

2.2.3 Finite graph

2.2.4 Connected graph

2.2.5 Bipartite graph

2.2.6 Path graph

2.2.7 Planar graph

2.2.8 Cycle graph

2.2.9 Tree

2.2.10 Advanced classes

3 Properties of graphs

4 Examples

5 Graph operations

6 Generalizations

7 See also

8 Notes

9 References

10 Further reading

11 External links

Graph storage or graph databases

A group of NoSQL databases that provide full transactional CRUD support for data modelled as graphs.

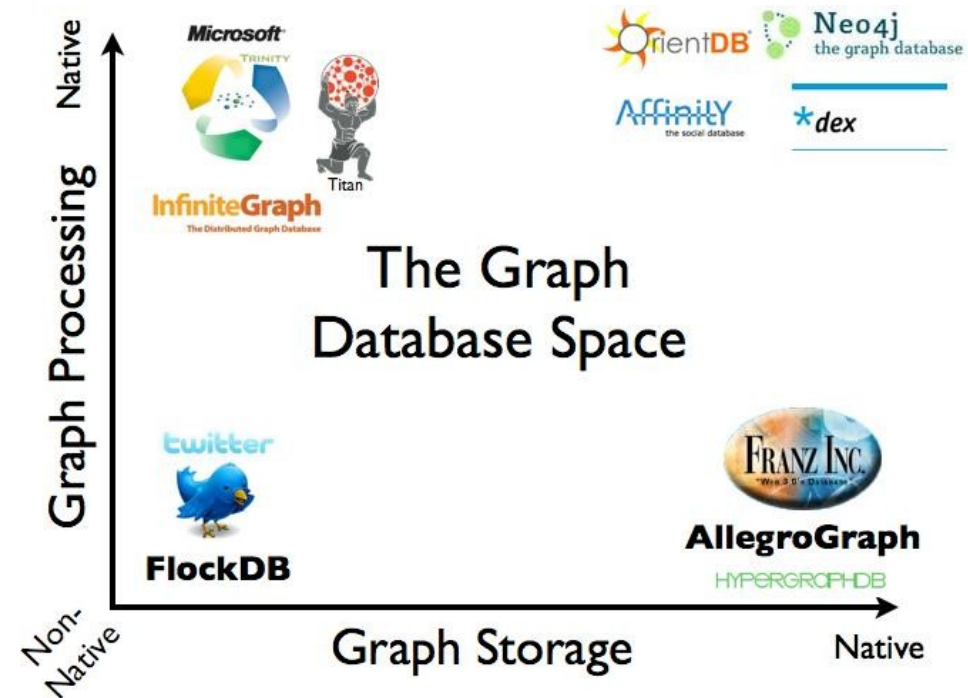
They

- Offer real time query support
- Are engineered for transactional integrity
- Can served as data sources for big data real time systems.

•

They can be used as a store for final or intermediate result of graph processing frameworks.

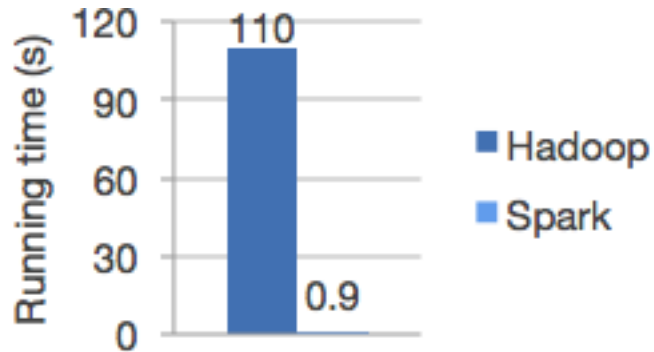
An example is MazeRunner extension that enables Spark application to read data from Neo4J database, process and store them back to the graph database.



Our next hackerday will cover [Modelling and Thinking in Graphs \(Neo4J\)](#).

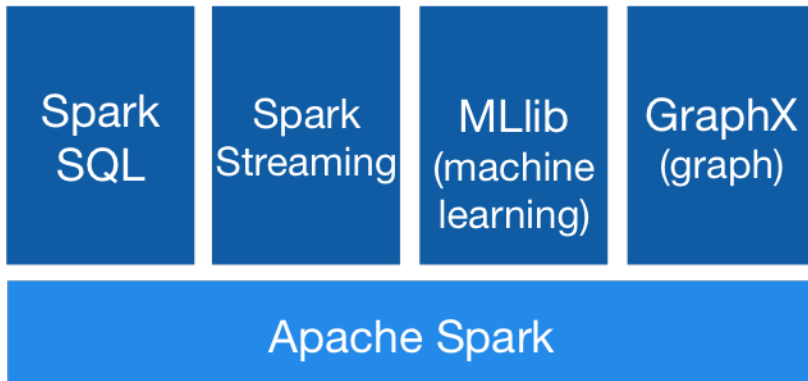
Graph processing with Spark

Apache Spark



An engine for fast and efficient data processing. It is a natural success of Hadoop's MapReduce framework that comes with in-memory data structure (RDD), directed-acyclic graphs of steps in the job pipeline and multi-language platform to write data processing application on and beyond Hadoop.

It is seriously becoming a sub-ecosystem as it runs other sub framework for data processing like Spark SQL, Spark Stream, Spark ML and many more.



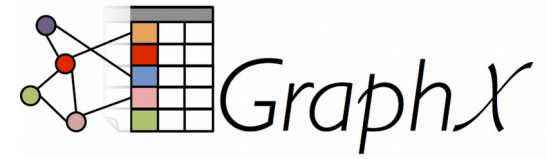
Scala and its build tool (Sbt)



Scala is a general purpose language that mixes various paradigms and runs on the java virtual machine (JVM). With Scala, you can perfectly mix functional and reactive styles of programming with Object-oriented programming.

SBT or Scala build tool - the maven for Scala.

Spark GraphX



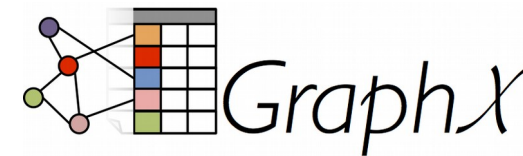
GraphX is the graph processing layer on top of Spark that brings the power of Big data to Graph processing.

It extends the RDDs of Spark computing framework as well creates a new abstraction called Graph.

Graphs in GraphX are treated as directed property multigraphs. The purpose of GraphX is enabled graph processing on big data.

To support this processing, GraphX helps us with a some fundamental graph operators as wells built-in algorithms for graph analysis and operation

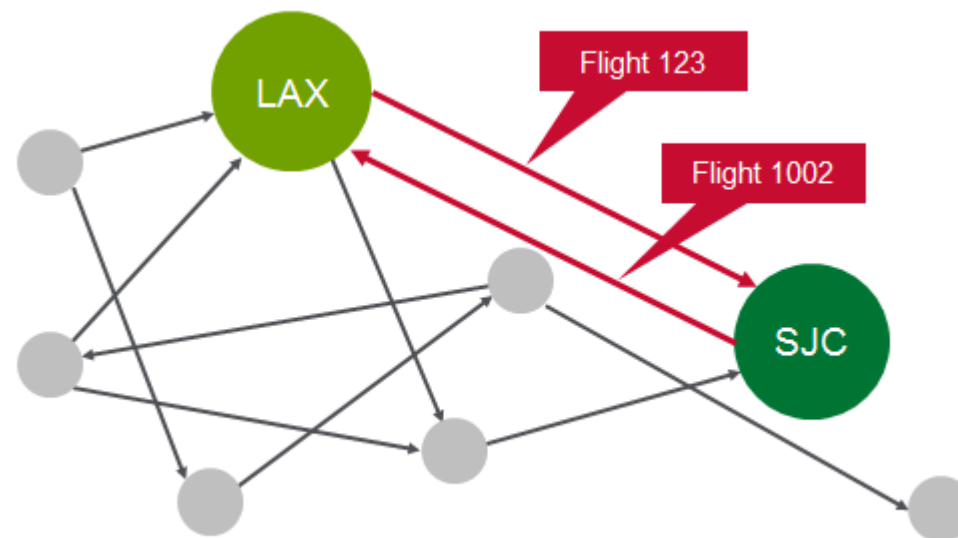
GraphX APIs



Spark graphX provides us with fundamental graph operators that make it easy to quickly assess or perform some operations on our graphs.

A sample list of operators include:

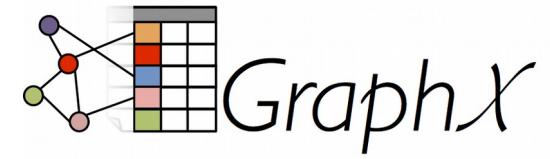
- inDegrees, outDegrees
- mapping edges, vertices or triplets
- subgraph
- aggregateMessages
- joinVertices
- groupEdges
- collectNeighbors



Also, some iterative algorithms like pageRank, triangle counting, connected components are available in Graphx as a single api call.

Thus we can perform iterative parallel-computations on big data modelled in a graph structure.

Spark GraphX Visualization



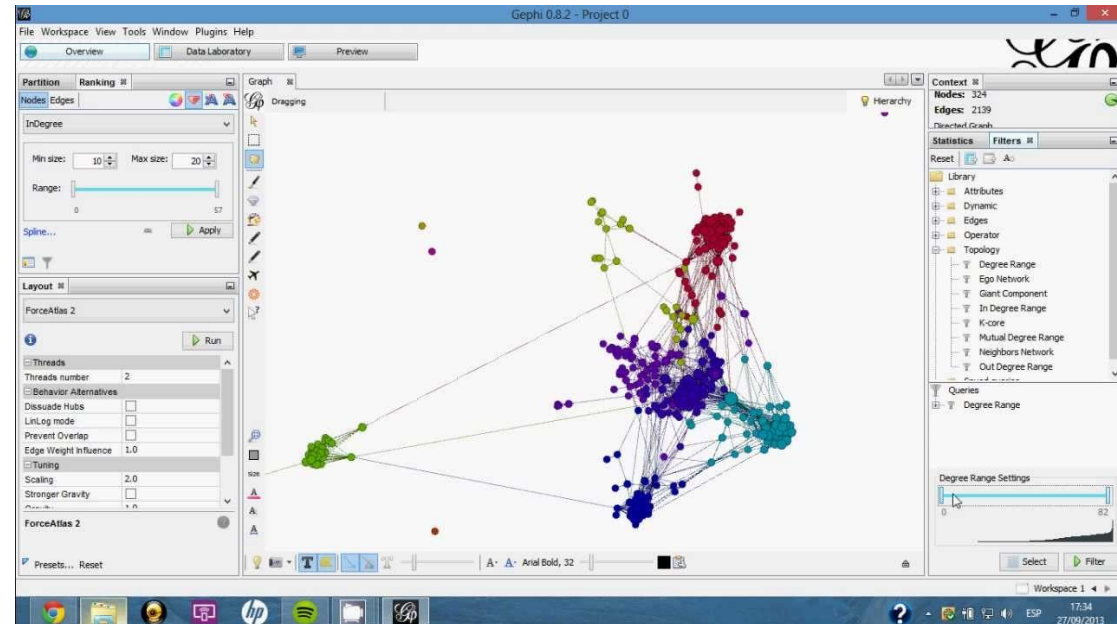
There are no visualization features with GraphX.

- GraphX is geared towards processing of large dataset such that visualization is not readily possible. Example, is visualizing a graph of 1 millions nodes and 200 millions edges.
-
- Just as Spark is for data processing, so is GraphX for graph processing.

My thoughts

To visualize graph, other tools have to be brought into the mix. They include

1. Gephi
2. Graphstream
3. Zeppelin



Case study

Requirement

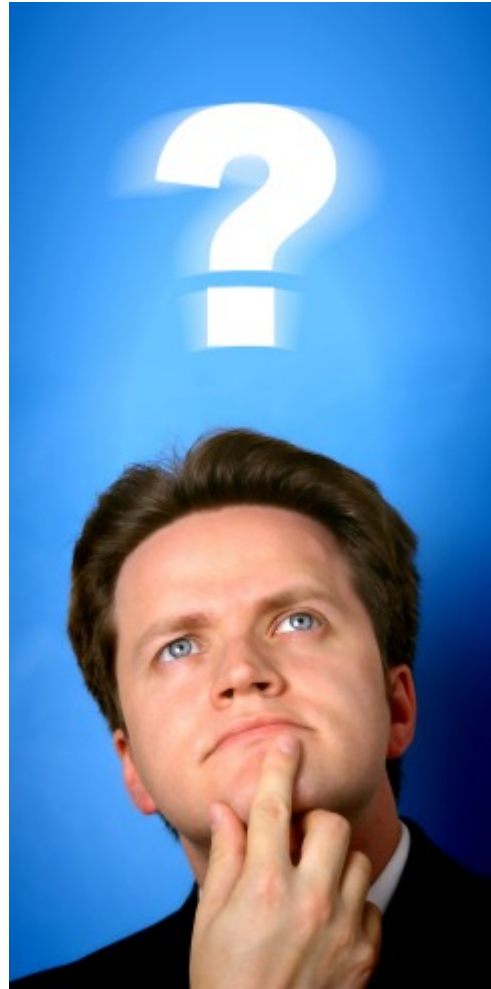
Given the dataset, XYZ telecommunication company has limited resources to invest in improving their quality of service. Hence the upgrade across the city of Milan will be done in phases.

Given that it is important that we rank traffic flows across the city so as to know optimize call quality with overall target on profit, a decision has to be made as to which grid station will be upgraded in what batch. This decision is going to be made based on daily call traffic in both direction.

To make which parts of the city should benefit from the upgrade, XYZ will like to know

- the grid connections with most traffic
- the grid with most incoming and outgoing call
- the grid with the most interactions to all other grid (the center of gravity)

.... and as much more information you can provide



**But I am not a Data
Scientist.**

**Where do this fit in the data
processing pipeline
???**

Dandelion Open Big Data [link](#)

Milano Grid

Some of the datasets referring to the Milano urban area are spatially aggregated using a grid. We refer to this grid as the Milano Grid.

<https://dandelion.eu/datagems/SpazioDati/milano-grid/description/>

Inter-city interaction

This dataset provides information regarding the directional interaction strength between the city of Milan different areas based on the calls exchanged between Telecom Italia Mobile users.

<https://dandelion.eu/datagems/SpazioDati/telecom-mi-to-mi/resource/>

Dandelion API

Open Big Data / Milano Grid

DescriptionTabular PreviewAPIResources

Some of the datasets referring to the Milano urban area are spatially aggregated using a grid. We refer to this grid as the Milano Grid.

Milano Grid schema

Name	Type	Description
cellId	Integer	the cell ID
geometry	Geometry	the cell geometry expressed as geoJSON and projected in WGS84 (EPSG:4326)

Description

The Milano Grid has the following spatial description:

[x₁,y₁]

9901	9902	9999	10000
9801	9899	9900
...								

[x₂,y₂]

Open Big Data / Telecommunications - MI to MI

DescriptionTabular PreviewAPIResources

Schema

1. **Square Id1**: the id of the square of the Milano GRID which is the origin of the interaction; TYPE: numeric
2. **Square Id2**: the id of the square of the Milano GRID which is the destination of the interaction; TYPE: numeric
3. **Time Interval**: the beginning of the time interval expressed as the number of millisecond elapsed from the Unix Epoch on January 1st, 1970 at UTC. The end of the time interval can be obtained by adding 600000 milliseconds (10 minutes) to this value. TYPE: numeric
4. **Directional Interaction Strength**: the value representing the directional interaction strength between the Square Id1 and the Square Id2. This value is proportional to the number of calls exchanged between callers located in the Square Id1 and receivers located in the Square Id2. TYPE: numeric

Important notes

If for a given combination of the Square Id1, the Square Id2 and the Time Interval no call is recorded the corresponding record is missing from the dataset. This means that a record of the form:

Square Id1, Square Id2, Time Interval, 0

will not be stored in the dataset.

See more details in the [Description tab](#)

Filename	Size	Created	Modified	Get It
december/full.zip	45.6 GB	Tue 21 Jan 2014 11:16	Tue 21 Jan 2014 11:16	Download
MitoMI-2013-11-01.zip	1.1 GB	Fri 10 Jan 2014 21:01	Fri 10 Jan 2014 21:01	Download
MitoMI-2013-11-02.zip	1.2 GB	Fri 10 Jan 2014 21:01	Fri 10 Jan 2014 21:01	Download
MitoMI-2013-11-03.zip	1019.3 MB	Fri 10 Jan 2014 21:01	Fri 10 Jan 2014 21:01	Download
MitoMI-2013-11-04.zip	1.8 GB	Fri 10 Jan 2014 21:01	Fri 10 Jan 2014 21:01	Download
MitoMI-2013-11-05.zip	1.9 GB	Fri 10 Jan 2014 21:01	Fri 10 Jan 2014 21:01	Download
MitoMI-2013-11-06.zip	1.9 GB	Fri 10 Jan 2014 21:01	Fri 10 Jan 2014 21:01	Download
MitoMI-2013-11-07.zip	1.9 GB	Fri 10 Jan 2014 21:01	Fri 10 Jan 2014 21:01	Download

Dataset Analysis

*We must note that we are using only a day's worth of data which in the real-world, is not worth making decision on.
We are doing this because of the sheer size.*

Milano Grid dataset

The city of Milan is spatial aggregated using a grid.

The dataset for this grid is a complex, multiple-nested json collection. Each item in the json collection contain the id, cellTowerId, geometry, shape, etc for that grid.

We will use features of Spark SQL to read the json file and construct an RDD from this dataset. Each item in this dataset will serve as a node in our graph.

One practical way to think of a grid is a region i in the city that will be served by cell tower $i - 1$.

[illegible]

Dataset Analysis

We must note that we are using only a day's worth of data which in the real-world, is not worth making decision on.

We are doing this because of the sheer size.

Aggregated CDR dataset

According to the description page, this dataset provides information regarding the directional interaction strength between the city of Milan different areas based on the calls exchanged between Telecom Italia Mobile users.

The directional interaction strength between the area A and the area B is proportional to the number of calls issued from the area A to the area B.

It is a tab delimited file that contains the following fields

- Timestamp
- GridID1 for the grid with the calling mobile number
- GridID2 for the grid with the receiving mobile number
- Directional interaction strength between the Grid id1 and the Grid id2. This value is proportional to the number of calls exchanged between callers located in the Grid id1 and receivers located in the Grid id2.

```
hadmin@hadoop-ss-lab: ~/sample-data/dandelion_mitomi
hadmin@hadoop-ss-lab:~/sample-data/dandelion_mitomi$ head MitoMI-2013-11-03.txt
-n 40
1383468000000 1 1 1.415182132451453E-4
1383474600000 1 1 3.6471276544197196E-4
1383476400000 1 1 1.5103247348560638E-4
1383477000000 1 1 3.6349656926791334E-4
1383483000000 1 1 1.478153689175954E-4
1383485400000 1 1 7.395236931786679E-5
1383499200000 1 1 2.172210647570028E-4
1383501600000 1 1 7.551623674280319E-5
1383507600000 1 1 7.22991321747922E-5
1383508200000 1 1 4.73328824751369E-6
1383510000000 1 1 7.22991321747922E-5
1383519000000 1 1 7.22991321747922E-5
1383464400000 1 10 2.814243229598746E-5
1383467400000 1 10 2.6943522612461817E-5
1383473400000 1 10 7.588125267884338E-5
1383475800000 1 10 1.8573731606172118E-4
1383489600000 1 10 1.6248613856696166E-4
1383492000000 1 10 2.591626890161992E-4
1383498000000 1 10 8.322838720443674E-5
1383498600000 1 10 1.0402368497483085E-4
1383500400000 1 10 7.925775136252494E-5
1383501000000 1 10 1.0620127397498675E-4
1383507000000 1 10 1.340085494001695E-4
1383460200000 1 1001 3.6865383222377156E-6
1383468600000 1 1001 5.63104352596541E-5
1383487200000 1 1001 6.740586716042852E-6
1383493800000 1 1001 1.5927026876965018E-4
1383496200000 1 1001 3.6865383222377156E-6
1383502200000 1 1001 5.9996973581891815E-5
1383469200000 1 1002 5.391152654393628E-5
1383472200000 1 1002 3.6865383222377156E-6
1383480600000 1 1002 3.6865383222377156E-6
1383505200000 1 1002 3.6865383222377156E-6
1383505800000 1 1002 4.314569789933472E-5
1383470400000 1 1003 5.3761298308481875E-5
1383487800000 1 1003 2.915793600063867E-6
1383506400000 1 1003 3.5196511189015013E-6
1383466200000 1 1004 4.6649398250341986E-5
1383488400000 1 1004 5.63104352596541E-5
1383489000000 1 1004 5.03359365725797E-5
hadmin@hadoop-ss-lab:~/sample-data/dandelion_mitomi$
```

Environment & tools

- Cloudera Quickstart VM 5.7 or 5.8
- Scala SDK and Runtime
- Scala build tool
- Gephi (optional)



Lab