# An Introduction of Gephi

CITI Visualization Group

February 21, 2017

# Contents

# 1 Introduction

Gephi(https://gephi.org/) is an open-source visualization platform for graphs and networks. This documents focuses on giving readers a brief overview of this software by introducing its configurations, concepts, and functions. To access more detailed tutorials and demonstrations we refer readers to Gephi's official website: https://gephi.org/users/.

# 2 File Formats and Basic Data Structures

While Gephi can import data from various file formats as listed below, GEXF, CVS, and Spreadsheet are the widely adopted because of their simple structures.

- GEXF
- GDF
- GML
- GraphML
- Pajek NET
- GraphVis Dot
- CVS
- UCINET DL
- TLP
- Netdraw VNA
- Spreadsheet(Excel)

GEXF is a XML based structure which is able to encode almost all of the data structures supported by Gephi, while CVS and Spreadsheet can only include a few of them. However, CVS and Spreadsheet can be read and edited by a wider range of applications.

In addition to these formats, Gephi can import and export its own state files having an extension .gephi, that not only include original data but also encode operations and filters have been applied to the data. In other words, a .gephi file stores an actual visualization.

The essential data structure supported by Gephi are node, edge, and attribute, where nodes and edges are referred to as the network topology while attributes are referred to as network data. A node represents a basic entity in the data set and two nodes can be connected by an edge. Attributes are data associated to nodes or edges, e.g. string, integer, boolean, etc..

Figure 1 illustrates how the aforementioned data structures are encoded in a GEXF file. In this example, we declared a directed graph with a "source" node

```
<gexf xmlns:viz="http:///www.gexf.net/1.1draft/viz" version="1.1" xmlns="http://www.gexf.net/1.1draft">
<meta lastmodifieddate="2010-05-17+11:28">
<creator>Gephi 0.7</creator>
<description>Example of Node, Edge, and Attribute</description>
</meta>
        <graph defaultedgetype="directed">
                <attributes class="node" mode="static">
                        <attribute id="size" title="Size" type="double"/>
                </attributes>
                <nodes>
                        <node id="0" label="source">
                                <attvalues>
                                        <attvalue for="size" value="10"/>
                                </attvalues>
                        </node>
                        <node id="1" label="target">
                                <attvalues>
                                        <attvalue for="size" value="40"/>
                                </attvalues>
                        </node>
                </nodes>
                <edges>
                        <edge id="0" source="0" target="1" label="S to D"/>
                </edges>
        </graph>
</gexf>
```
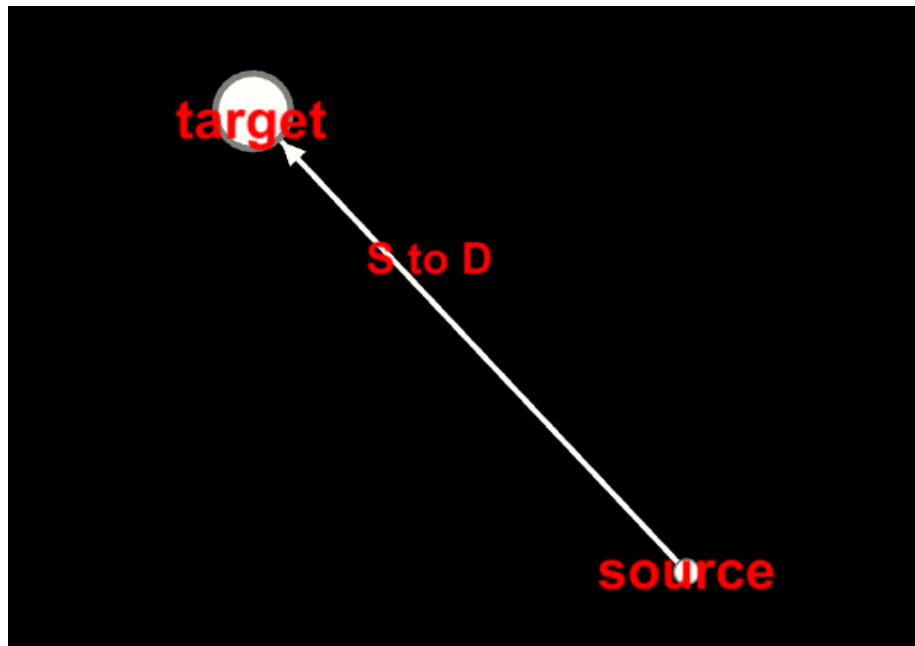
Figure 1: A GEXF file encoded two nodes attached with one attribute, and a directed edge connects them.



Figure 2: The corresponding visualization of the example illustrated in Figure 1, where sizes of nodes were determined by their "size" attribute.

and a "target" node that were distinguished by their unique node ids, i.e. 0 and
1. They were also superimposed with a "size" attribute which was represented
by a scalar value. An edge between them was declared to indicate a directed
path from source to target. The corresponding visualization of this example is
shown in Figure 2, where the sizes of nodes were assigned based on the "size"
attribute.

# 3   Creating Visualizations

This section will introduce the basic usage and functions of Gephi by demon-
strating the process of exploring an Airline sample data set.

To load the data, we can simply click the **File** drop-down menu, press **Open**,
and select *airlines.gexf*. Once the data has been imported successfully, Gephi



Figure 3: A default view the airline data set in Gephi.

will create a default display which is similar to Figure 3. As seen in this fig-
ure, Gephi contains four primary panels: **Appearance**, **Layout**, **Graph**, and
**Filters and Statistics**. Appearance panel provides the user interface to ma-
nipulate rendering options of nodes and edges, i.e. color, size, label color, and
label size. Layout panel lists a set of built-in algorithms to adjust the graph
layout. As shown in Figure 4, we can configure the properties and parameters
of a particular algorithm by selecting it in the drop down menu, and apply the
algorithm by clicking the **Run** button. Graph panel shows the actual visual-
ization, and also provides a suite of user interfaces, e.g. selection of entities,
creation of nodes and edges, configuration of displays, etc.. Filters tab provides
a list of filters and operations that can be applied for supporting specific queries
of the dataset, while the Statistics tab implements statistical calculations that
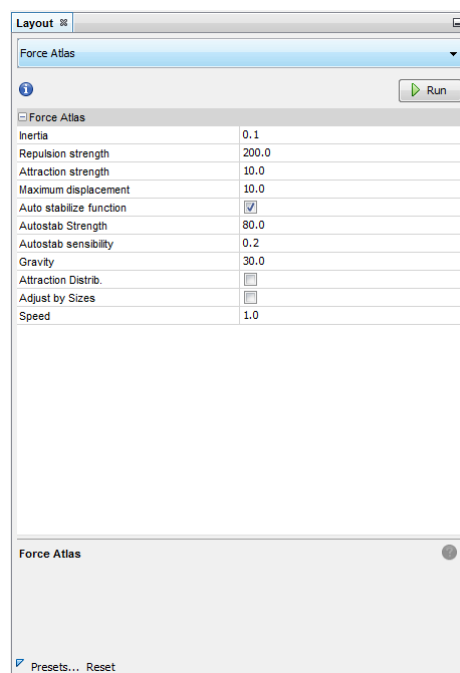can be used for analysis of the dataset.

Figure 4: Selecting Force Atlas algorithm in Layout panel.

Even though the default visualization shown in Figure 3 may lead viewers to some superficial interpretations of the data set, for instance there are many airports and they established flights between one another, it can hardly support more detailed explorations accurately, such as where are the airports, which airport maintains the most airlines, etc.. The first step to enhance the visualization is to superimpose a geographical map by applying a **Map of Countries** layout algorithm. Figure 5 shows our specified configurations of the layout, meaning

Figure 5: Configurations of Map of Countries layout algorithm.

that the display region of the map is restricted to U.S. After running the algorithm, we obtained an incorrect result, as shown in Figure 6, because of the inconsistence between geological information encoded by the dataset and that of the map. This can be fixed by applying a **Geo Layout** algorithm, with Latitude and Longitude parameters setup correctly, as shown in Figure 7.

In Figure 7, each circle represents an airport on the map, and a straight forward idea to improve this visualization is to encode the number of airlines maintained by individual airports by sizes of the circles, i.e. an airport has more airlines will be rendered as a bigger circle, and vice versa. We can simply make the modification in Appearance panel by changing the size of nodes. However, since the map shown in the visualization is also constructed by nodes and edges, a direct manipulation of all nodes will also change appearance of the map, resulting in undesired display. This can be avoided by creating filters to extract nodes only associated with airlines.s To achieve the proper filters, we first examine properties of all nodes and edges contained in the visualization by clicking the **Data Laboratory** button. Figure 8 exemplifies the data table, and we notice that a *backgroundmap* attribute is associated to each node
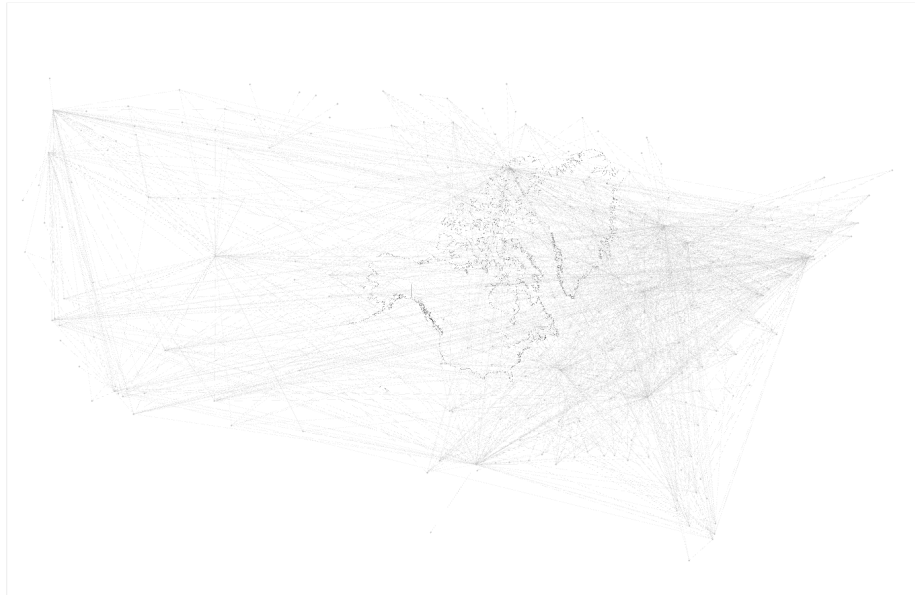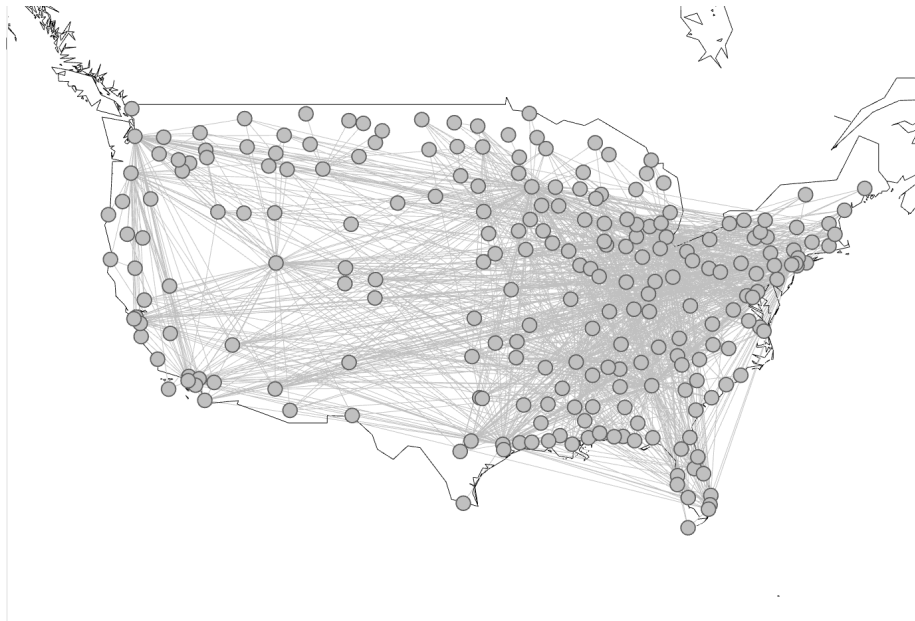
Figure 6: An incorrect map layout.



Figure 7: A correct map layout adjusted by Geo Layout algorithm.

Figure 8: Data laboratory shows attributes of all nodes and edges of the visualization.

to indicate whether it belongs to the background map or not. Now we can go back to **Filter** tab and extract all background nodes by clicking *Attributes*, *Non-null*, *backgroundmap*, and the filter will be listed in the **Queries** text box. Then we attached this filter to a *NOT(nodes)* operation, meaning that we select all nodes except those we obtained from the filter. Figure 9 shows an updated visualization where the *Queries* box at the lower right corner listed the filters and operations been conducted.

With the filters running we are able to manipulate appearance of airport nodes by clicking **Nodes** button in **Appearance** panel. We select size option and click **Ranking** button to vary node size based on one of its attached scalar attributes. In this example, we associate the *degree* attribute, which indicates the number of airlines connected to individual airports, to node size, meaning that an airport shown by a larger circle is considered to be busier. Once we apply modifications of appearance of nodes, we stop the filter created before to restore the background map to the display, as shown in Figure 10.

Similarly, we can also adjust color codings to represent additional information of the graph in **Appearance** panel. For instance, we can associate node color to **betweenness centrality** attribute, which measures the number of shortest paths between any two nodes that passes through a particular node. In other words, an airport node has higher betweenness centrality is connecting various different parts of the airline network, being considered to be more important globally. Figure 11 exemplifies this modification, where darker green represents higher betweenness centrality, and vice versa. Not surprisingly, we can notice that bigger nodes tend to have darker colors, meaning that busier airports are also playing more important roles in this network.

Figure 9: An updated visualization only shows nodes that are not background nodes.



Figure 10: An updated visualization where the node size is varied based on the **degree** attribute.
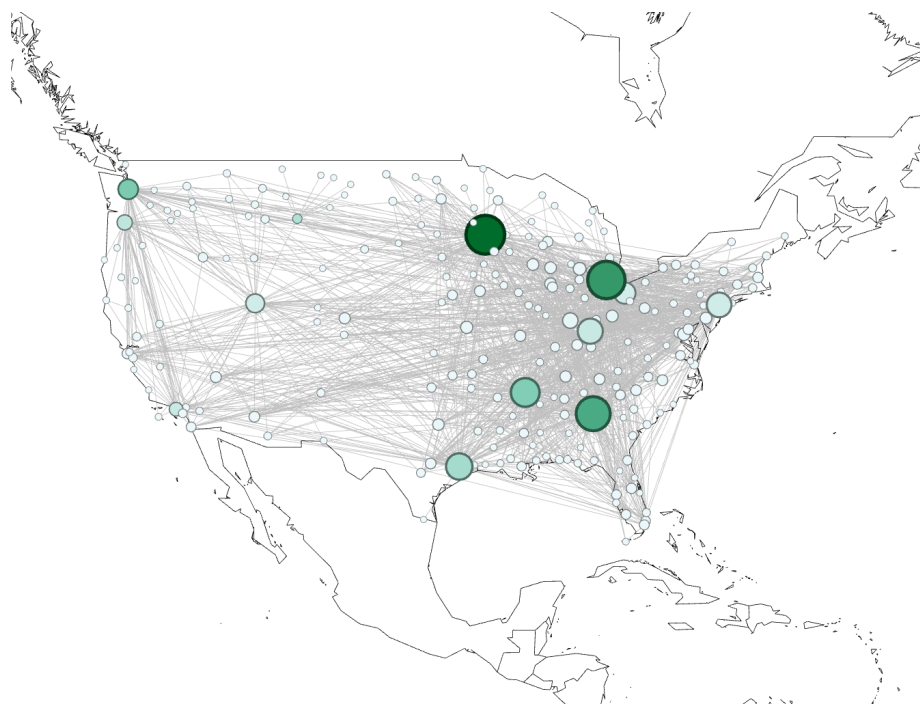
Figure 11: An updated visualization where the node color is assigned based on the **betweenness centrality** attribute.

The next straight forward question viewers would like to ask about this visualization is "what is the biggest and darkest airport?" This query can be answered by using Gephi's label and selection tools. We first open the label
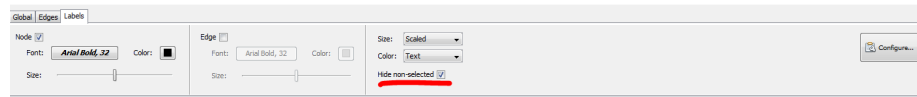


Figure 12: Configurations of label.

configurations panel, as shown in Figure 12, by clicking the small triangular icon located at the lower right corner of Graph panel, and then check the *Hide non-selected* option, as highlighted in the figure. We now can enable the *Direct*
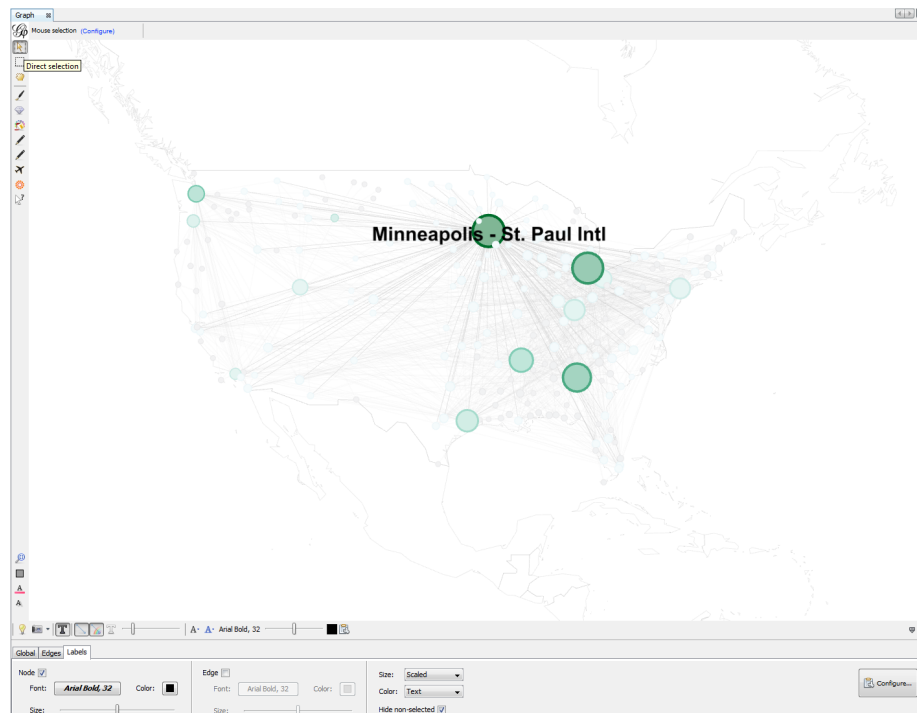


Figure 13: An example of selecting the biggest node in the network.

*selection* tool that located on the right side of Graph panel, and move the cursor onto a node we are interested in to see name of the corresponding airport, as shown in Figure 13.

Above we introduced the basic procedures of creating a visualization in Gephi, the final step is to output our visualization to an image file by using Gephi's **Preview** panel, as shown in Figure 14. It allows user to adjust rendering configurations, e.g. opacities of nodes and edges, showing labels or not,

Figure 14: Preview panel in Gephi.

drawing bounding boxes of texts or not, etc., and provides a preview of the configuration. After all settings have been specified, we can click **Export** button to output the display with supported file formats.