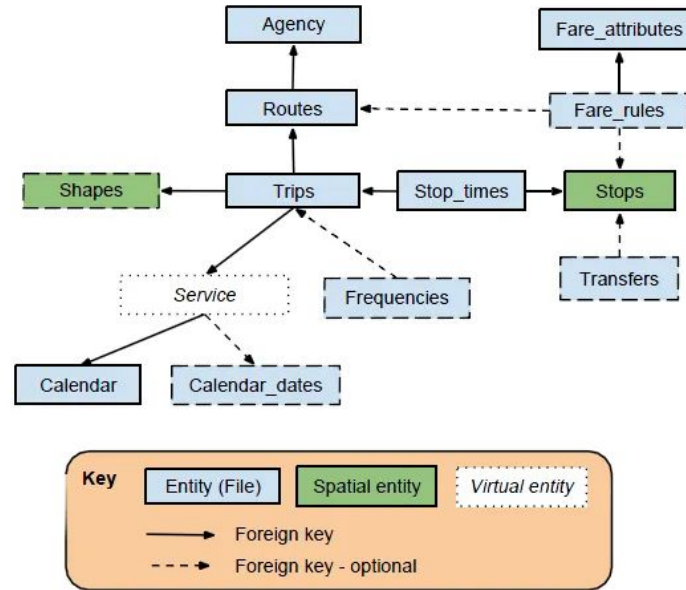# GTFS graph building and its path finding algorithms

# Summary

- Project overview.
- Our project architecture.
- Used technologies and external dependencies.
- How to load/build/persist the GTFS data graph ?
- How to search for a path in the GTFS graph ?
- How to cluster our GTFS graph ?
- Conclusion.

isep
École d'ingénieurs du numérique

# Project overview - what is GTFS ?

- GTFS stands for *General Transit Feed Specification*.

- A set of .csv (or .txt) files, containing data about a transportation network.

- Each file has potentially a foreign key field for it to be linked to another file. Just like the SQL tables.

# Project overview - our goals

What we needed to do :

- Collection of data and construction of the graph
  - Unweighted graph
  - Weighted graph
- Calculation of shortest paths
- Shortest paths for graph clustering

isep
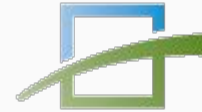École d'ingénieurs du numérique

# Project architecture

- **algorithms :** with all the advanced algorithms algorithms.
  - **clustering :** find the clusters, their barycenters and the distance between one another.
  - **shortestpath** : contains the algorithms to find the shortest paths.

- **nio** : for non-blocking input output ; it contains the logic for loading, building and eventually serializing the GTFS data into a graph. It supports remote operations with a GCP bucket.

- **resources** : the local data are stored in the resources folder.
  - **mbta** : A folder to store the GTFS files of the mbta dataset.



```
∨ src
  ∨ main
    ∨ java / gtfs / corev2
      ∨ algorithms
        > clustering
        > shortestpath
      > nio
      J BFSFrame.java
      J DijkstraFrame.java
      J EdgeType.java
      J GraphVisualizer.java
      J GTFSEdge.java
      J GTFSEdgeTemp.java
      J GTFSVertex.java
      J MainFrame.java
      J PathVisualizer.java
    > resources
  > test / java / gtfs / corev2
```
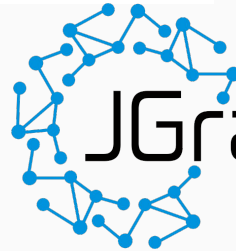
# Used technologies and external dependencies
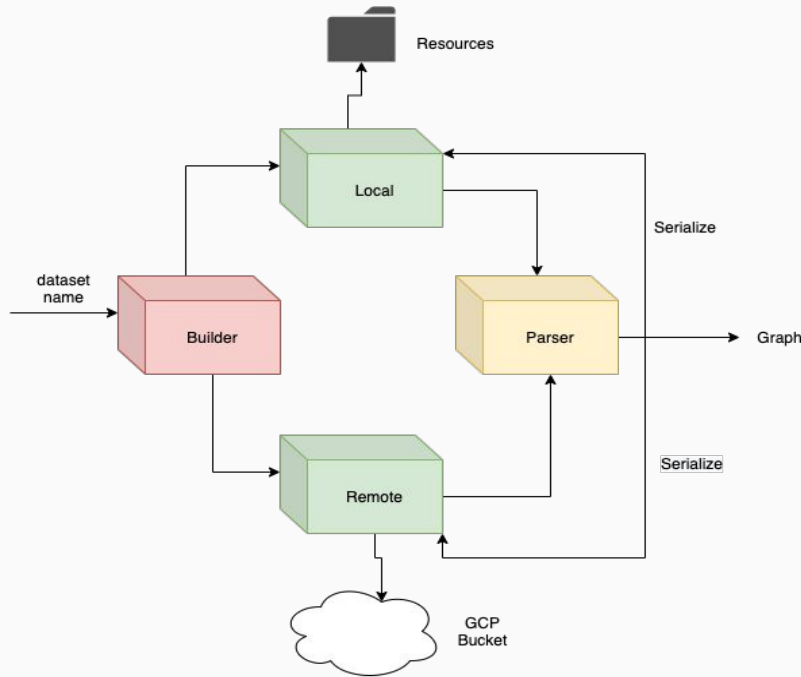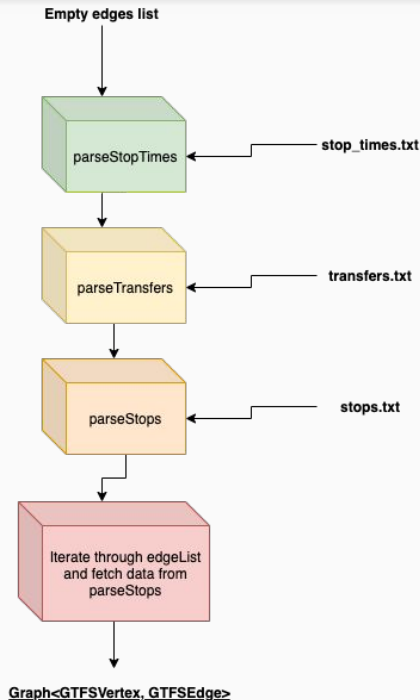
# How to load/build/persist the GTFS data graph ?



- **Builder** design pattern (*GTFSGraphBuilder class*)

- Can fetch data from local/remote destination in a multithreaded way.

- Parse the raw data loaded in memory and generate the *GTFSGraph<GTFSVertex, GTFSEdge>* object.

- Ability to serialize/unserialize (keeping the data integrity) the graph object under JSON format for a faster retrieval the next time (**100x** speed in local, **1000x** speed in remote)

# How to parse GTFS data ?



- Generate a List<List<String>> representing the edge list in parsing **transfers.txt** and **stop_times.txt**
  - use **parseStopTimes(edgeList)** function to mutate a Map<String, GTFSEdgeTemp> where the key is a code containing the starting stop id and the target stop id.
  - use **parseTransfers(edgeList)** function to mutate a Map<String, GTFSEdgeTemp> where the key is a code containing the starting stop id and the target stop id. The transfers add the travel within the same station whereas **parseStopTimes** is for two different stations.
- Generate a Map<String, List<String>> which is a representation of **stops.txt** where the key is the stop id and the value is the remaining data in the line
- Iterate through **edgeList** and for each edges find the corresponding vertex data (latitude, longitude and names in our simplified version) using the Map<String, List<String>> computed earlier. Create the real GTFS graph with weighted edges computed with the distance in km between two vertex.
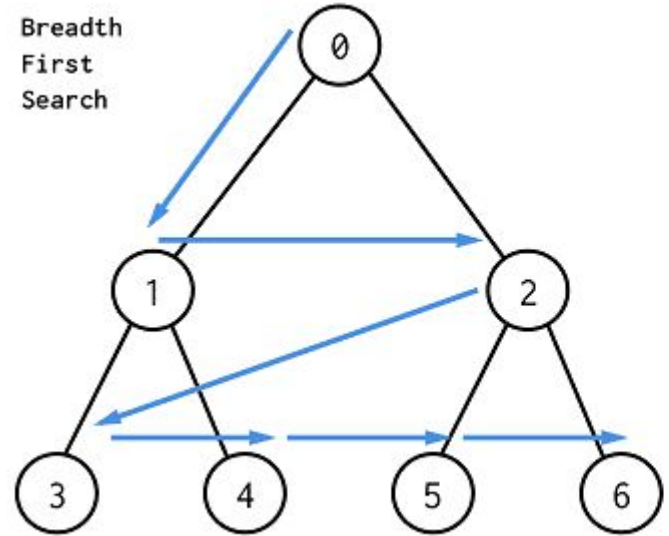
# How to search for a path in the GTFS graph ?

- Breadth First Search (BFS) in order to search for the best path in an Unweighted Graph
- Dijkstra to find shortest path in terms of distance in Weighted Graph
- Find path between two nodes or clusters

isep
École d'ingénieurs du numérique

# BFS

Nodes that have already been visited are marked to prevent the same node from being explored more than once.
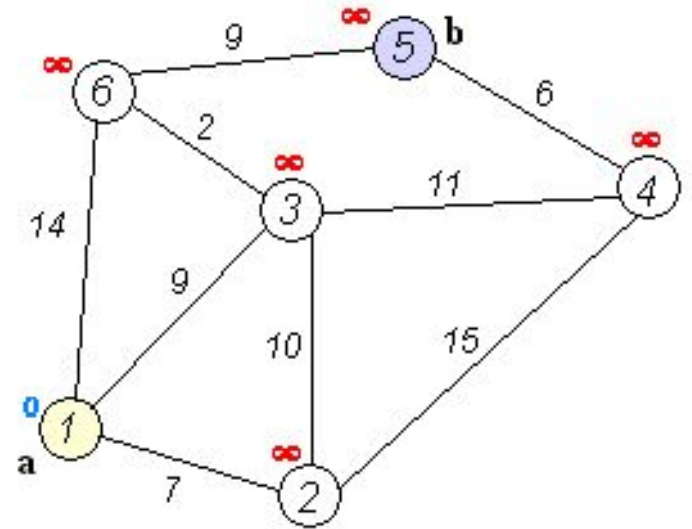
Algorithm Steps :

- Put the source node in the queue.
- Remove the node from the beginning of the queue for processing.
- Put all unexplored neighbors in the queue (at the end).
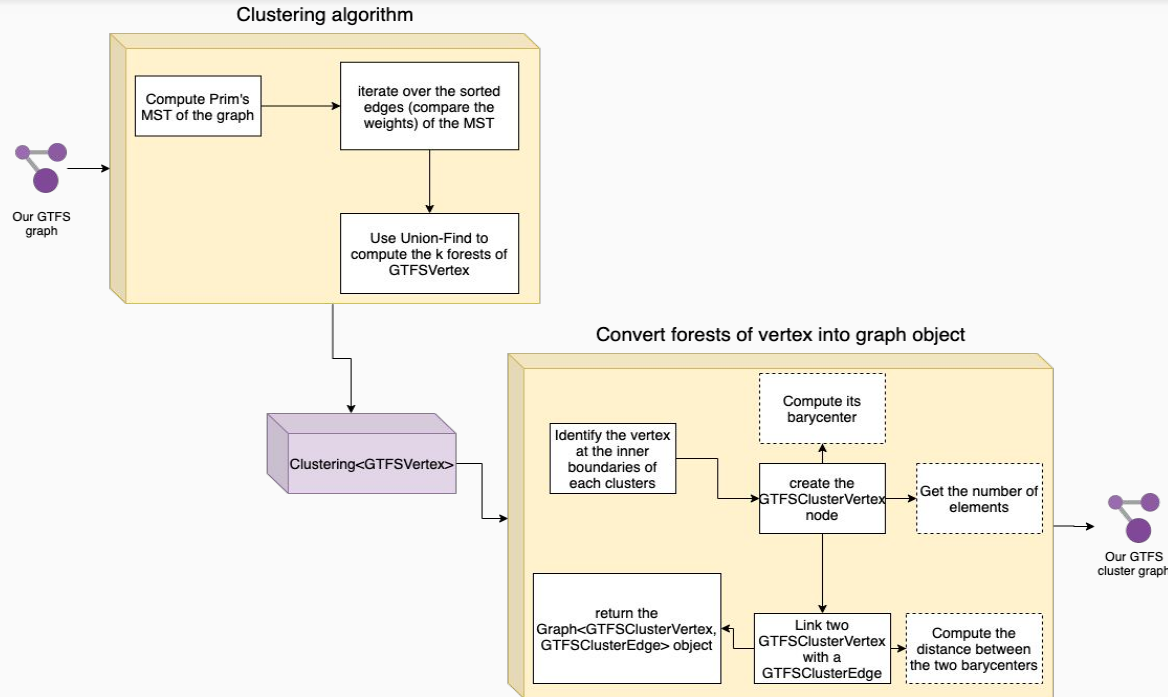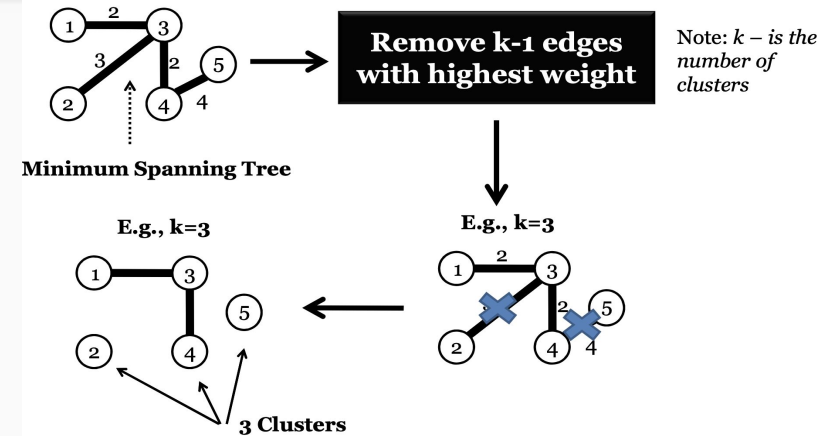- If the queue is not empty, resume step 2.

# Dijkstra

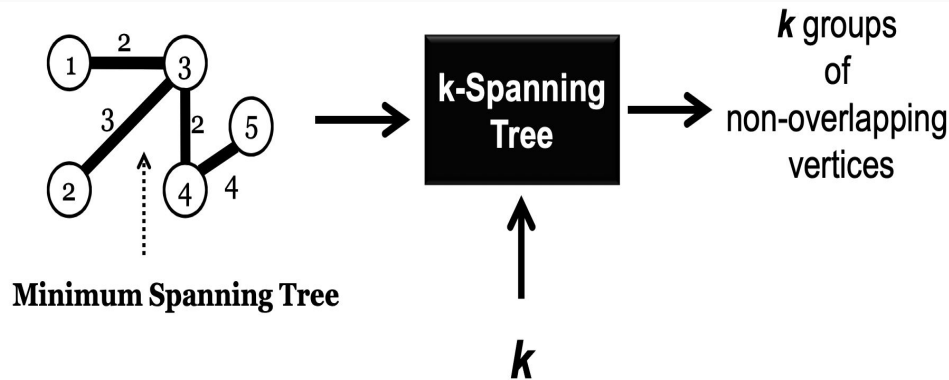- Build a sub-graph in which the different vertices are classified in ascending order of their minimum distance from the starting vertex.
- At the beginning, we consider that the distances from each vertex to the starting vertex are infinite
- During each iteration, a vertex of minimum distance is chosen outside the subgraph and added to the subgraph.
- The distances of the vertices adjacent to the added one are updated

# How to cluster our GTFS graph ?



**Clustering algorithm**

- Our GTFS graph
- Compute Prim's MST of the graph
- iterate over the sorted edges (compare the weights) of the MST
- Use Union-Find to compute the k forests of GTFSVertex

Clustering<GTFSVertex>

**Convert forests of vertex into graph object**

- Identify the vertex at the inner boundaries of each clusters
- Compute its barycenter
- create the GTFSClusterVertex node
- Get the number of elements
- return the Graph<GTFSClusterVertex, GTFSClusterEdge> object
- Link two GTFSClusterVertex with a GTFSClusterEdge
- Compute the distance between the two barycenters

Our GTFS cluster graph

**isep**
École d'ingénieurs du numérique

# How to cluster our GTFS graph ?



- Obtains the MST of the input graph.
- Obtains the list of sorted edges of the MST
- Remove k -1 edges from the MST. ( use of Union-Find )
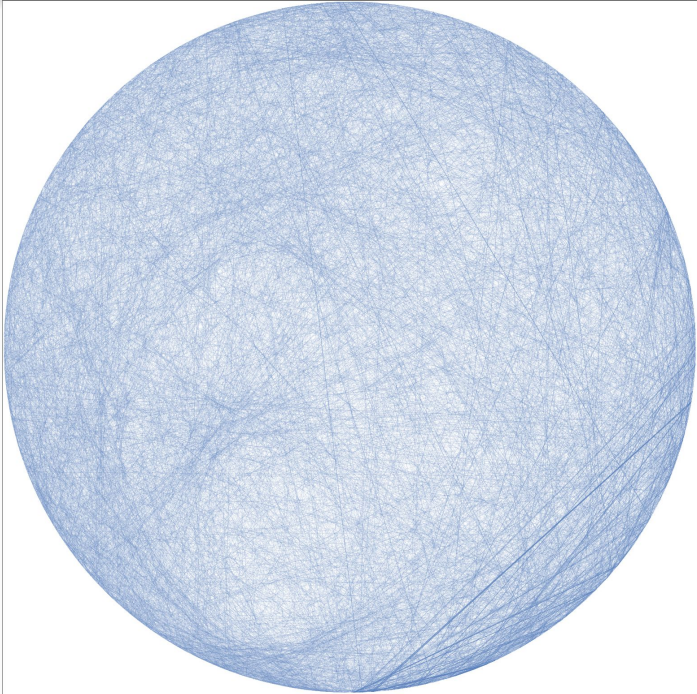- Results in k cluster. ( The k partition of the Union-Find )

isep
École d'ingénieurs du numérique

# Complexity of our cluster algorithm

$$O(m + nlog(n) + nlog(n) + klog^*(n) + (kn + (n-k)m))$$

Prim's MST     Sort MST's edges     Build the forests     Compute barycenters for each cluster     Link the disjoint clusters

$$O((n-k)m)$$

isep
École d'ingénieurs du numérique

# Conclusion



- Drawings of boston public transit network
- Almost 8000 vertices 11000 edges.

# Conclusion



- Drawings of the path between two nodes using BFS and Dijkstra