# FacebookAnalytics

Exploratory Data Analysis on the dataset of Facebook's Altoona Data Center

Thamizhiniyan Pugazhenthi

22/12/2021

## INTRODUCTION

Since IT operations are so important for business continuity, they usually incorporate redundant or backup components and infrastructure for power, data communication connections, environmental controls (such as air conditioning and fire suppression), and other security systems. A huge data center is a large-scale activity that consumes as much energy as a small town[1]. Data center modernization and transformation improve performance and energy efficiency. Information security is also a concern, which is why a data center must provide a secure environment that reduces the risk of a security breach. As a result, a data center must maintain high standards for ensuring the integrity and functionality of the computer environment it hosts. The average age of a data center, according to industry research firm International Data Corporation (IDC) is nine years [2].

## DATASET DESCRIPTION

There are three clusters in the data center i.e., Cluster A is for Database, Cluster B is for Web servers and Cluster C is for Hadoop servers. Since we are working only with Cluster A and Cluster C, both of them have same set of columns. Each has 273 bz2 files in compressed format. Upon decompressing, each file is of tsv type with columns - timestamp, packet length, anonymized source IP, anonymized destination IP, anonymized source L4 Port, anonymized destination L4 Port, IP protocol, anonymized source hostprefix, anonymized destination hostprefix, anonymized source Rack, anonymized destination Rack, anonymized source Pod, anonymized destination Pod, intercluster and interdatacenter. All the data fields are anonymized for confidentiality. The timestamp column has values of type 'Unix timestamp'. The IP protocol column has three values - 6, 17 and 58. The value '6' points to Transmission control protocol(TCP), '17' points to User Datagram protocol (UDP) and '58' points to Internet Control message protocol(ICMP) for IPV6(Internet Protocol version 6).The last two columns plays major role to instigate analysis. The value '0' in intercluster denotes that the transmission of packets is within the cluster and '1' denotes that the transmission was between cluster. Same goes with the column 'datacenter'. Since facebook uses TCP segmentation offload, packet length can be larger than 64 kb because outbound packets are sampled in the kernel. It is not guaranteed to be a 1:1 mapping since the content has been hashed and get a subset of the hash value. The prefix of a hostname is called the host prefix. A computer named "web102.prn1.facebook.com" has the hostprefix "web." It's a very rudimentary classification of machine types. However, keep in mind that numerous programmes can operate on the same machine.

## EXPLORATORY DATA ANALYSIS

For this analysis, we took 100 samples from each cluster as subset and we will analyse the dataset with respect to each column and derive results from many perspectives especially with intercluster and interdatacenter. First, we need to cleanse the data since some of the fields have corrupted values. After cleansing, we need to

fix the values of timestamp column. The column has values of type 'unix timestamp' which need to converted to normal time value - BST value. Over the course of analysis, we will compare different characteristics of the data with respect to each column and analyse the data traffic in depth. Before step into the analysis, we will categorize the data into three types - Intra cluster and Intra datacenter, Inter cluster and Intra datacenter, Inter cluster and Inter datacenter.

**TIMESTAMP**

In database cluster, We can see the level of traffic in the data center over a period of time from Figure 1. As we can see that the level of traffic is higher in machines within cluster and datacenter followed by the machines between clusters within datacenter. We cannot see much traffic in machines between clusters and between datacenters. Moreover, we can see the intensity of internet traffic is higher between 18.00 hrs and 24.00 hrs in a day in all three classifications. From this plot, we can understand that the usage of facebook is higher in this timeframe than the rest of the day.

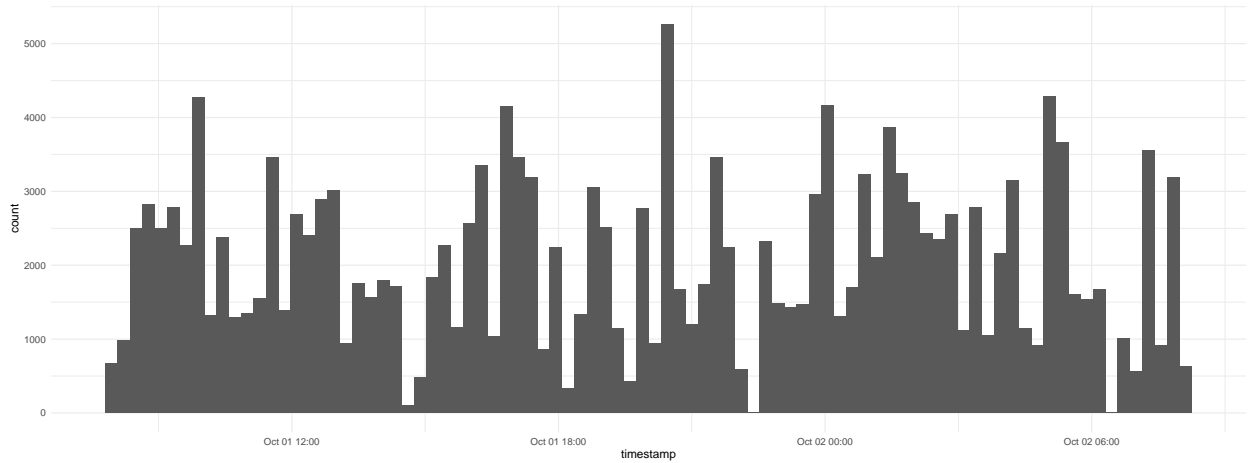Figure 1: Traffic characteristics - Database cluster - Overall

Figure 2: Traffic characteristics - Database cluster - Intra cluster - Intra datacenter

2

From Figure 2, Figure 3 and Figure 4, we can see the level of traffic in all three categorizations with separate plot for a clear picture.
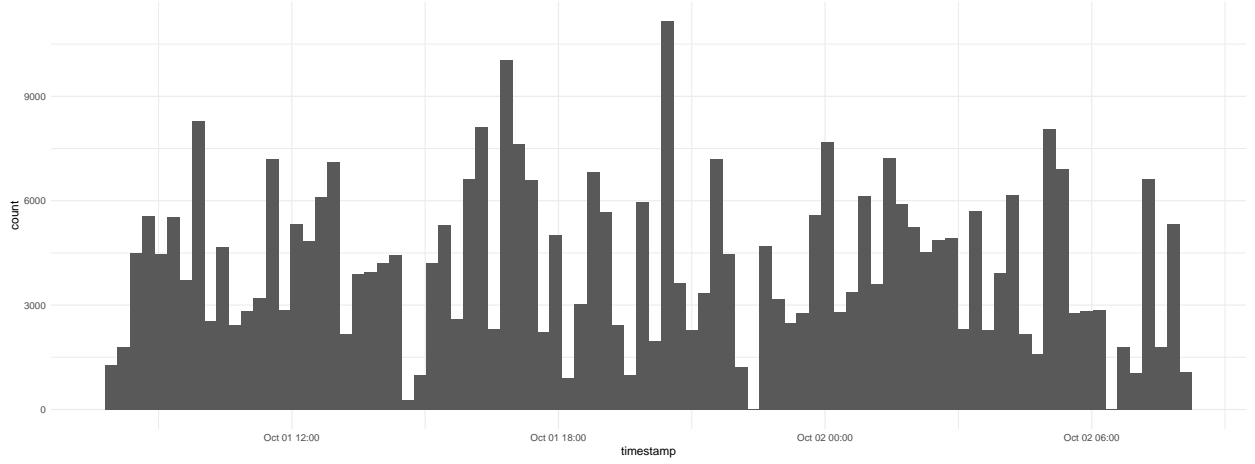


Figure 3: Traffic characteristics - Database cluster - Inter cluster - Intra datacenter
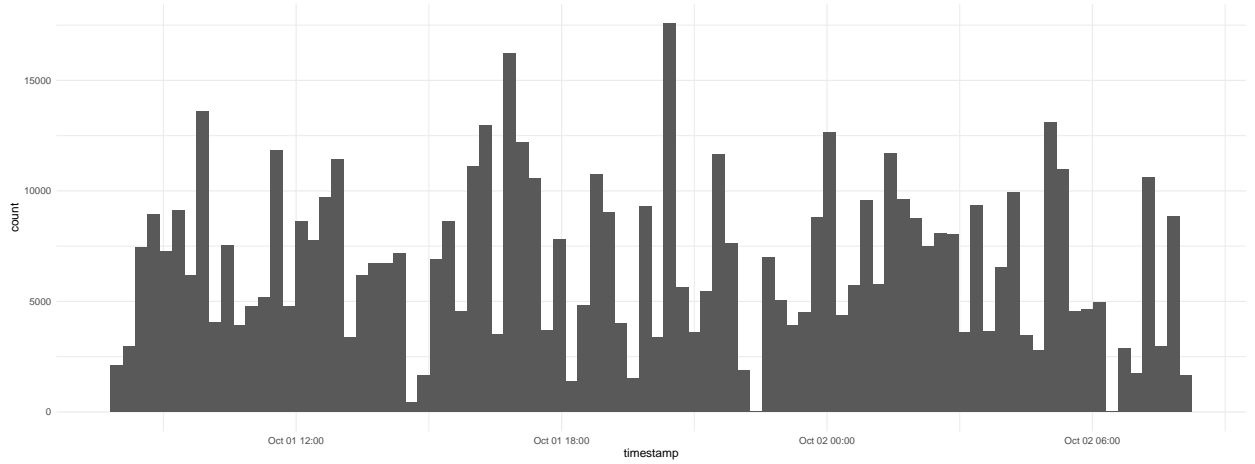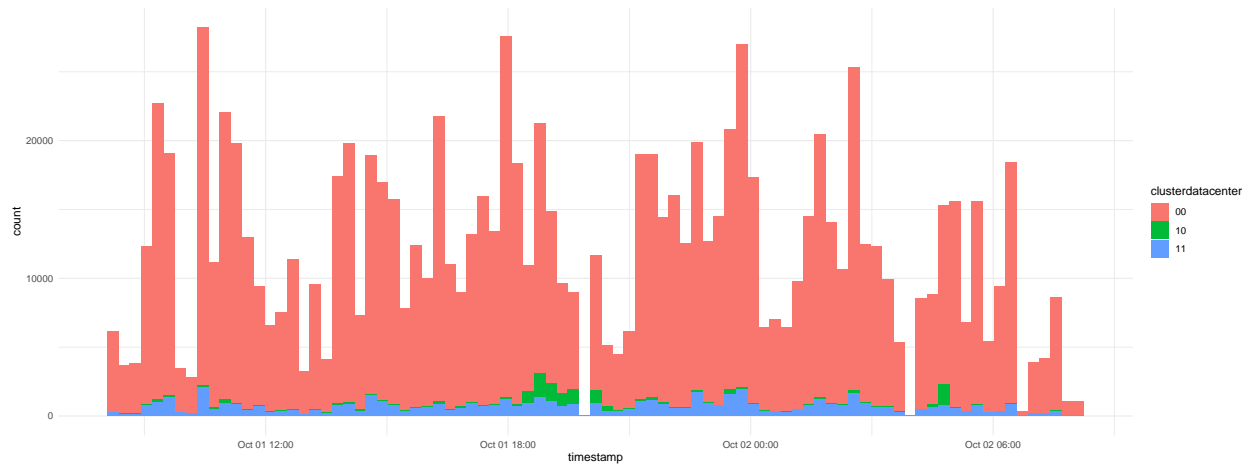


Figure 4: Traffic characteristics - Database cluster - Inter cluster - Inter datacenter

In Hadoop cluster, We can see the level of traffic in the data center over a period of time from Figure 5. As we can see that the level of traffic is very much higher in machines within cluster and datacenter followed by the machines between between clusters and between datacenters. We cannot see much traffic in machines between clusters within datacenter.

From Figure 6, Figure 7 and Figure 8, we can see the level of traffic in all three categorizations with separate plot for a clear picture.

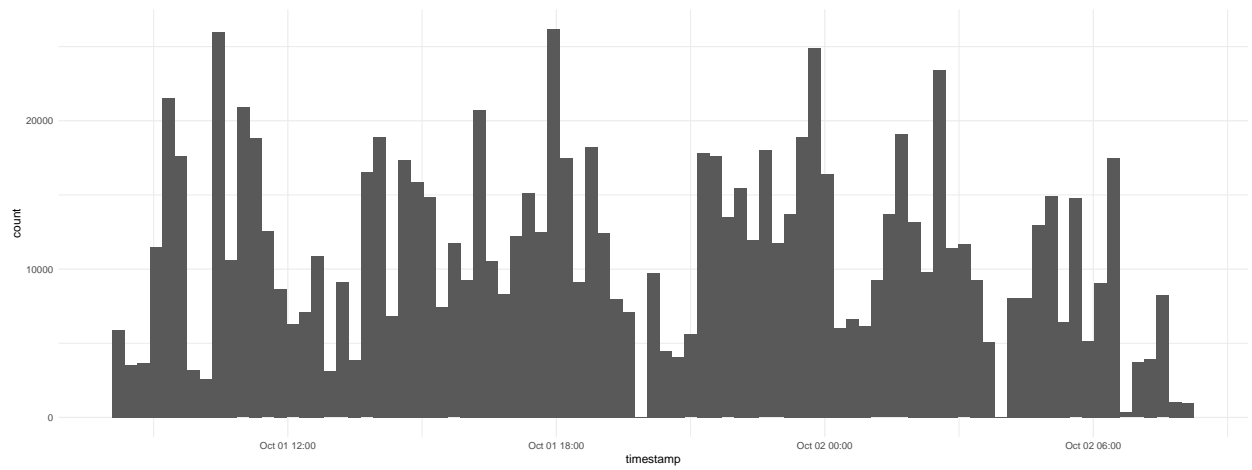Figure 5: Traffic characteristics - Hadoop cluster - Overall



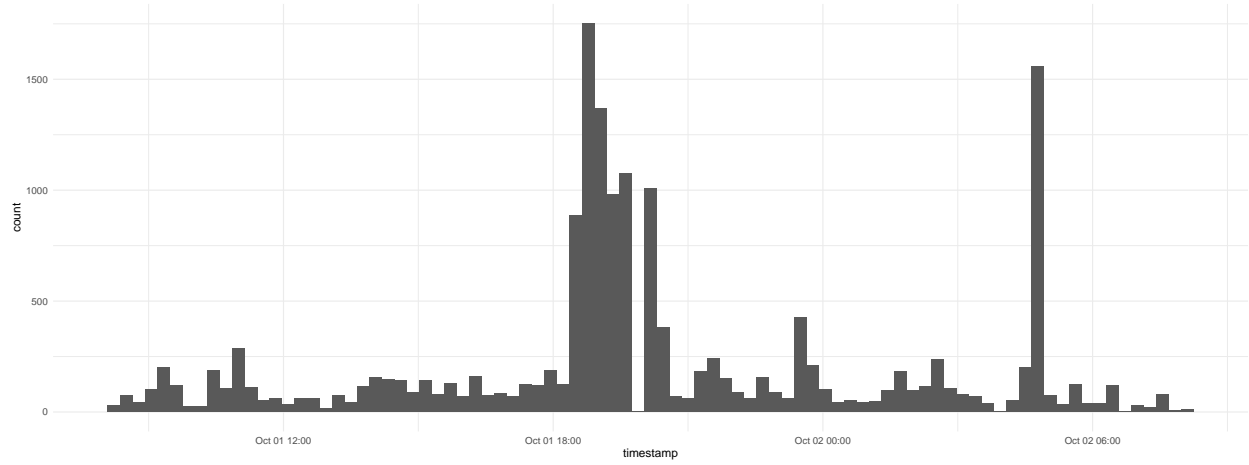Figure 6: Traffic characteristics - Hadoop cluster - Intra cluster - Intra datacenter

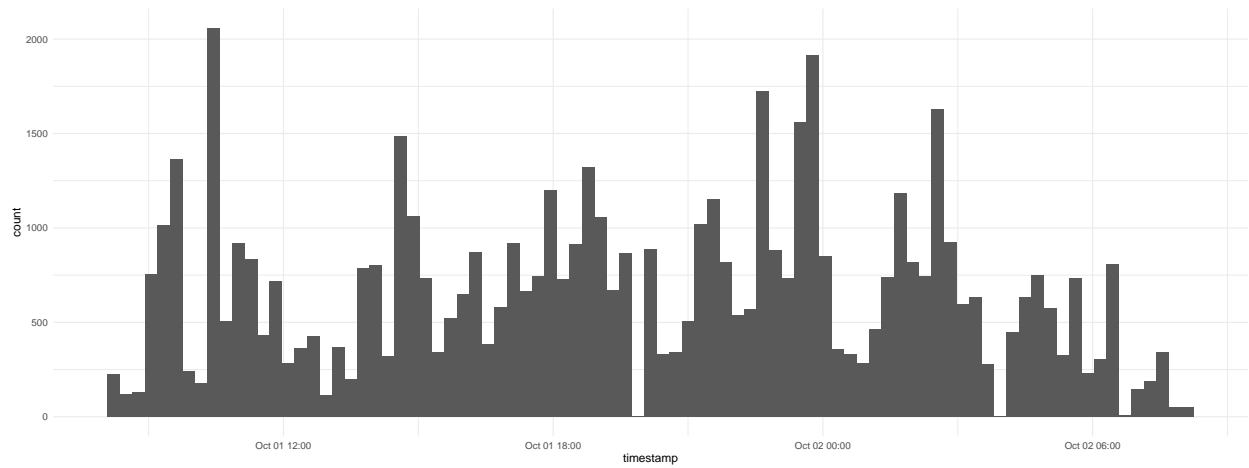Figure 7: Traffic characteristics - Hadoop cluster - Inter cluster - Intra datacenter



Figure 8: Traffic characteristics - Hadoop cluster - Inter cluster - Inter datacenter

In this second level of analysis, we will compare the level of traffic between database and hadoop cluster. It is evident that the traffic is higher in the database cluster than in hadoop cluster. We can identify this with the Figure 9.

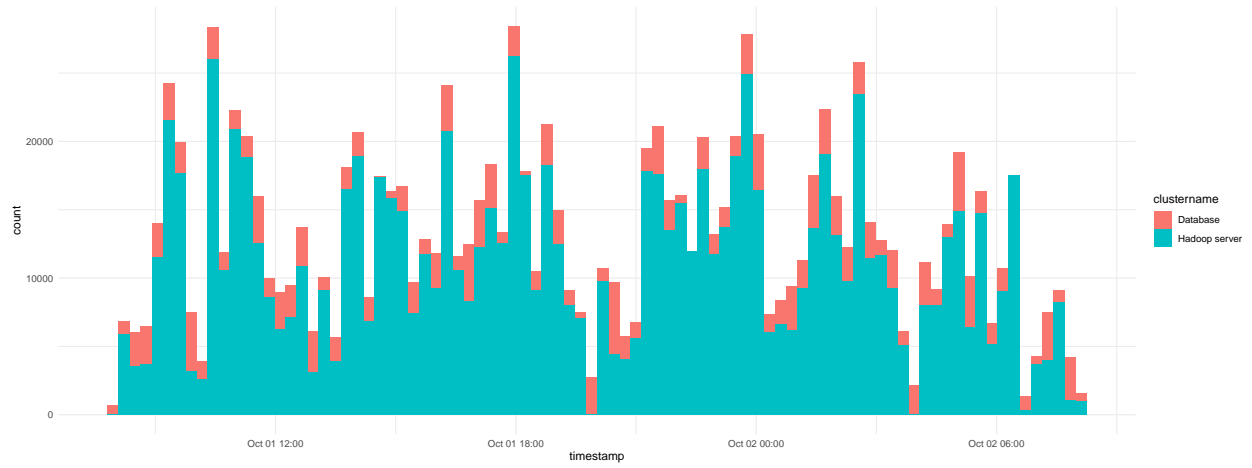Figure 9: Traffic characteristics - Database and Hadoop server - Overall



Figure 10: Traffic characteristics - Database and Hadoop server - Intra cluster and Intra datacenter
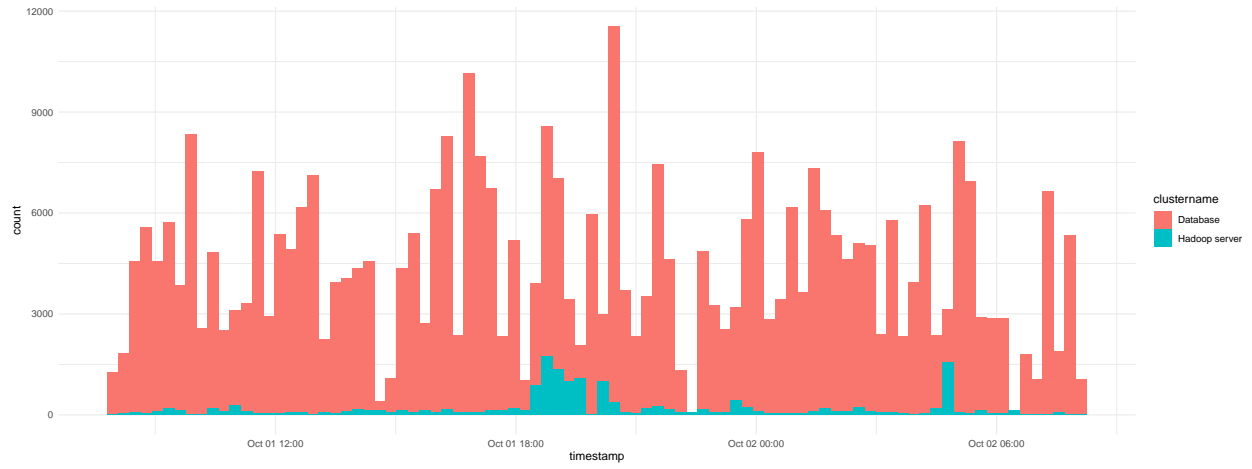
6

Figure 11: Traffic characteristics - Database and Hadoop server - Inter cluster and Intra datacenter



Figure 12: Traffic characteristics - Database and Hadoop server - Inter cluster and Inter datacenter

From Figure 10, Figure 11 and Figure 12, we can see the comparison between clusters in all three classifications with separate plot for each.

## IP ADDRESS

We will analyse the data with respect to the source and destination IP address. From Figure 13, we can see the top 10 combination of source and destination IP address which contributes more to the traffic between the machines overall without any classifications in database cluster.
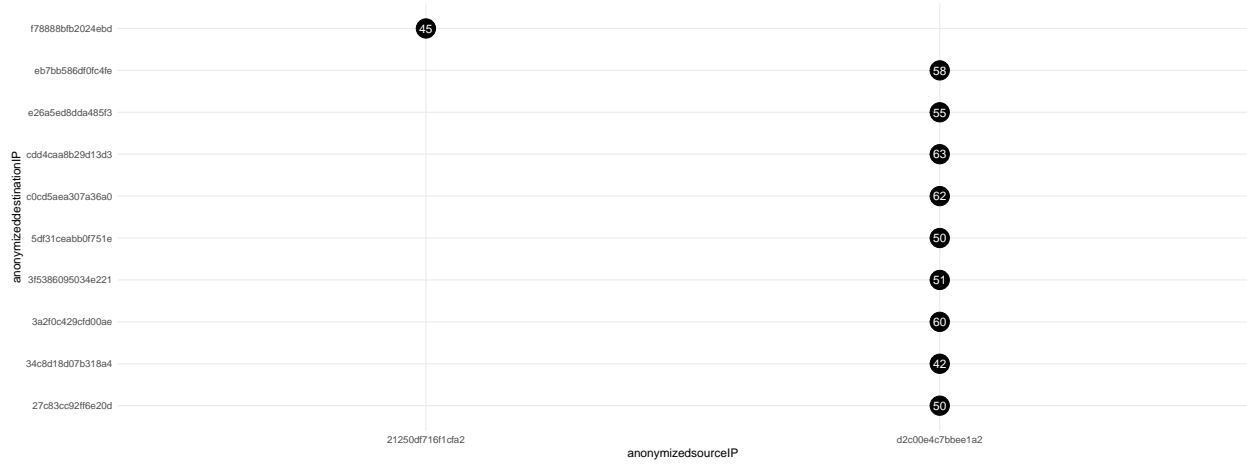


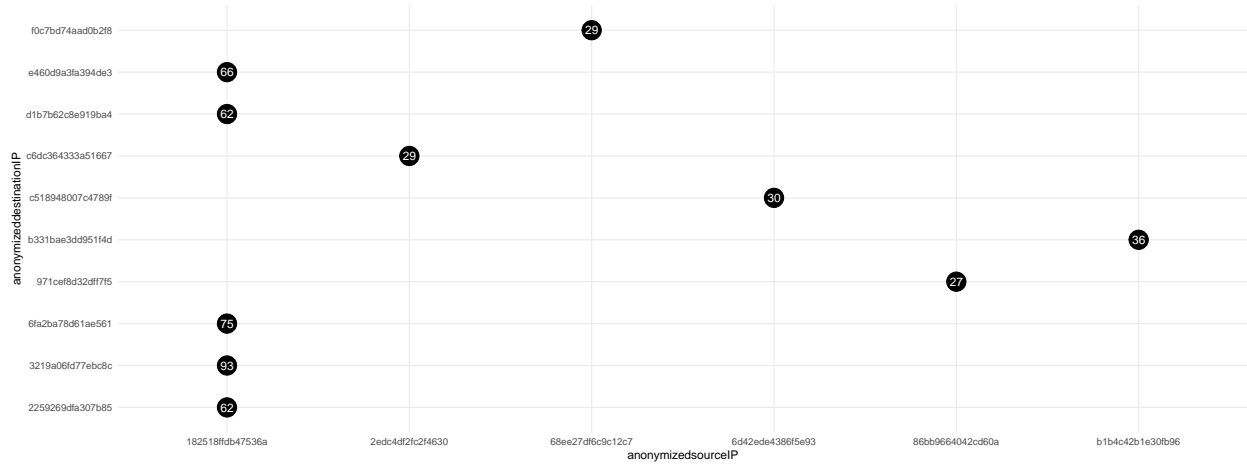Figure 13: Traffic between source and destination IP - Database Cluster



Figure 14: Traffic between source and destination IP - Hadoop Cluster

From Figure 14, we can see top 10 combinations of source and destination IP address which contributes more traffic in Hadoop cluster. To explain this plot, X axis denotes source IP address and Y axis denotes destination IP address with the number in the circle in align with the axis denotes the level of traffic between these addresses.

## RACK AND POD

In Altoona data center, they divided the network up into small identical pieces – server pods – instead of huge devices and clusters, and created consistent high-performance connectivity between all pods in the data

center. Each pod is served by a set of four fabric switches, which maintain the benefits of our present 3+1 four-post architecture for server rack TOR linkups while also allowing for future expansion. Each TOR presently has four 40G up links, allowing a rack of 10G-connected computers to access 160G of total bandwidth. Each pod has only 48 server racks, and this form factor is consistent across all pods. Another major difference is the manner in which the pods are linked to form a data center network.They constructed four different "planes" of spine switches to achieve building-wide connection, each scalable up to 48 independent devices within a plane. Within its local plane, each pod's fabric switch links to each spine switch. Pods and planes combine to produce a modular network topology capable of supporting hundreds of thousands of 10G-connected servers, growing to multi-petabit bisection bandwidth, and providing non-oversubscribed rack-to-rack performance across their data center facilities [3].

From Figure 15 and Figure 16, we can see the combination of source and destination Rack address which contribute more to the traffic between machines in Database and Hadoop clusters respectively. If we see the plot, we can identify that four source and three destination racks contributes more traffic to the entire system in Figure 16. This constitutes top 10 of the entire combinations. In Figure 17, we can see that it is 1-1 mapping in terms of traffic, each source and destination combination does not collide with any other combination within the cluster. Even this figure displays the top 10 combinations of the entire list.
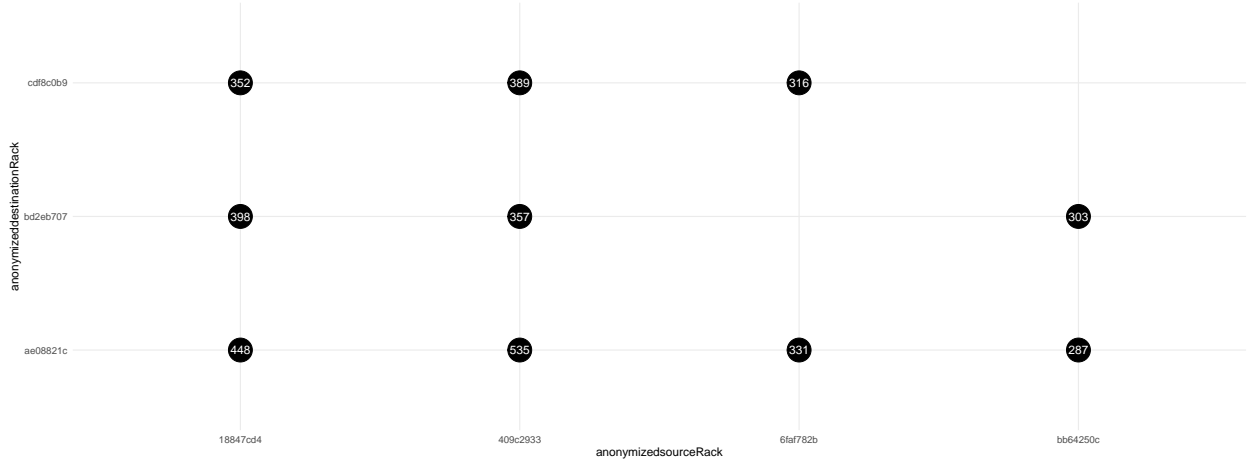


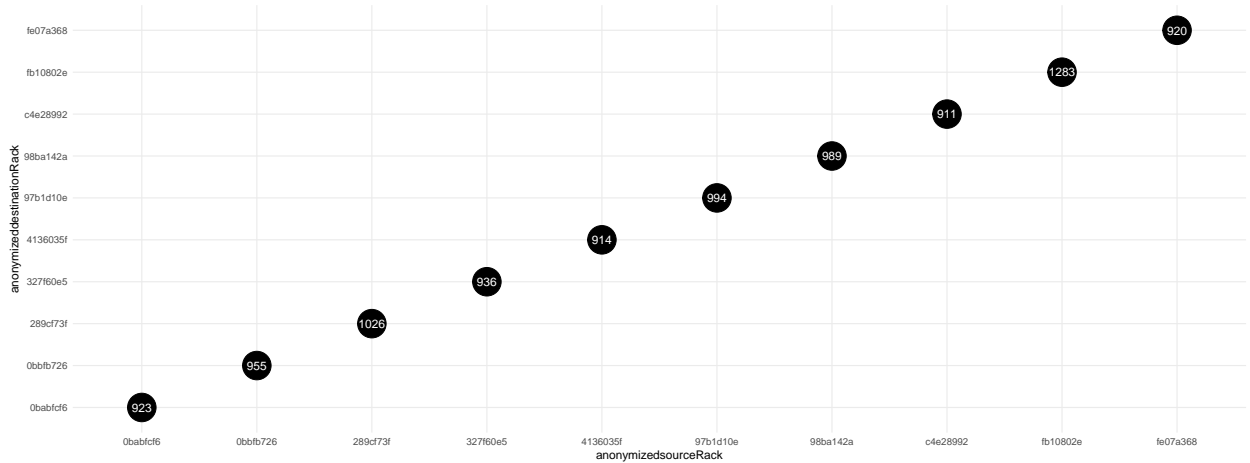Figure 15: Traffic between source and destination Rack - Database Cluster



Figure 16: Traffic between source and destination Rack - Hadoop Cluster

In Figure 17 and Figure 18, we can identify the combination of source and destination pairs of pod which
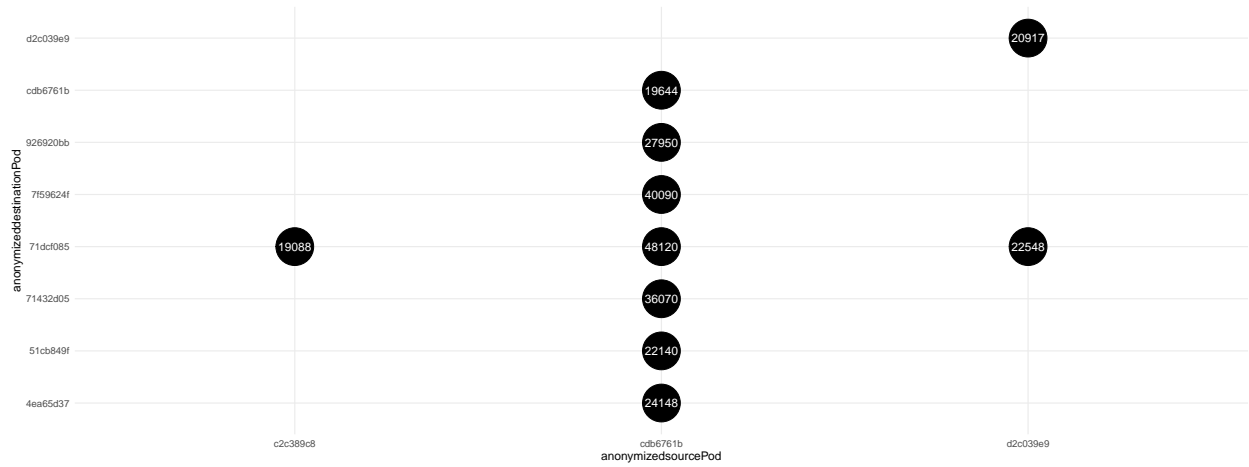
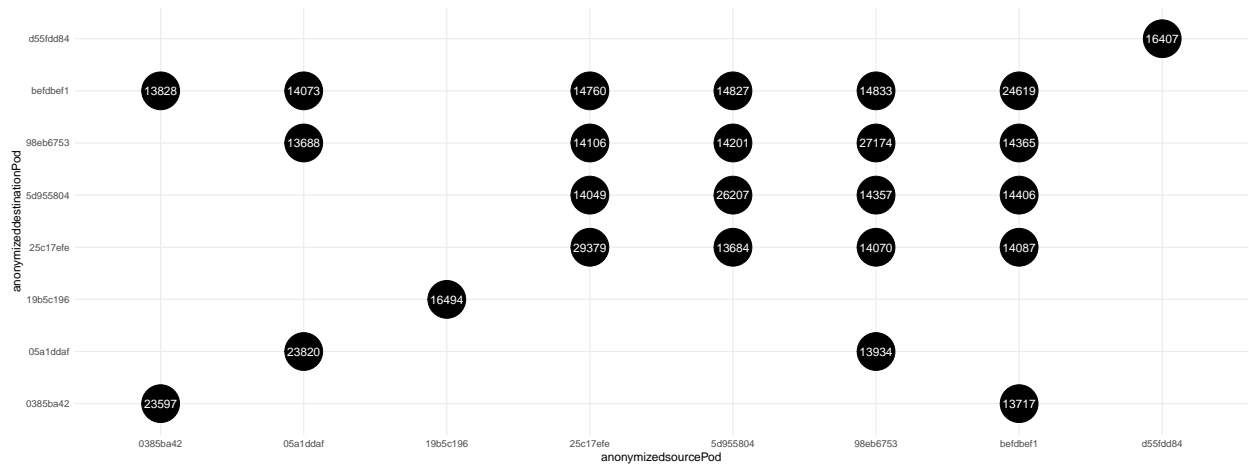Figure 17: Traffic between source and destination Pod - Database Cluster



Figure 18: Traffic between source and destination Pod - Hadoop Cluster

contributes more to the traffic intensity. From Figure 17, we can see that only three source pod contributes more to the traffic but in Figure 18, there is an equal distribution of source and destination combinations in Hadoop cluster.

## HOST PREFIX

From Figure 19 and Figure 20, we can see the top 10 combination of source and destination host prefix in database and Hadoop cluster respectively. As we can see three source host prefix and two source host prefix contributes more to the traffic in database and Hadoop cluster. The intensity of the traffic keeps increasing gradually from the plots of IP address to host prefix.
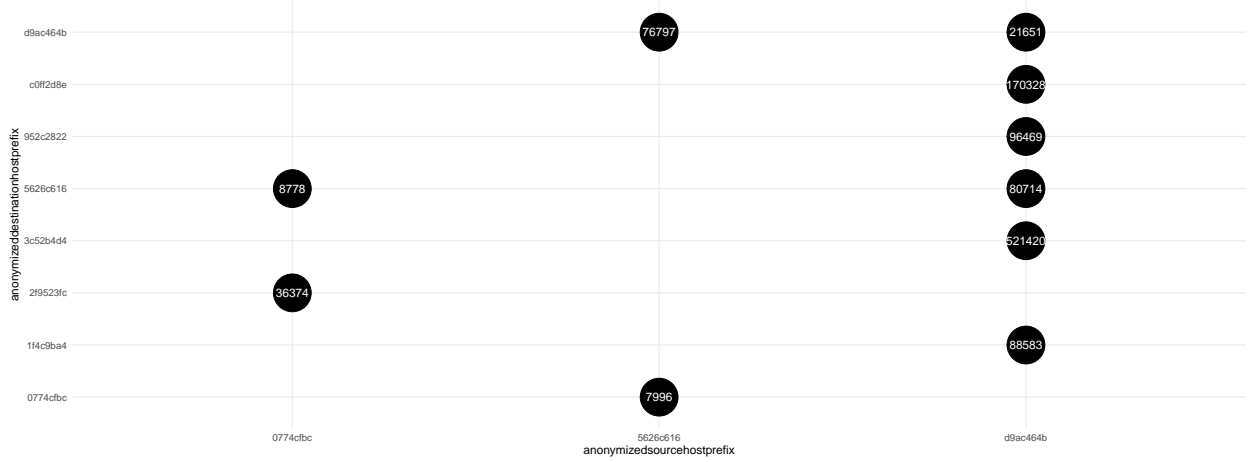


Figure 19: Traffic between source and destination hostprefix - Database Cluster



Figure 20: Traffic between source and destination hostprefix - Hadoop Cluster

## L4 PORT

From Figure 21 and Figure 22, we can see the top 10 combination of source and destination host prefix in database and Hadoop cluster respectively. As we can see three source host prefix and two source host prefix contributes more to the traffic in database and Hadoop cluster.
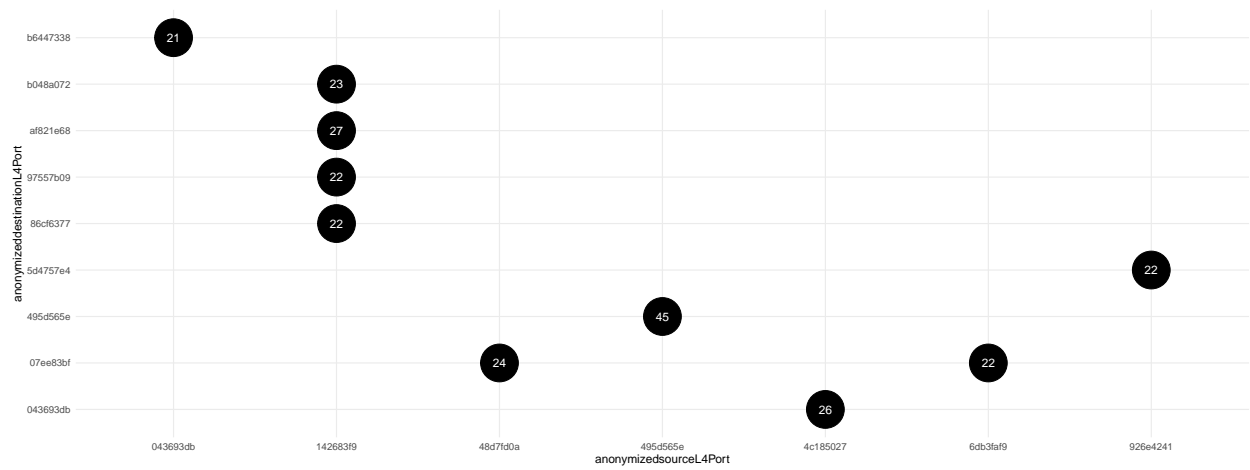
11

Figure 21: Traffic between source and destination L4 port - Database Cluster
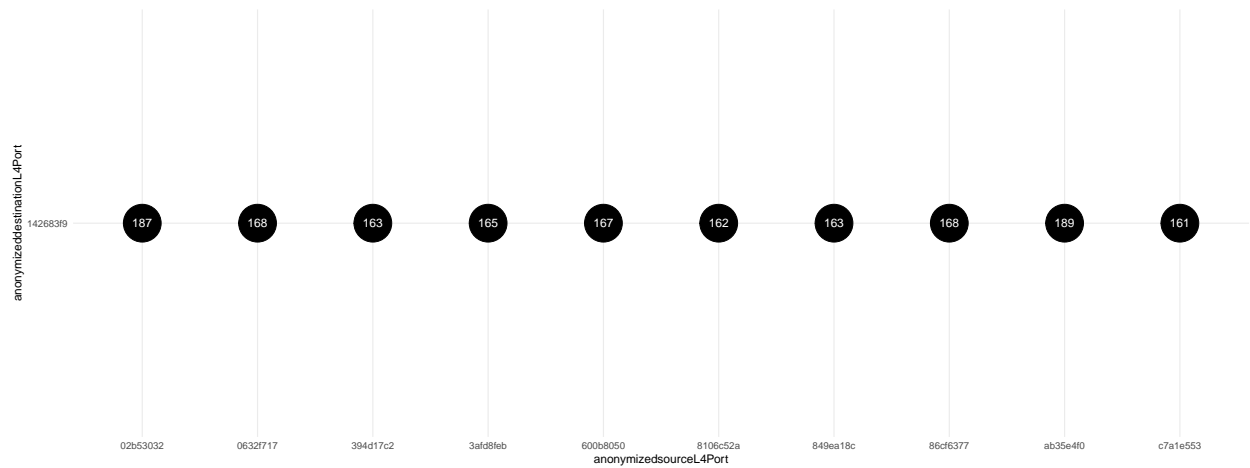


Figure 22: Traffic between source and destination L4 port - Hadoop Cluster

## IP PROTOCOL

All these request and reply happen within three IP protocol number - 6, 17 and 58. 6 is for Transmission Control protocol (TCP), 17 is for User Datagram protocol (UDP) and 58 is for Internet Control message protocol(ICMP) for IPV6(Internet Protocol version 6) (ICMP - IPV6). We can see that the lot of traffic intensity is from TCP which is way higher than the rest of the protocols. We will categorize this with respect to hostprefix and will analyse in depth in the following figures.

```
##   IPprotocol       n
## 1          6 1137919
## 2         17     760
## 3         58      45


##   IPprotocol       n
## 1          6 1036337
## 2         17    5150
## 3         58      50
```

From Figure 23 to Figure 34, we will see the traffic among different characteristics of dataset in depth. The top level will be Database and Hadoop cluster, the second level will be the classification we made before and final level will be the protocol number.
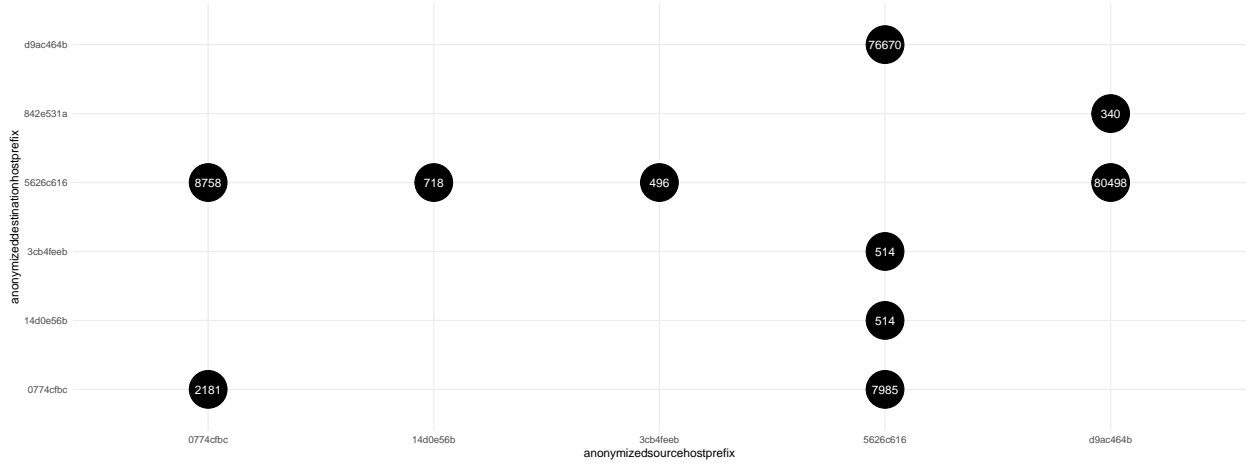


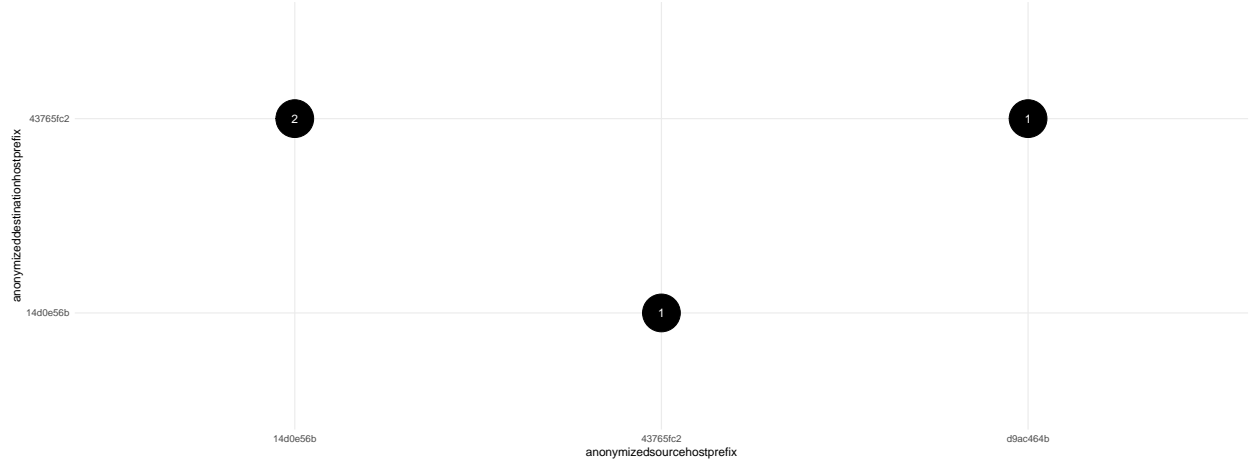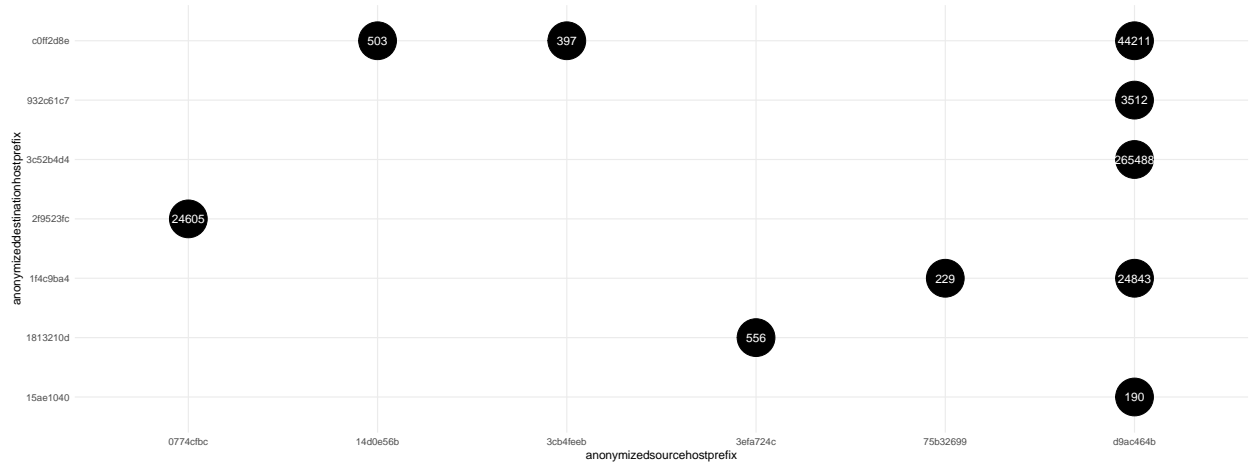Figure 23: Source and destination hostprefix - Database Cluster - Intracluster - Intradatacenter - TCP
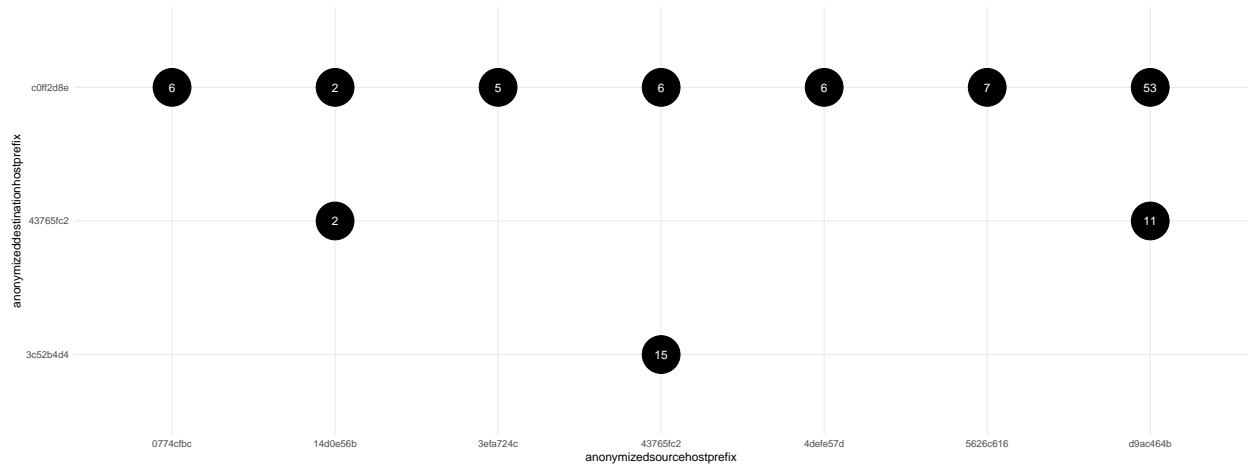
Figure 24: Source and destination hostprefix - Database Cluster - Intracluster - Intradatacenter - UDP



Figure 25: Source and destination hostprefix - Database Cluster - Intercluster - Intradatacenter - TCP



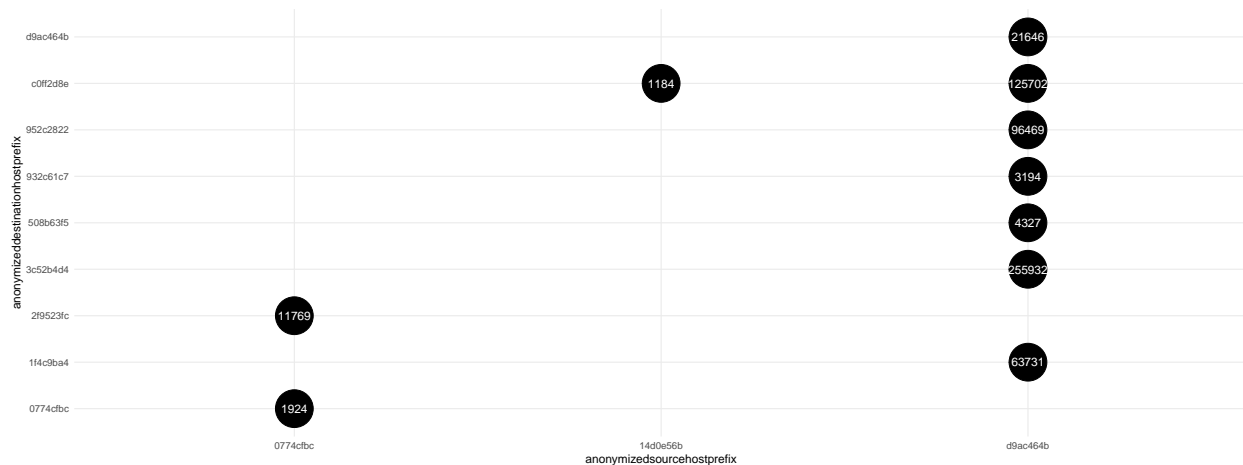Figure 26: Source and destination hostprefix - Database Cluster - Intercluster - Intradatacenter - UDP

Figure 27: Source and destination hostprefix - Database Cluster - Intercluster - Interdatacenter - TCP
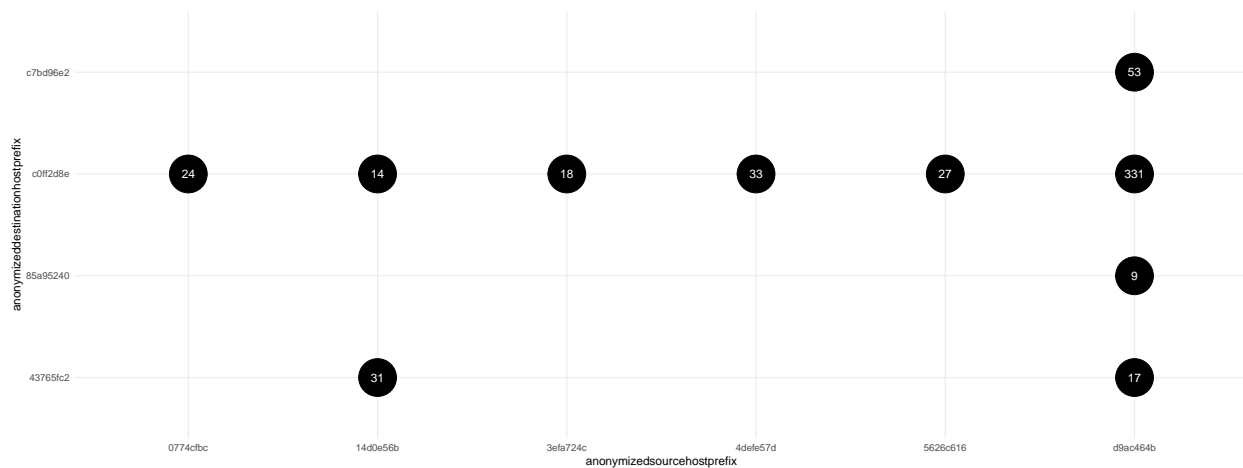


Figure 28: Source and destination hostprefix - Database Cluster - Intercluster - Interdatacenter - UDP
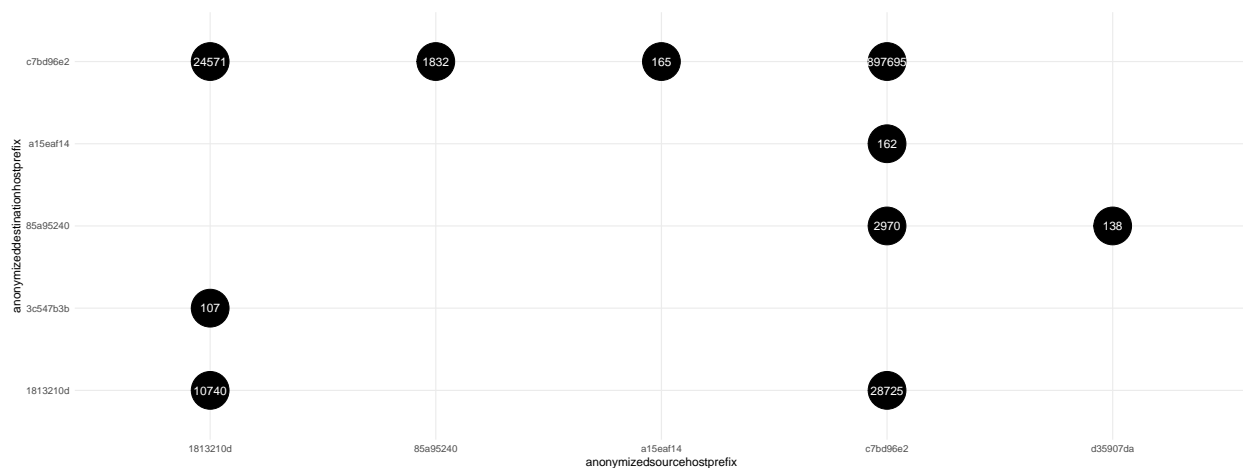


Figure 29: Source and destination hostprefix - Hadoop Cluster - Intracluster - Intradatacenter - TCP
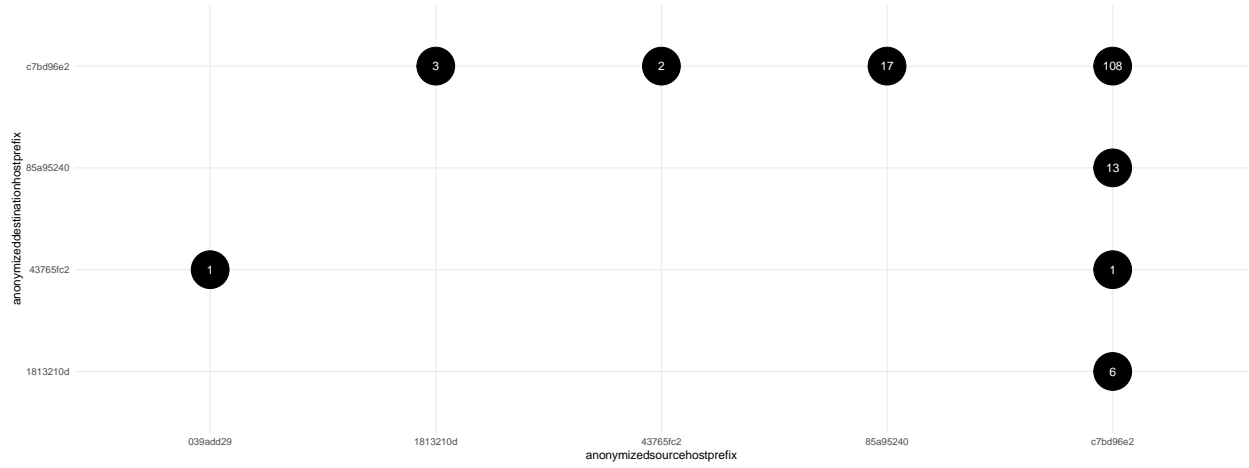
Figure 30: Source and destination hostprefix - Hadoop Cluster - Intracluster - Intradatacenter - UDP



Figure 31: Source and destination hostprefix - Hadoop Cluster - Intercluster - Intradatacenter - TCP



Figure 32: Source and destination hostprefix - Hadoop Cluster - Intercluster - Intradatacenter - UDP
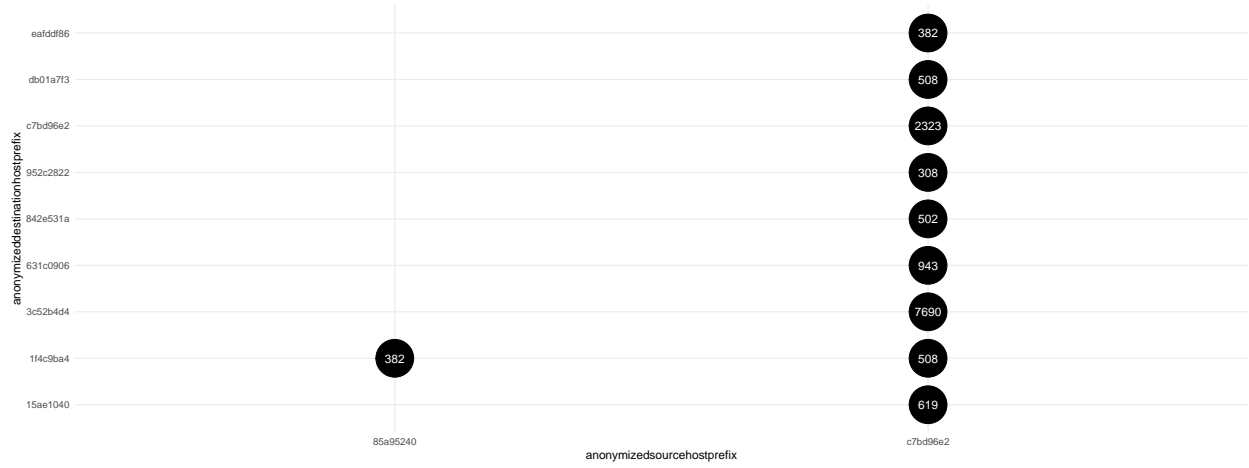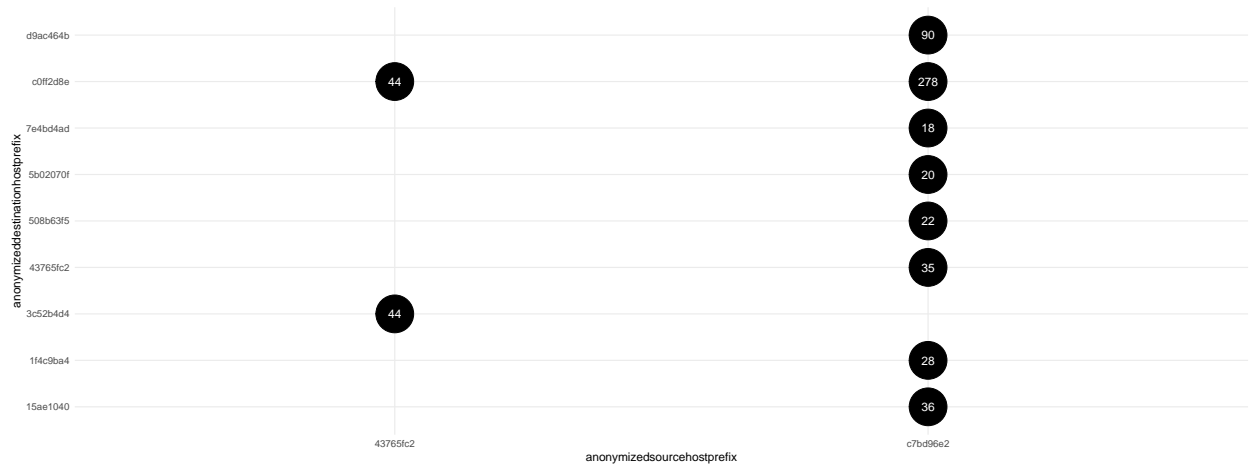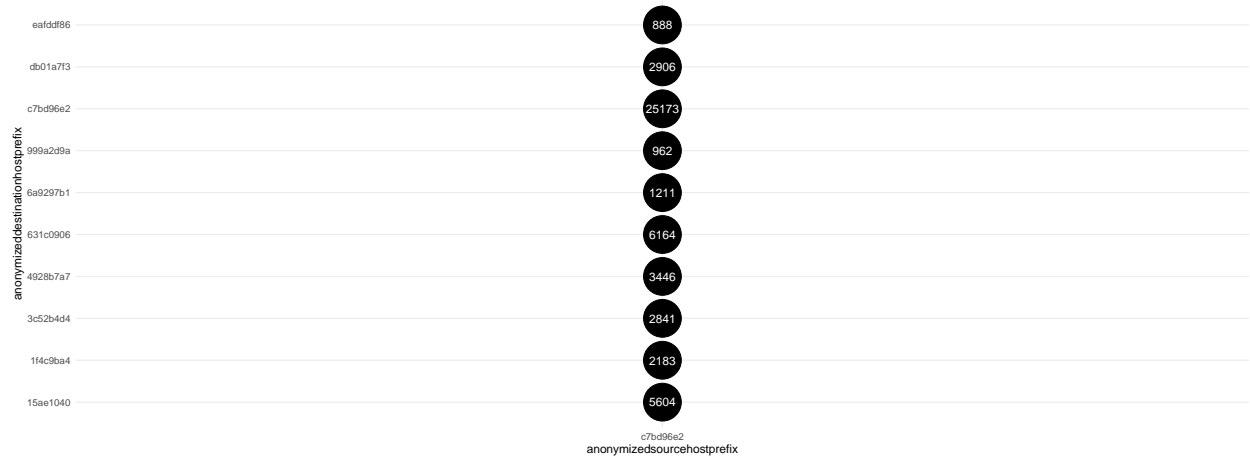
Figure 33: Source and destination hostprefix - Hadoop Cluster - Intercluster - Interdatacenter - TCP



Figure 34: Source and destination hostprefix - Hadoop Cluster - Intercluster - Interdatacenter - UDP

From these plots, we can identify the source and destination hostprefix addresses which contribute more to the machine to machine traffic in the Altoona data center. Since the plots are self explanatory, we can derive serious conclusions from the figures itself.

## CONCLUSION

From the analysis, we did analyse the dataset from various perspectives and measurements. From IP address to hostprefix, we can see the level of intensity in traffic increases gradually. The people at Altoona data center should look at this carefully and do the necessary steps to maintain load balancing to reduce the machine to machine traffic. Since the machine to machine will be higher in Altoona data center compared to user to machine traffic[4], they need to find out new strategy to tackle the traffic. Also they need to spend their time on the effect of facebook data center to the environment since it needs tons and tons of electric power to run which will sure have negative impact on the environment. But the facebook concentrates on renewable energy, we can say that they are heading in the right direction in terms of sustainability[5].

## REFERENCES

1.https://en.wikipedia.org/wiki/Data_center

2.http://www.mspmentor.net/2011/08/17/hp-updates-data-transformation-solutions/

3.https://engineering.fb.com/2014/11/14/production-engineering/introducing-data-center-fabric-the-next-generation-facebook-data-center-network/

4.https://www.cs.unc.edu/xcms/wpfiles/50th-symp/Moorthy.pdf

5.https://datacenters.fb.com/wp-content/uploads/2021/10/Altoona-Data-Center.pdf