# WTO Merchandise Trade Analysis using Neo4j*

Deepak Khirey[†]
*Indiana University, Bloomington*

Gautam Matkar[†]
*Indiana University, Bloomington*
(Dated: July 20, 2019)

Merchandise Trade in the modern world is typically driven by the difference in technology, resources, demand, Existence of Government policies and the existence of economies of scale in production. Merchandise trade [1] is characterized by merchandise exports and imports. Typically, there are two systems to record merchandise trade referred to as general trade and special trade. Special trade excludes certain trade flows. Merchandise trade is typically defined by general trade which covers inward and outward flows of goods through country or territory. Goods are merchandise which adds or reduces the stock of materials because of inward or outward flow through country or territory. The flow of good is valued at transaction values [2]. The trade values can be a good indicator of overall trade in International market, the trend for trade flows among developed and developing countries, growth gap between various economies, trade shares of different economies [3]. We are proposing here to build a analyse merchandise trade transactions over the years as per data collected by WTO using Graph Database Neo4j.

## I. INTRODUCTION

In the last few decades, the global economy has opened and trade has increased with a rapid rate. As per report Latest trends in world trade 2017-2018 Chapter III [4], World merchandise exports increased to US$ 17.73 US trillion in 2017, up from US$ 16.03 US trillion in 2016, World commercial services exports grew by 8% in 2017, reaching US$ 5.28 trillion, Asia was the top contributor to trade growth in volume terms in 2017, growing by 8%, World merchandise trade grew by 4.7% in 2017 in volume terms, driven by a rising demand for imports across the world. These numbers are massive but at the same time do highlight that these trends will change at a faster rate than any time before. These trends can provide great insight about the global economy. This data is quite intuitive for statisticians and economist but not intuitive for another genre. This information can be visualized in an intuitive way at a single place rather than residing in multiple manuals or papers. Another aspect is the capability of localizing or analyzing local trends rather than just global trends. These trend scan be analyzed based on time series and sliced diced by products, type of merchandise goods, export and import. In this study, we aim to create visualizations which bring together trends at a single place. We aim to analyze this dataset further by using Graph Algorithm techniques which will explore regional trades as Network among different conutries.

## II. RESEARCH QUESTION

In this study, we will analyze WTO merchandise trade values annual dataset indicating the trade pattern and present indicators of Political policies, Trade expansion, Regional cooperation and trade balance using Graph Algorithms.

## III. RELATED WORK

Official website of the World Bank [5] has shown non-hierarchical treemap based visualization. We feel it does a decent job in representing the data. However, it doesnt clearly depict the relationship between various attributes of the data.It also lack's the geolocation artifacts.Treemaps are good in representing the portion to the whole ration but might not capture the explicit comparison between data attributes. In this project, we hope to solve these issues. One way to look at this data is to represent in terms of trade balance which gives a fair idea of the gap between import and export for each country.While this is excellent graphics provided by WTO itself,it fails to recognize regional cooperation.It also does not show categorically as to how the center of world trade has shifted from the USA to China over the years.Here are examples of various visualizations provide by WTO which represents the data: [6] [7] [8]

## IV. DATA AND METHOD

### A. Dataset Summary

The data for this study is sourced from 'WTO — Trade Statistics - Bulk download of trade datasets' with a specific focus on dataset WTO merchandise trade values an-

nual dataset'[9]. The data is in CSV format. The following table represents various columns and their sample distinct values to get the feel of the data. Highlighted columns are the target artifacts on which visualization would be focused. We have close to 959,601 data points. [FIG 1]

## B. Data Preparation

WTO has provided this dataset in CSV format with delimiter as comma.It does not have any NULL or NaN values. All fields are present in all rows.It is a very clean and neat data source from Data Quality perspective.However we need to make few changes in the dataset to make it suitable for our visualization.

**Data Anomalies -**
When we attempted to import this dataset into a dataframe, few rows were found to have more tokens than 14. This was happening because there was a comma present in the country name itself ie in Reporter Description and Partner Description fields. eg "Yemen, Republic of" So we had to replace such comma in country name with dash character. We did this activity manually with help of text editors. This helped us to resolve issue of tokenization.

**Data Filters -**
Once we imported dataset into dataframe, with further insights we realized that "Flow_Code","Unit","Flag","Source_Description","Note" fields do not have much data in it and we can safely drop those fields without affecting visualization outcome. We also saw that few rows have trade value as zero. We decided to drop those rows as well because they will not show up on visualization. Both these filters helped us to reduce overall data size which we can manage for visualization.

## C. Data Visualization

We performed Exploratory Data Analysis of WTO Merchandise data with help of Tableau Desktop visualization software.Tableau provides very intuitive and quick way to represent data in various visualization formats.We decided to go for Heatmap and Line Chart options.

**Line Charts**[FIG 4] [FIG 5] -

Since WTO dataset contains information from 1947 to 2018, it becomes very interesting see the trend in global trade over the years. Line chart was most effective way to visualize time series data. We decided to make separate line charts for Export and Import transactions.

Also since these transactions are categorized into multiple buckets like 'Food', 'Fuels','Chemicals' etc., we decided to define separate line for each trade category in the same plot. This provides an opportunity to visually compare trends in each of these category over the years. Both line charts depict 3 major milestones of world economy. We can see that there is significant increase in activities after 1990 which is when globalization started taking place. As a result of globalization, there were many avenues of increased trade activity. This was manly lead by Machinery category and Fuel category, which are basic infrastructure for any kind of business.
Then we see another spurt in activities by year 2000. We can see a sudden jump in trade activities by many folds. This can be attributed to growth in Telecom and Electronics field.
Then we see a major dip in businesses in the year 2008. It represents 'The Great Depression' in the USA after the fallout of of sub-prime crisis. We can see that trade is affected across all categories during this period.
However, the world seems to have recovered from the after-shocks of recessions and again back to growth trajectory. We can see that after 2016, there is a small jump in trade activities. We can see Fuels and Mining products are showing an anti-pattern here which may be an indication of another recession in coming years.

**Heatmaps** [FIG 2] [FIG 3] -

While line charts are great for trends, they do not shade light on the comparative volume of trade by various countries. Since this dataset provides data of almost 400 economic entities, it is not convenient to draw lines for all entities in a line chart. It will make the visualization 'crowded' and non-intuitive. Hence, we decided to go with heatmap which perfectly shows the trade volume comparison. Here, the block sizes correspond to trade value, so we can easily see significant entities in global trade going by bigger block sizes. Again, for simplicity of understanding, we chose to split our visualization for Import transactions and Export transactions.
We can see from heatmap, that 'world' is biggest trade entity. It does not make much sense because world is kind of abstract, but this can be attributed to the way WTO has received data from member countries. Unfortunately, all countries do not provide granular details of their trade with each country. In both Import and Export heatmaps, we can see that global trade is dominated by Machinery category and fuel category. We can also see that Economic groups like Euro, NAFTA,ASEAN are biggest trade partners of each other.This is followed by China, which is biggest individual Exporter and Importer of goods. This is consistent with the fact that China is leading economic activities on all fronts. China is closely followed by USA, Japan in some categories of trade.

| | Column Name | Distinct Count | Sample Values |
|---|---|---|---|
| 0 | Reporter_code | 269 | ['AF' 'afr' 'AFR' 'ACP' 'AL'] |
| 1 | Reporter_description | 268 | ['Afghanistan' 'Africa' 'African, Caribbean and Pacific States (ACP)' 'Albania' 'Algeria'] |
| 2 | Partner_code | 52 | ['WL' 'OAF' 'CIS' 'NI4' 'IN'] |
| 3 | Partner_description | 49 | ['World' 'Other Africa' 'Commonwealth of Independent States (CIS), including associate and former member States' 'Four East Asian traders' 'India'] |
| 4 | Indicator_code | 31 | ['TO ' 'AGFO ' 'MA ' 'MAMTOF ' 'AG '] |
| 5 | Indicator_description | 31 | ['Total merchandise' 'Food' 'Manufactures' 'Office and telecom equipment' 'Agricultural products'] |
| 6 | Flow_Code | 5 | ['X ' 'M ' 'RX' 'DX' 'RM'] |
| 7 | Flow_Description | 5 | ['Exports' 'Imports' 'Re-exports' 'Domestic exports' 'Retained imports'] |
| 8 | Year | 71 | [1948 1949 1950 1951 1952] |
| 9 | Unit | 1 | ['million USD'] |
| 10 | Value | 642296 | [49. 58. 47. 52. 53.] |
| 11 | Flag | 4 | [' ' 'E ' 'B ' 'E'] |
| 12 | Source_Description | 1 | ['WTO'] |
| 13 | Note | 1 | [' '] |

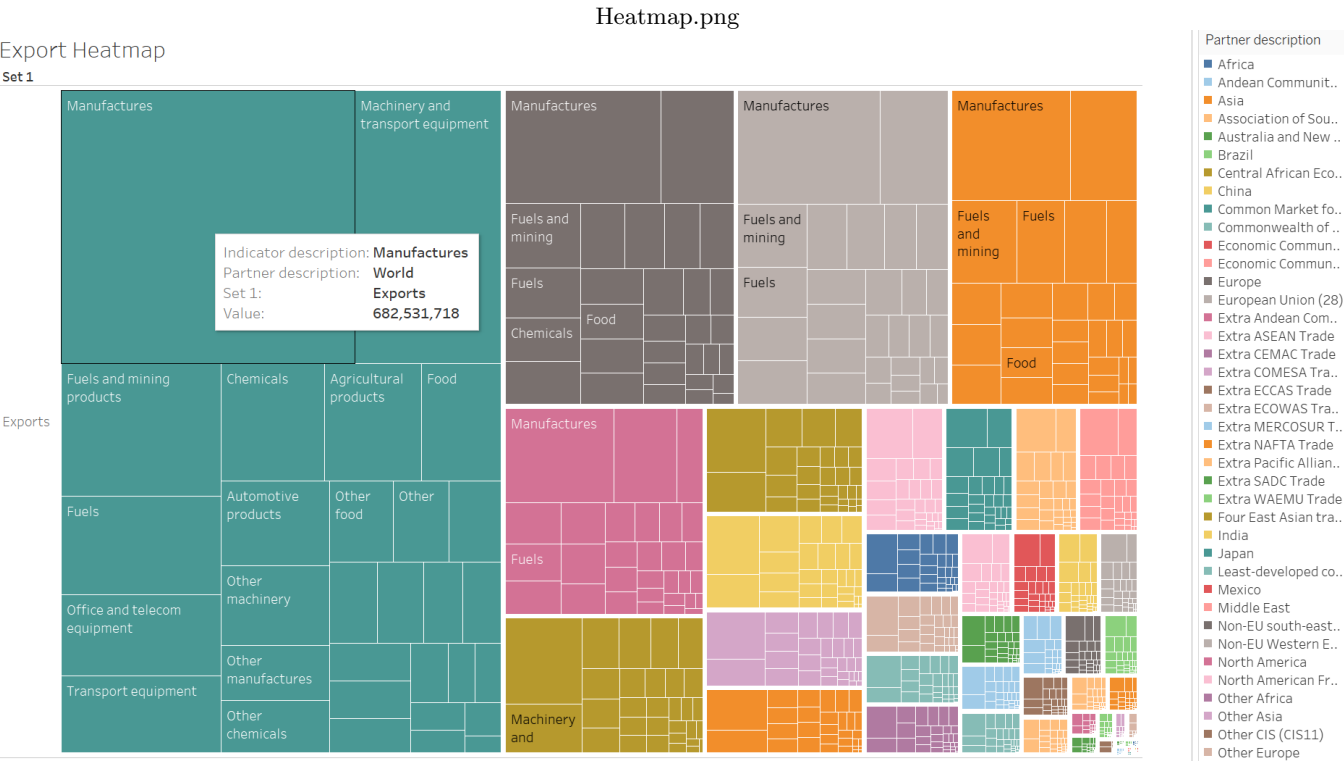FIG. 1. Data Quality summary

Heatmap.png



FIG. 2. Export Trade Heatmap

## D. Graph Analysis

**Why Split into Two Networks?**

Merchandise trade behaviour over time series provide lot of information of a well being of economy and global economy as a whole.For economist,Export and import pattern are of utmost important and provides critical in-formation across various trade categories.So research was divided in two categories one being generating and an-alyzing network graph of Export and other being Im-port.Moreover number of nodes and link are large so dataset was split logically and physically in Export and Import category and resulted in separate Network cre-ation of Import and Export.
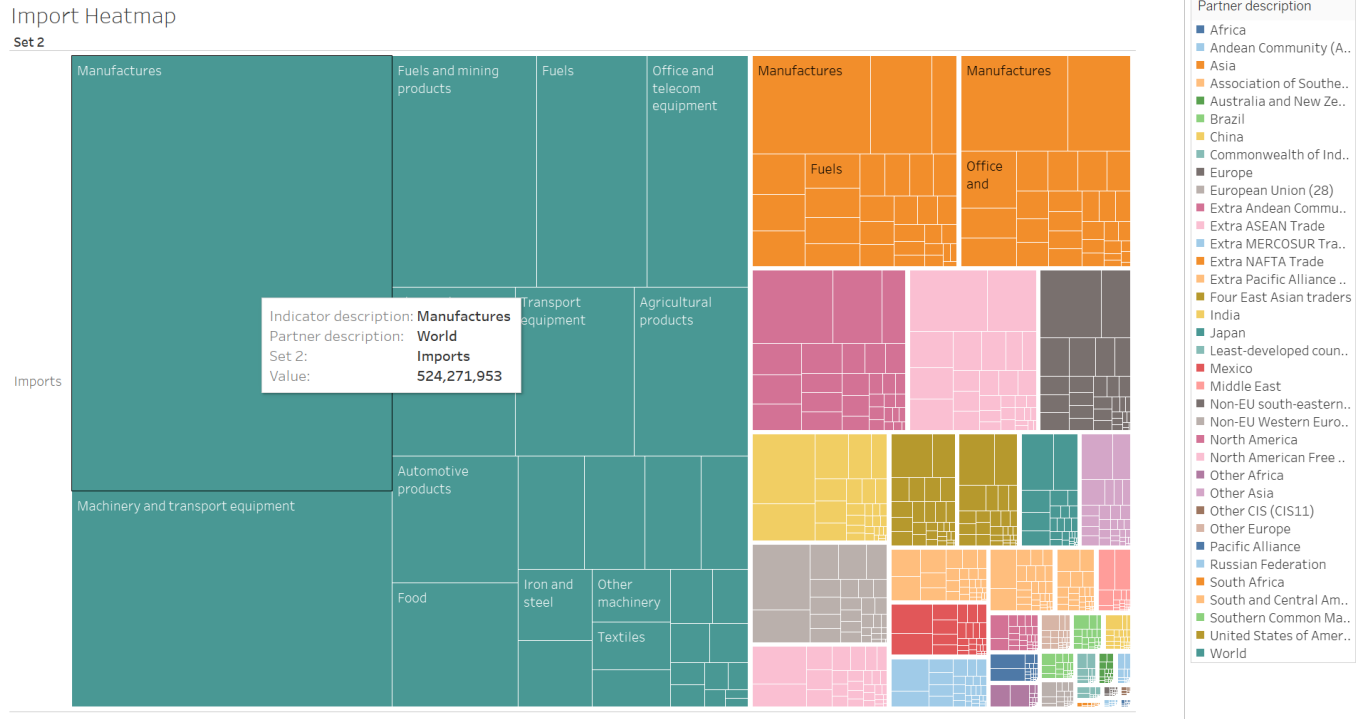
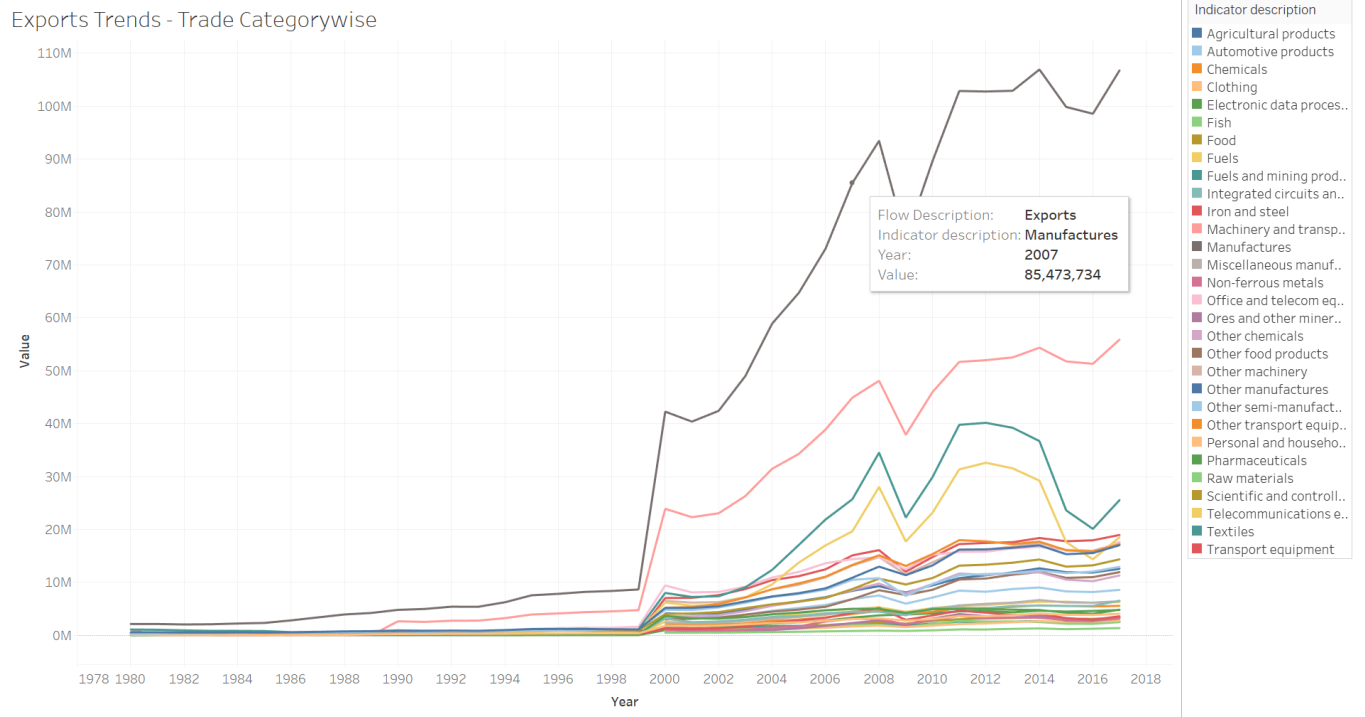Heatmap.png



FIG. 3. Import Trade Heatmap

Trends.png



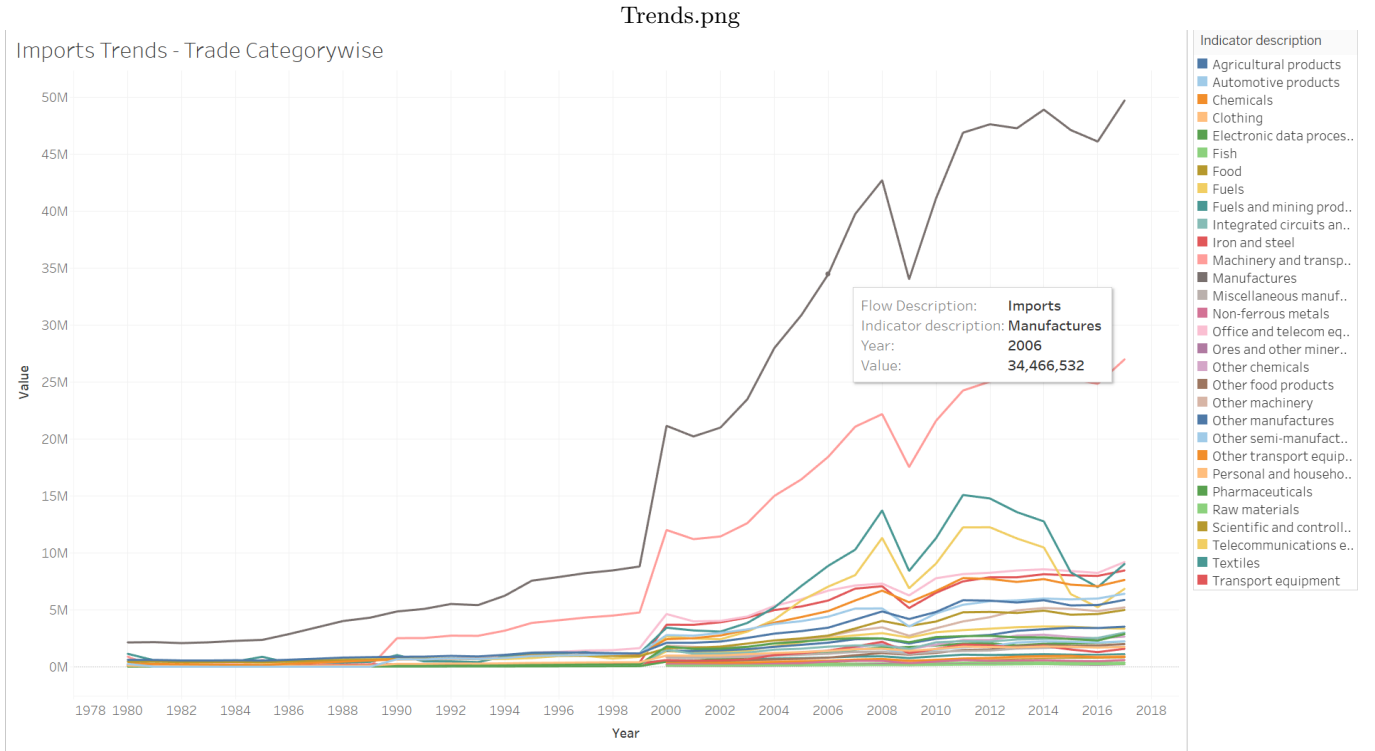FIG. 4. Export Trade Categorywise Trends

Trends.png



FIG. 5. Import Trade Category wise Trends

### 1. Export Network

**Network Creation**

For network definition in Neo4j, we need essentially information about nodes and links which, in our case, are countries or economic groups and trade transactions between them. We defined one CSV file which contains list of all countries in the dataset.It also includes countrycodes which will be used as uniqueness attribute in the network for node identification. For links creation, we filtered out all 'EXPORT' transactions from WTO dataset and then used it to create an 'EXPORT' link between Reporter Country and Partner country. Each link was assigned additional attributes of Type, Year and Value of that particular transaction. Type describes the category of trade like Manufacturing, Agriculture etc. Year describes the transaction year. It ranges between 1947 to 2018. Value is transaction amount expressed in million US dollars. This way we got a directed, multilink graph in neo4j. Our Export network contains 419 nodes and 617121 links. Following are the commands we used for network creation-

1-using periodic commit load csv with headers from "file:/merchandise_values_annual_dataset_country.csv" as line create (:Country countryname: line.Reporter$_d$escription, countrycode : line.Reporter$_c$ode)

2-using periodic commit load csv with headers from "file:/merchandise_values_annual_dataset_export.csv" AS input MATCH (from:Country countrycode:input.Reporter$_c$ode), (to : Countrycountrycode : input.Partner$_c$ode)

3-CREATE (from)-[:EXPORT type: input.Indicator$_d$escription, year : input.Year, value : input.Value]− > (to);

Queries

Once we have network in place, we can perform further analysis on it by various Cypher queries. We tried to see many different aspects of this datasets using queries, some of which are to see transactions between specific countries, all transactions for a given country etc. For example -

MATCH p= (c1:Countrycountrycode:'IN')-[r:EXPORT]-¿(c2:Countrycountrycode:'CN') Where r.year='2017' RETURN p

This query provides all transactions between India and China in the year of 2017. We can see that trade relationship between these countries is very deep and it involves export in many trade areas which is indicated by multiple links between the nodes. Each link represents a trade type with certain value of transaction. [FIG 7]

MATCH p= (c1:Countrycountrycode:'IN')-[r:EXPORT]-¿(c2:Country) Where r.year='2017' RETURN p limit 20 We can also see all trading partner for given country. We tried to see all trade partners

of India for year 2017. As we can see that India has Exported lots of goods to various countries across the world. This suggest that India is a economically very active in the year 2017. Each link represent trade with partner country with certain trade value.

### Algorithms

Next level of study was to get essence of network as whole. This can be done by applying various networking algorithms and get relevant details which lead to conclusions about economic activities. Many of such algorithms are available in Neo4j readily and can be deployed using Cypher queries. We decided to use Centrality algorithms and Community detection algorithms as they are more suitable for this dataset.

### Centralities

There are various centrality measure in Network Science depending on what they convey. Degree centrality describes how many nodes are connected to a given node, thus it tells us about most popular node of the network. Betweenness centrality describes the number of shortest paths passing through a particular node, which means that a node with higher betweenness centrality measure is very important for communication across network. PageRank centrality lets us know the probabilistic nature of connections of one node to other nodes in the network. Closeness centrality tells how close a given nodes is to the hub of the network.

Most of these centrality measures are usually referred in context of social network which are very complex in nature. There are connections of friends, friends of friends etc. However, in our case, WTO dataset is pretty straightforward and we don't have many connection beyond second order. So centrality measures like Betweenness centrality, Closeness centrality do not make much sense to study. But, Degreee centrality is very important for this network as it will tell us the most connected node which in turn is the 'Hub' of economic activity. So we calculated Degree centrality using below query-
CALL algo.degree("Country", "EXPORT", direction: incoming", writeProperty: "degree") MATCH (c:Country) WHERE exists(c.degree) RETURN c.countryname AS name, c.degree AS degree ORDER BY degree DESC LIMIT 15

We calculated Degree centrality of each node and stored it into an attribute called 'degree'. As we can see from the result, it shows all countries with most economic influence in descending order. 'world' is at number one, because most of the countries are reporting their total export and not the granular export with each country. It is followed by China, Russia, Europe and India. These results are very much in line with expectations based on economic reality. We can see USA at the bottom of list because it is mainly an import oriented country.[FIG 8]

### Community Detection

Another important aspect of network properties is communities within the network. Every network has unique community structure associated with it which is inherently developed due to the links between nodes. In case of WTO network, we can think of community as the group of countries which have frequent trade transactions among themselves as compared to other countries in the network. We know that , world had many economic groups such as NAFTA, ASEAN, EURO which consists of regional countries and they follow free trade policies. As a result of that, there is increased trade activity among them and it gives rise to economic community.

There are various ways to detect community structure in the network. Community detection algorithms typically evaluate links of a given node and rate them based on within and outside community. This action happens recursively to for community pattern. We decided to use one such algorithm readily available in Neo4j and its called 'Louvain Method'. It is based on free walk in the network process. It can be called by below cypher query and we stored detected community result as 'community' attribute-
CALL algo.louvain('Country', EXPORT', write:true, writeProperty:'community')YIELD nodes, communityCount, iterations, loadMillis, computeMillis, writeMillis; We can see that, Louvain algorithm has successfully detected 3 communities in the network.

### Network Visualization

Once we are done with network analysis using Graph algorithms, there is a need to visualize the results for better understanding of the outcomes. Neo4j browser does a limited job in providing basic visualization of nodes and links for a small portion of network. But when it comes to visualize network at large scale with many attributes, there are Javascript based options like in built Neo4j tool like Neovis.js,Popoto.js or embedded libraries like D3.js,vis.js etc. Biggest benefit of these libraries is that they are platform independent and can be embedded into a web based page for interactive visualization. We can also deploy third party standalone products with Neo4j drivers like Kineviz, Linkurious etc. However they are licensed products for enterprise applications.

We decided to go with neovis.js for our project because it is built in tool with Neo4j and its very simple to deploy. [10] It basically works within HTML page where we can define a Draw() function and define how we want to see our network. First it has Neo4j connection information. It connects to Bolt URL with username and password. This comes very handy if we have a Neo4j Sandbox spin up online. Than it has information about node and link attributes. We have chosen to represent node size by Degree centrality measure and Link thickness by Trade value. We also want communities in different colors

driven by community attribute in our network. We have defined countryname as tool tip which can show up when we hover over a node.

We can see that our network visualization clearly shows communities detected by Louvain algorithm. Roughly, we can see that it has clubbed together all countries which have chosen to give Export details only with 'world' in one community and countries which have chosen to give granular details of Export with each country in other community.Hence Blue community has only World at its center and Yellow community is much more interconnected.[FIG 6] [FIG 9]

### 2. Import Network

**Network Creation**

Neo4j follows approach of property graph model, where data is organized as nodes, relationships, and properties.Nodes are the entities in the graph whereas links provide directed,named,semantically-relevant connections between two node entities (data stored on the nodes or relationships).In this research nodes and links are economic groups or countries and and trade transactions within them.This helps to identify flow of merchandise trade between different entities.Master CSV file representing list of all countries from WTO dataset (which has countrycode as unique identifier for node),was generated. Next step was link creation which was achieved by filtering 'IMPORT' transactions from WTO dataset.This dataset was used to create 'IMPORT' link between Partner country and Reporter country. Link has additional attributes of Transaction's Year,Type(which classifies trade like Manufacturing etc) and dollar value(in million) of each merchandise transaction.Dataset spanned between 1947 to 2018. This results in directed,named,semantically-relevant multi-connections between nodes.The resultant Import network has 17 nodes and 33344 relationships. Neo4j network was generated using below set of commands:

1-using periodic commit load csv with headers from "file:/merchandise_values_annual_dataset_country.csv" as line create (:Country countryname: $line.Reporter_description, countrycode : line.Reporter_code)$ $2 - CREATE INDEX ON : Country(countrycode); 3 - CREATE(from) - [: IMPORT type : input.Indicator_description, year : input.Year, value : input.Value] -> (to);$

$4 - MATCH p = (c1 : Country) - [r : IMPORT] -> (c2 : Country) Where r.type =' Total merchandise' and r.year =' 2017' RETURN p limit 20$

**Algorithms**

Network generated can be analyzed by applying various networking algorithm which leads to drawing conclusions about global economic behavior. Neo4j provides out-of-box algorithms which can be utilized by using Cypher queries.In this research Centrality algorithms [11] and Community detection algorithms[12] were used to analyze the dataset.

**Centralities**

There are multiple Centrality algorithms[11] in the Neo4j Graph Algorithms library.There are various centrality algorithms which determine the importance of distinct nodes in a network like PageRank,ArticleRank,Betweenness Centrality,Closeness Centrality,Harmonic Centrality,Eigenvector Centrality and Degree Centrality.Each algorithm have a specific intended use,for example Betweenness Centrality is typically used where there is a need of detecting the amount of influence a node has over the flow of information in a graph ,specific use case being research the network flow in a package delivery process, or telecommunications network or identify influencer in legitimate, or criminal, organizations.Closeness centrality is another algorithm used where there is need of detecting nodes that are able to spread information very efficiently through a graph,real world use case being estimate the importance of words in a document etc. Based on Research paper use case where it is important to understand merchandise trade pattern i.e flow of transactions in and out of WTO member country;Degree centrality algorithm was apt algorithm.Theoretically,Degree centrality measures the number of incoming and outgoing relationships from a node.So by calculating Degree centrality it helps to identify top exporters or importers.This provides important as important input to economist,to focus on specific WTO countries economy and analyze specifics.

Degree centrality using below query-
CALL algo.degree("Country", "IMPORT", direction: incoming", writeProperty: "degree") MATCH (c:Country) WHERE exists(c.degree) RETURN c.countryname AS name, c.degree AS degree ORDER BY degree DESC LIMIT 15

Degree centrality of each node was calculated and stored it into an attribute called 'degree'.'world' is at number one,because most of the countries are reporting their total imports and not the granular import with each country.It is followed by some Asian economies. These results are very much in line with expectations based on economic reality.

**Community Detection**

Community detection algorithms [12] evaluate how a group is clustered and also indicator of its tendency to strengthen or break apart. This algorithm can also provide great deal of information about trading partners and evolution and possibly future partnerships.There are multiple community detection algorithm like Louvain,Label propagation, Connected components,Strongly connected components,Triangle counting,Balanced Triads.In case of WTO network, we can think of community

```html
<html>
    <head>
        <title>DataViz</title>
        <style type="text/css">
            #viz {
                width: 900px;
                height: 700px;
            }
        </style>
        <script src="https://rawgit.com/neo4j-contrib/neovis.js/master/dist/neovis.js"></script>
    </head>
    <script>
        function draw() {
            var config = {
                container_id: "viz",
                server_url: "bolt://localhost:7687",
                server_user: "neo4j",
                server_password: "changeidea",
                labels: {
                    "Country": {
                        caption: "countryname",
                        size: "degree",
                        community: "community"
                    }
                },
                relationships: {
                    "EXPORT": {
                        caption: false,
                        thickness: "value"
                    }
                },
                initial_cypher: "MATCH p= (c1:Country)-[r:EXPORT]->(c2:Country) Where r.type='Total merchandise' and r.year='2008' RETURN p"
            }

            var viz = new NeoVis.default(config);
            viz.render();
        }
    </script>
    <body onload="draw()">
        <div id="viz"></div>
    </body>
</html>
```

FIG. 6. Export Trade Communities HTML Settings



FIG. 7. Export Trade for India

```
$ MATCH (c:Country) WHERE exists(c.degree) RETURN c.countryname AS nam...
```

| name | degree |
|------|--------|
| "World" | 48337.0 |
| "China" | 235.0 |
| "Russian Federation" | 226.0 |
| "Europe" | 225.0 |
| "India" | 224.0 |
| "Four East Asian traders" | 223.0 |
| "Middle East" | 221.0 |
| "Japan" | 220.0 |
| "Australia and New Zealand" | 217.0 |
| "European Union (28)" | 214.0 |
| "Commonwealth of Independent States (CIS)- including associate and former member States" | 213.0 |
| "Other Africa" | 211.0 |
| "Other CIS (CIS11)" | 209.0 |
| "South Africa" | 208.0 |
| "United States of America" | 207.0 |

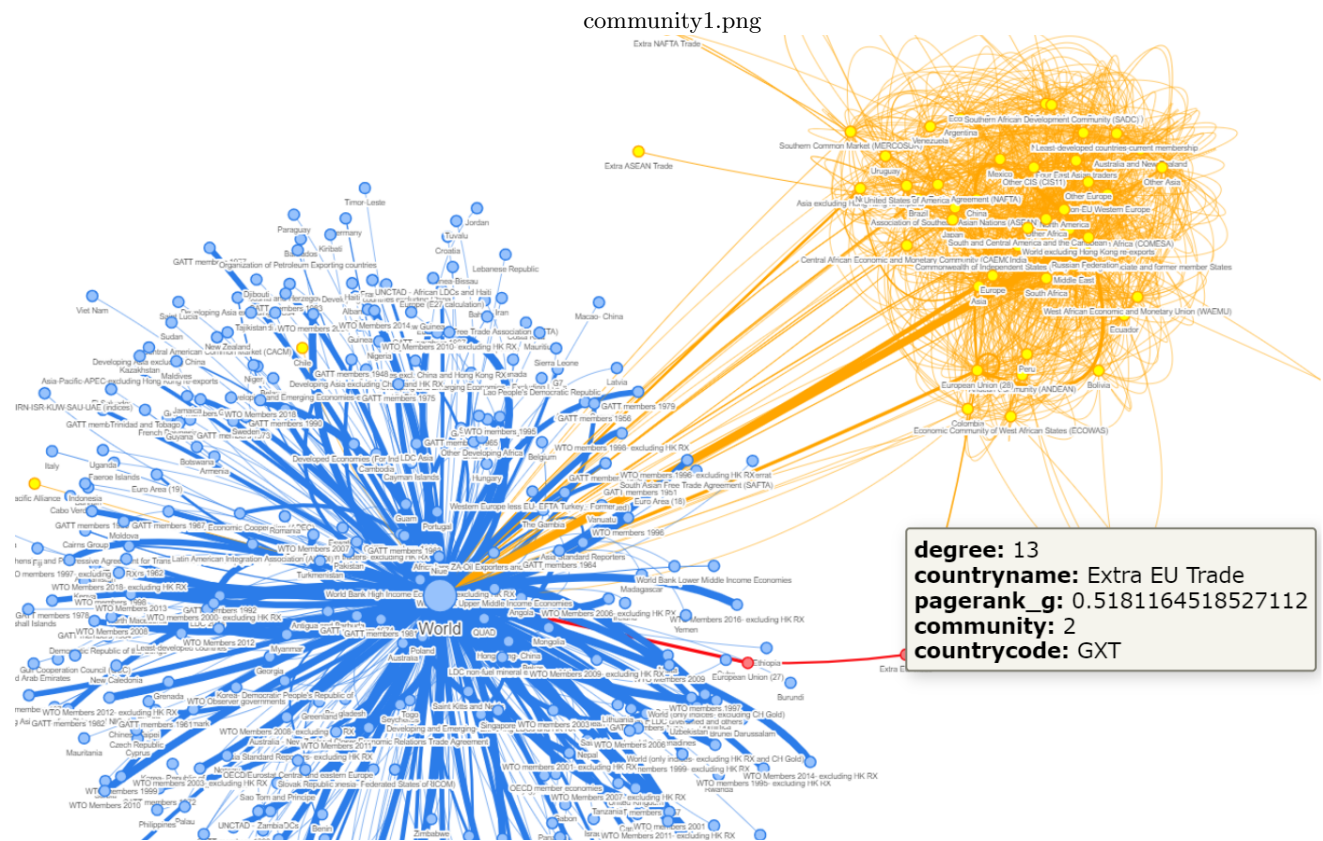FIG. 8. Export Trade Network Degree Centrality



FIG. 9. Export Trade Communities Visualization with NeoVis.js

FIG. 10. Import Network

as the group of countries which have frequent trade transactions among themselves as compared to other countries in the network. We scoped in 'Louvain Method' to evaluate community detection algorithm.It is based on random walk through the network process. It can be called by below cypher query and we stored detected community result as 'community' attribute-

CALL algo.louvain('Country', 'IMPORT', write:true, writeProperty:'community')YIELD nodes, community-Count, iterations, loadMillis, computeMillis, writeMillis; We observed that, Louvain algorithm has successfully detected 9 communities in the Import network.

## V. CONCLUSION

Many businesses are evolving beyond atomic intelligence, and making huge competitive gains by leveraging connected intelligence[13].Relational Databases has helped Business community for ages but with evolving connected intelligence use cases, graph databases like Neo4j is the best way ahead.As indicated in research paper,graph databases are quite efficient and helpful to perform analysis and identify trends.

## VI. TEAM CONTRIBUTION

For this project Deepak and Gautam worked as a team working coherently together. All data discussions, issue resolution were done collaboratively. However for sake of accountability, work was divided as below- Deepak handled exploratory data visualization with Tableau. He also worked on Export Trade graph analysis. Gautam worked on Intitial data analysis and Data preparation. He also handled Import Trade graph analysis. Presentation and Project report is done by common effort.

## ACKNOWLEDGMENTS

```
$ MATCH (c:Country) WHERE exists(c.degree) RETURN c.countryname AS name, c.degree AS degree ORDER BY d...
```

| name | degree |
| --- | --- |
| "World" | 54325.0 |
| "Non-EU Western Europe" | 26.0 |
| "Commonwealth of Independent States (CIS)- including associate and former member States" | 26.0 |
| "China" | 26.0 |
| "European Union (28)" | 25.0 |
| "Other CIS (CIS11)" | 24.0 |
| "India" | 24.0 |
| "South Africa" | 24.0 |
| "Other Africa" | 24.0 |
| "Australia and New Zealand" | 24.0 |
| "Europe" | 24.0 |
| "Four East Asian traders" | 24.0 |
| "Japan" | 24.0 |
| "Other Asia" | 24.0 |
| "South and Central America and the Caribbean" | 24.0 |

Started streaming 15 records after 33 ms and completed after 34 ms.

FIG. 11. Import Trade Degree Centrality

[1] W. T. Notes, Nicita, a. (2014). united nations conference on trade and development. united nations publication. geneva 10, switzerland: United nations conference on trade and development. (2019).

[2] W. T. Organizations, World trade organization . (2019). wto statistical data sets - metadata. (2019).

[3] W. T. Organizations, World trade organization. (2018). latest trends in world trade 2017-2018. world trade statistical review , 28-39. (2019).

[4] W. T. Organizations, Bulk download of trade datasets. retrieved from www.wto.org: (2019).

[5] W. T. Organizations, Statistics on merchandise trade. retrieved from www.wto.org: (2019).

[6] W. T. Organizations, Country analysis. retrieved from wits.worldbank.org: (2019).

[7] W. T. Organizations, Detailed product analysis. retrieved from wits.worldbank.org: (2019).

[8] W. T. Organizations, Merchandise trade balance with partners. retrieved from wits.worldbank.org: (2019).

[9] W. T. Organizations, World integrated trade solution. retrieved from wits.worldbank.org: (2019).

[10] W. Lyon, Graph visualization with neo4j using neovis.js (2018).

[11] neo4j, Chapter 5. centrality algorithms (2019).

[12] neo4j, Chapter 6. community detection algorithms (2019).

[13] J. W. Ian Robinson and E. Eifrem, Graph databasesbook to help solve common business problems (2019).

[14] A.-L. Barabsi, Network Science (Cambridge University Press, 2016).