

# Exploring IP Address Characteristics to Identify Router Deployments

Yue Xin

# IP Address Alias Resolution

## Why?

- Information is transmitted at the IP level, while AS/ ISP plan on Router level.

## Problems Szenario:

- Routing politics Differences
- Unstable probing Samples
- IPv4 and IPv6 barrier

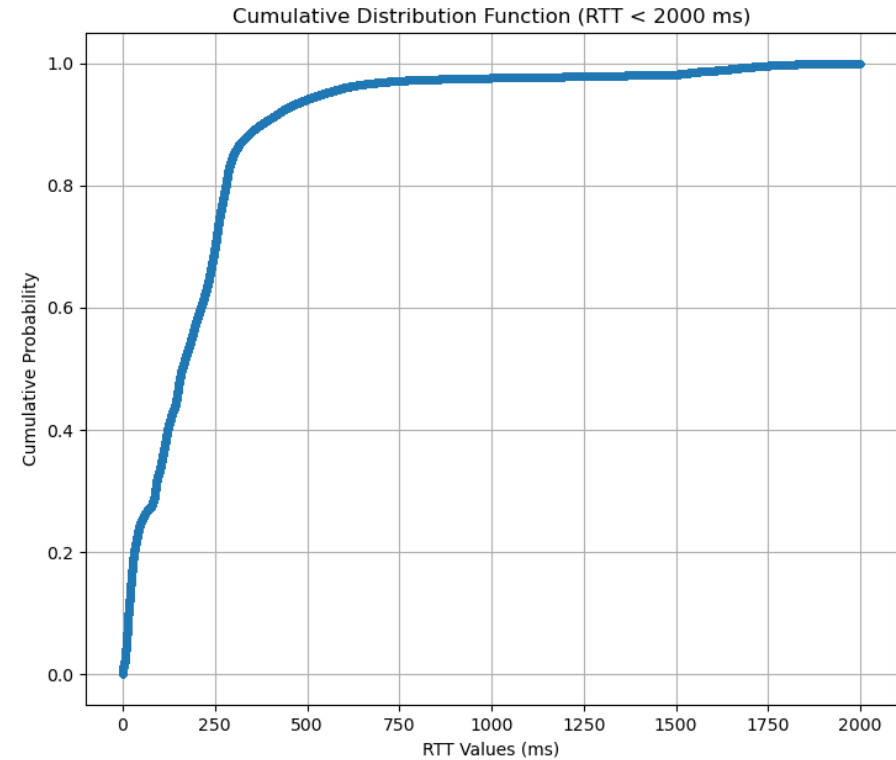
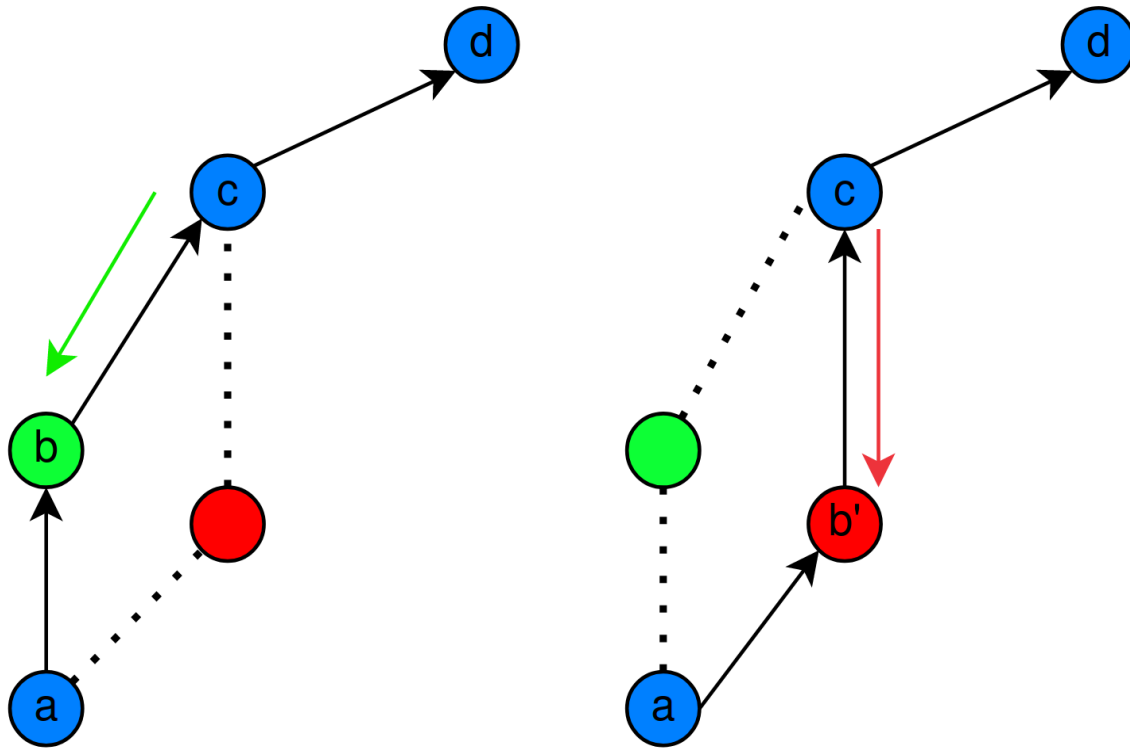
## Problems Clustering:

- Feature vector definition
- Distance definition
- Complexity of the algorithm

## How was it before?

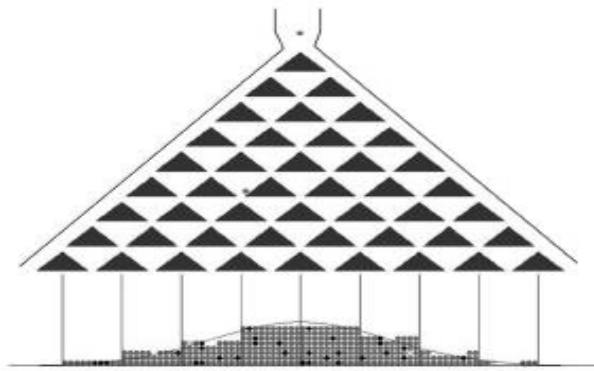
- Inference-based (prefix)
- Measurement based (APAR, MIDAR) (RTT, Loss Rate)
- Traffic-based

# Data Generation – ICMP - RTT curve balls

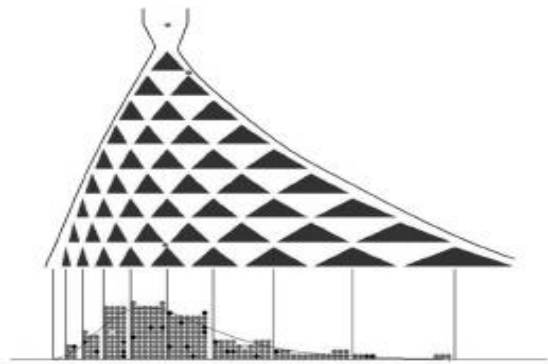


$199,861,638 \text{ m/s} \times 0.250\text{s} \approx 50,000 \text{ km}$

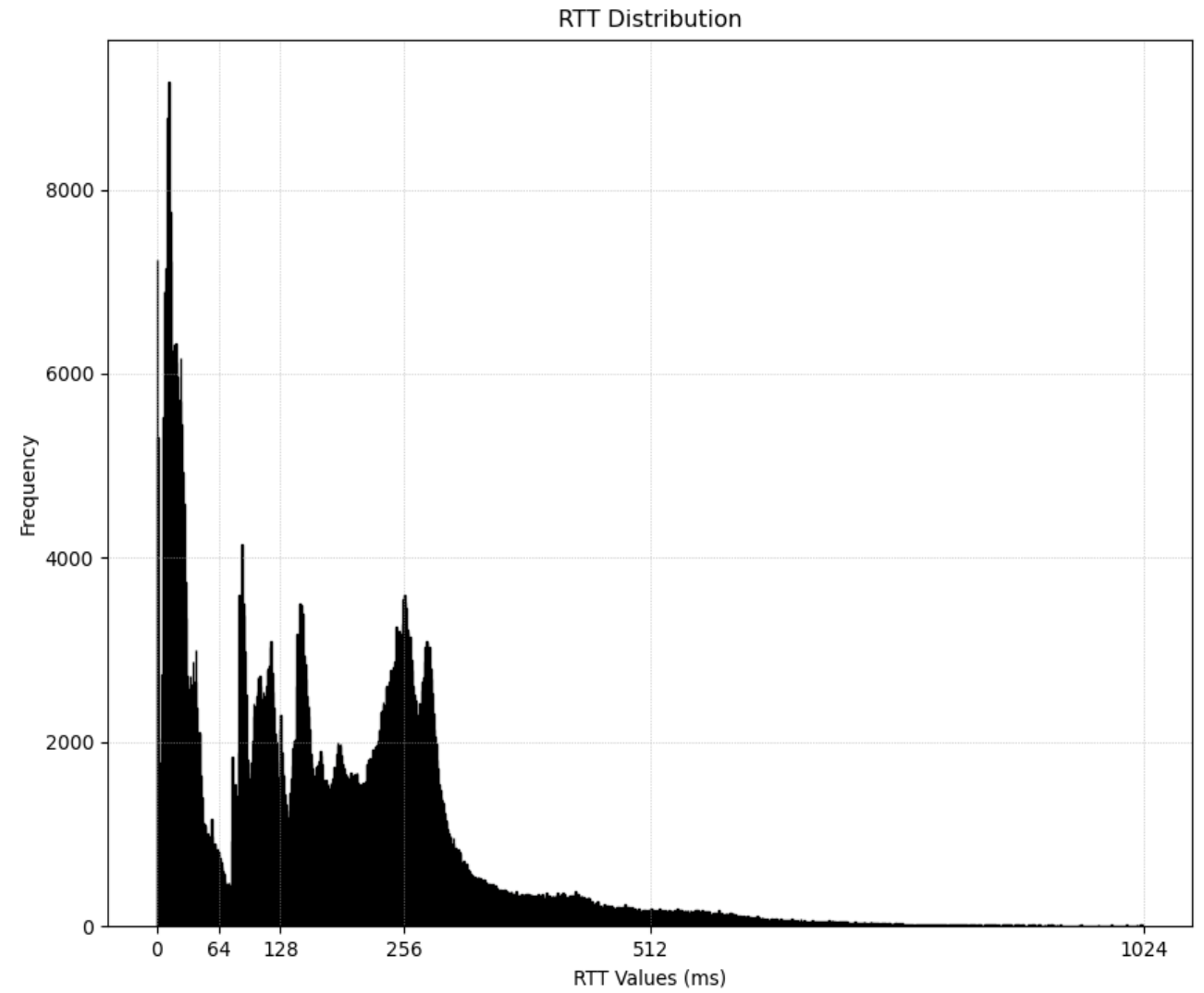
# RTT



(a)



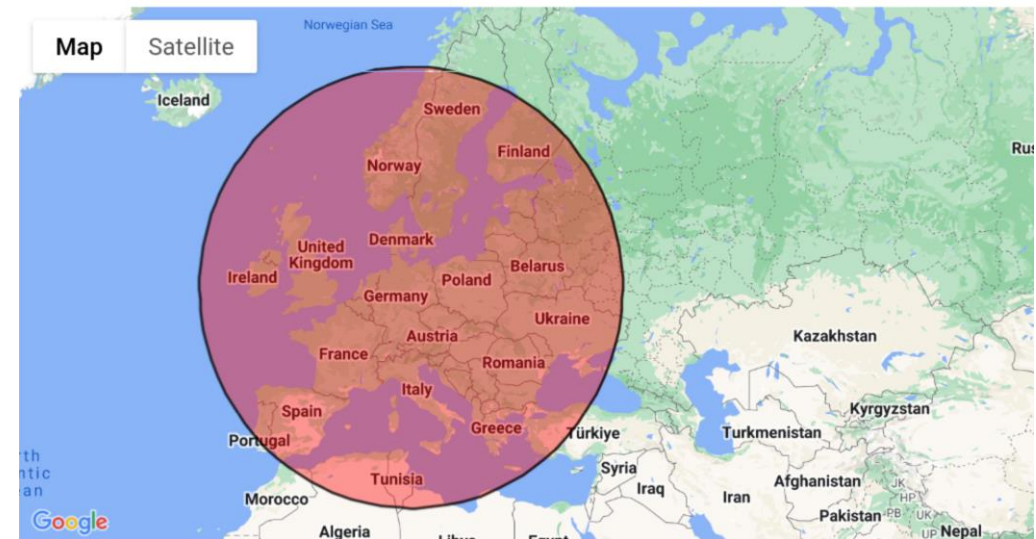
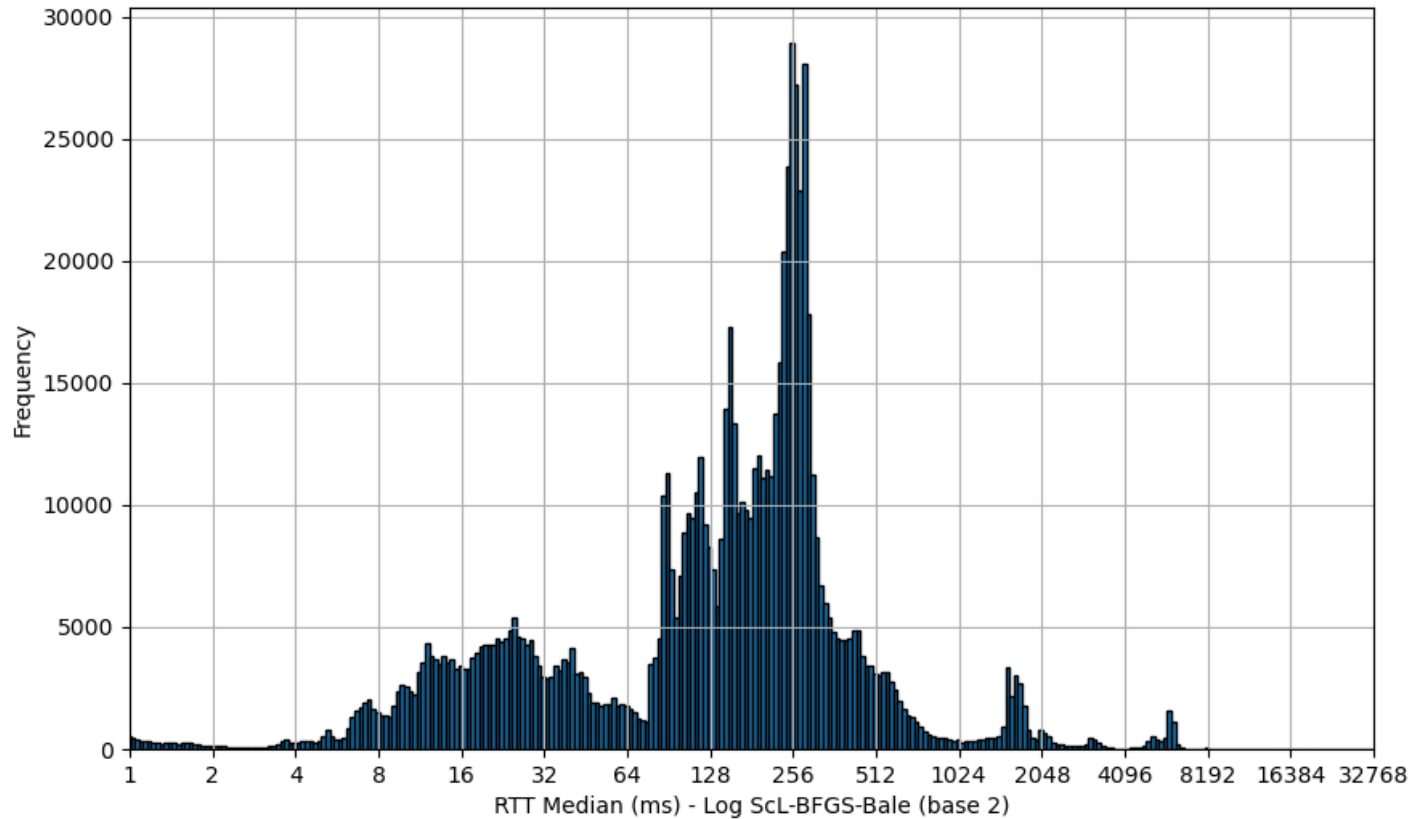
(b)



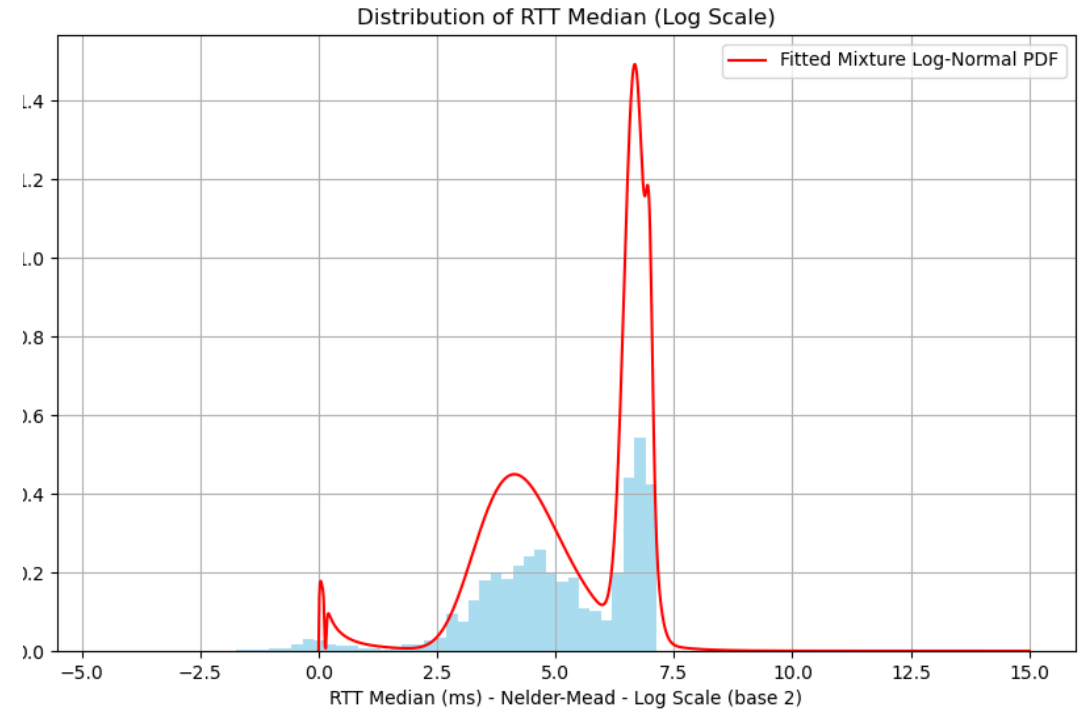
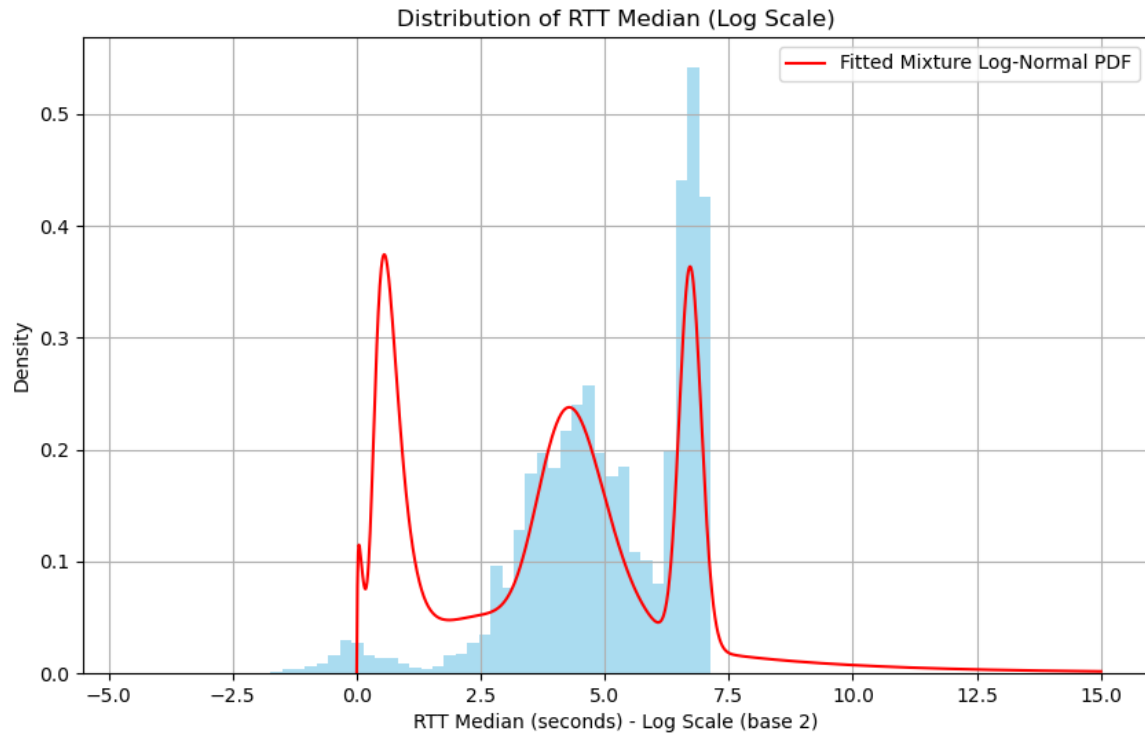
Population Distribution VS. RTT Distribution ?

# RTT - mixture of Log-Normal Distributions.

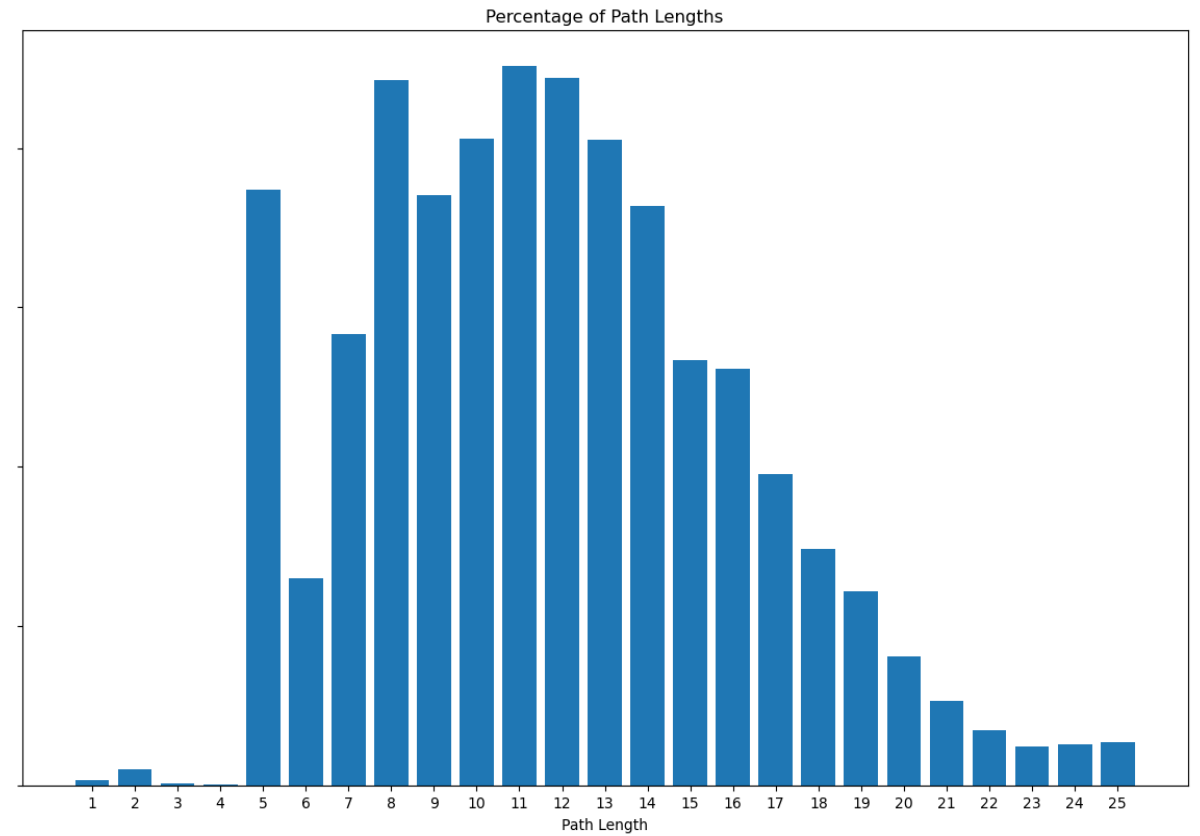
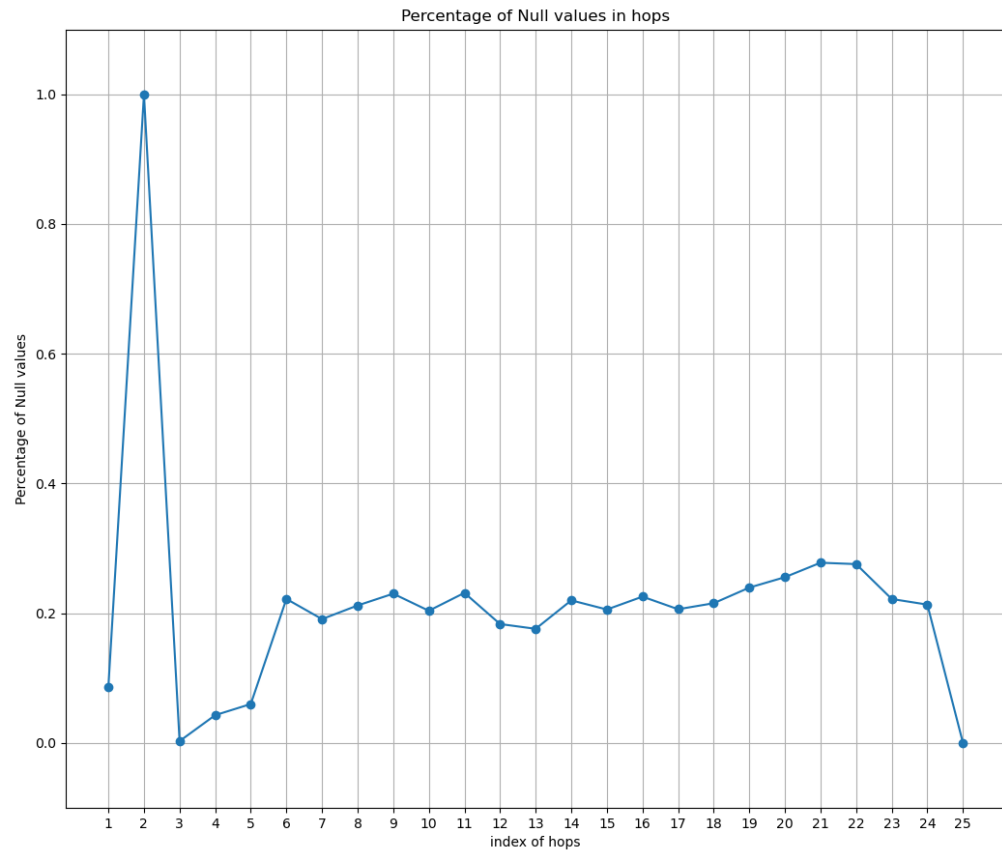
Distribution of RTT Median (Log Scale)



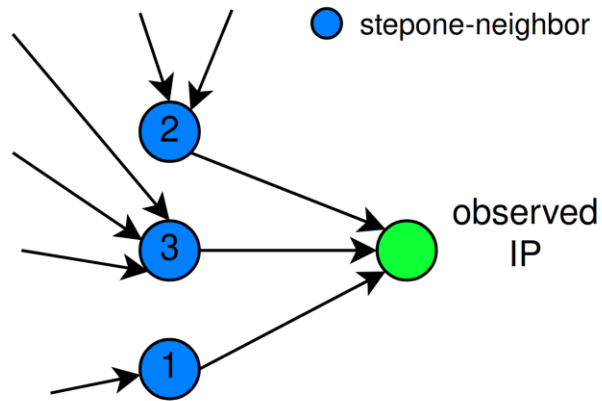
# Fit this distribution



# Null Values

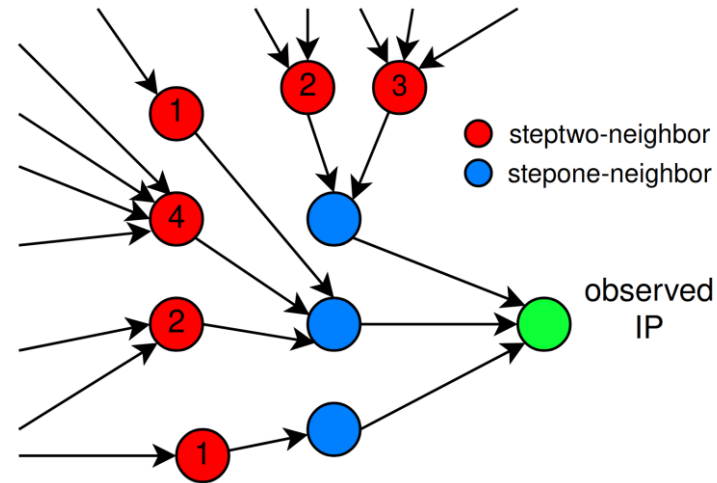


# My Solution - Local Topology



$\{ (2,1), (3,1), (1,1) \}$

Centroid:  $(2,1)$

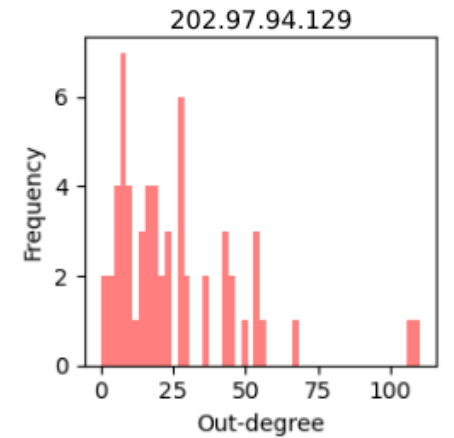
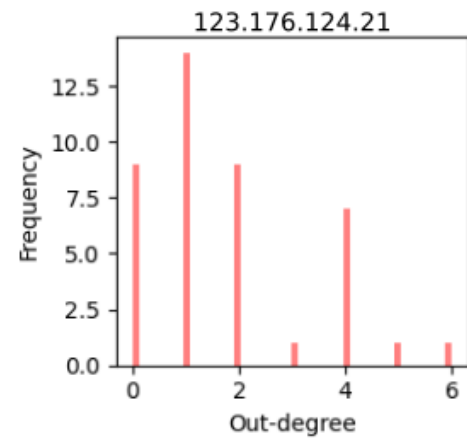
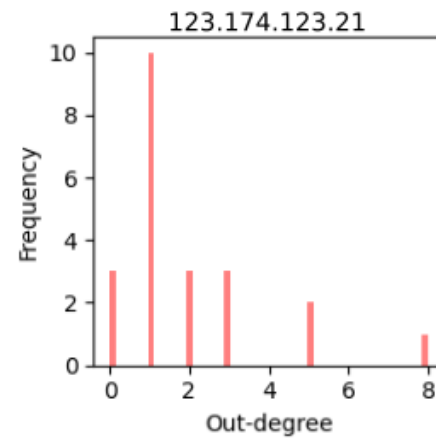
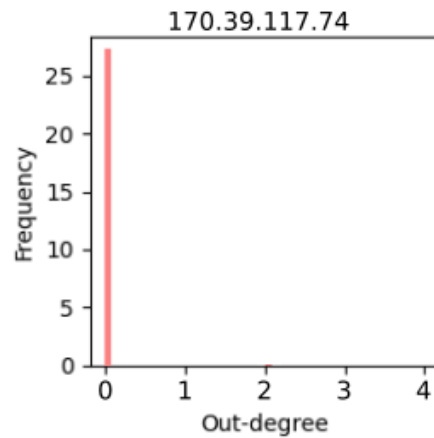
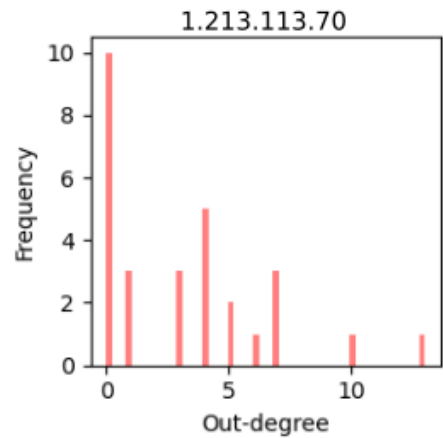
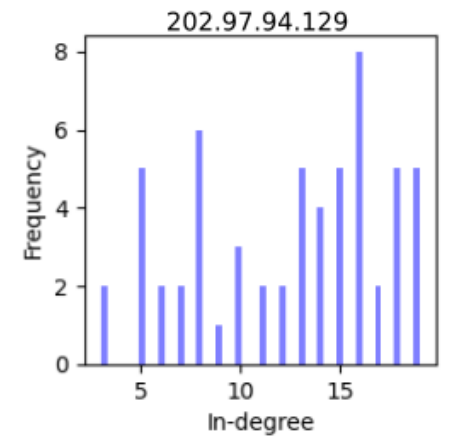
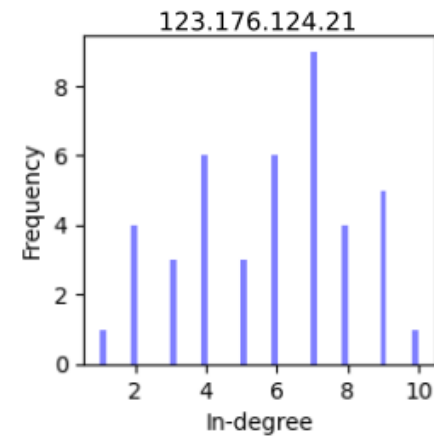
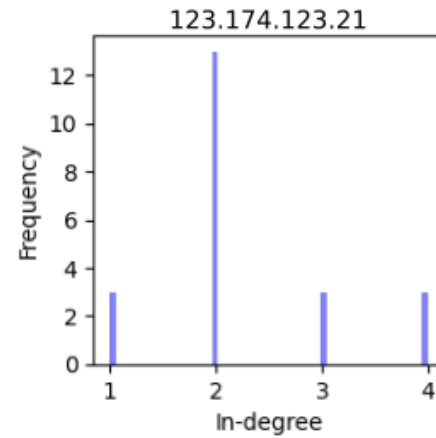
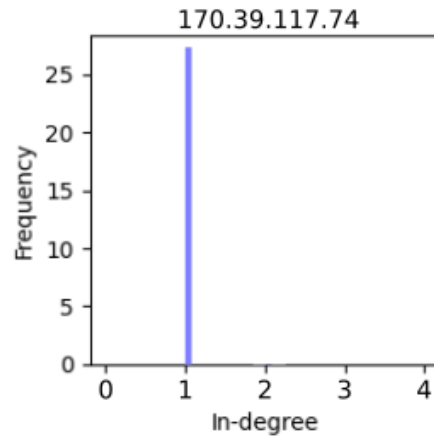
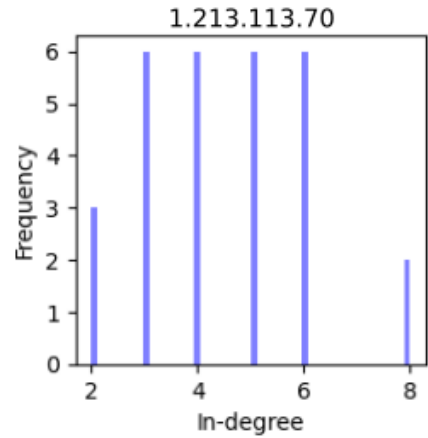


$\{ (3,1), (2,1), (1,1), (4,1), (2,1), (1,1) \}$

$(13/6, 1)$



# In-degree , Out-degree



# Distance

Distance:

- Euclidean Distance of Geometric Centroids
- „Distance“ of Distribution

$$a = [1, 1, 1, 2, 2]$$

$$b = [1, 1, 1, 2, 2]$$

$$c = [1, 1, 1, 2, 2, 99]$$

$$d = [1, 1, 1, 2, 2, 3, 5, 8]$$

$$e = [1, 1, 1, 2, 2, 9]$$

$$f = [8, 9, 9, 100]$$

# Wasserstein Distance

$$W(P, Q) = \sum_{i=1}^n |F_P(x_i) - F_Q(x_i)|$$

F = Cumulative Distribution Functions (CDFs)

Table 1: Wasserstein Distance

	a	b	c	d	e
b	0.0000				
c	16.2667	16.2667			
d	1.4750	1.4750	15.7917		
e	6.4824	6.4824	20.2157	5.0074	
f	30.1000	30.1000	13.8333	28.6250	23.8235

# Jensen-Shannon Divergence

$$JSD(P\|Q) = \frac{1}{2}(D_{KL}(P\|M) + D_{KL}(Q\|M))$$

$$D_{KL}(P\|Q) = \sum_i P(i) \log \left( \frac{P(i)}{Q(i)} \right)$$

Table 2: Jensen-Shannon Divergence

	a	b	c	d	e
b	0.0000				
c	34.3108	34.3108			
d	5.5452	5.5452	33.0887		
e	90.4557	90.4557	109.2773	82.1875	
f	36.3518	36.3518	70.6625	41.8969	126.8075

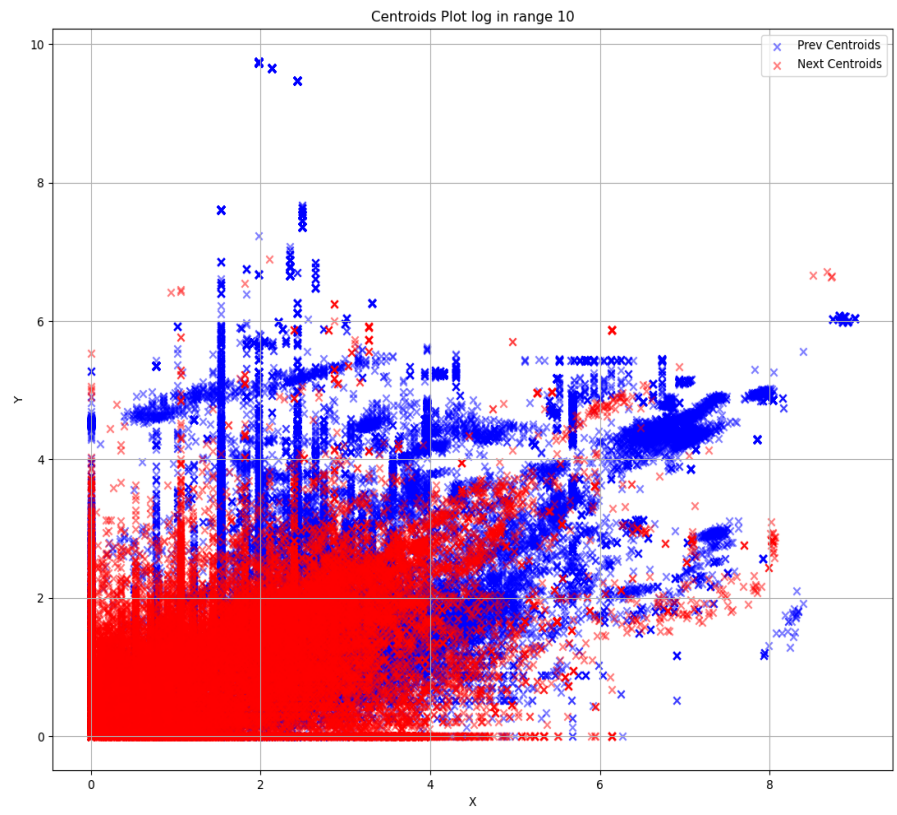
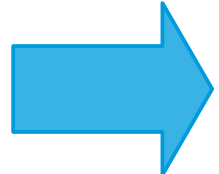
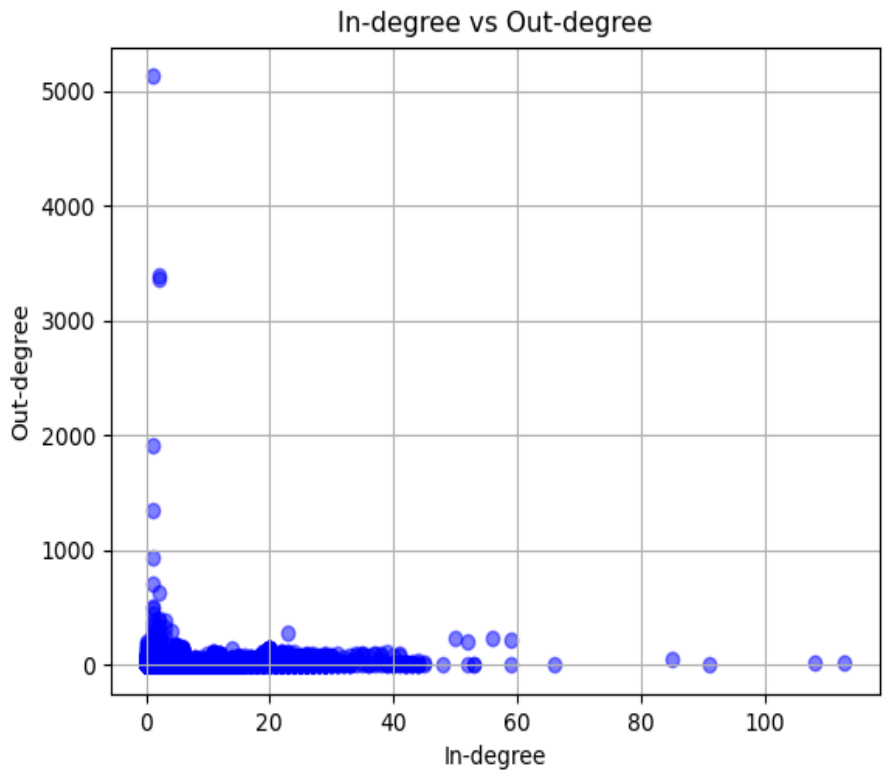
# Jaccard Distance

$$D_J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

Table 5: Customized Jaccard Distance

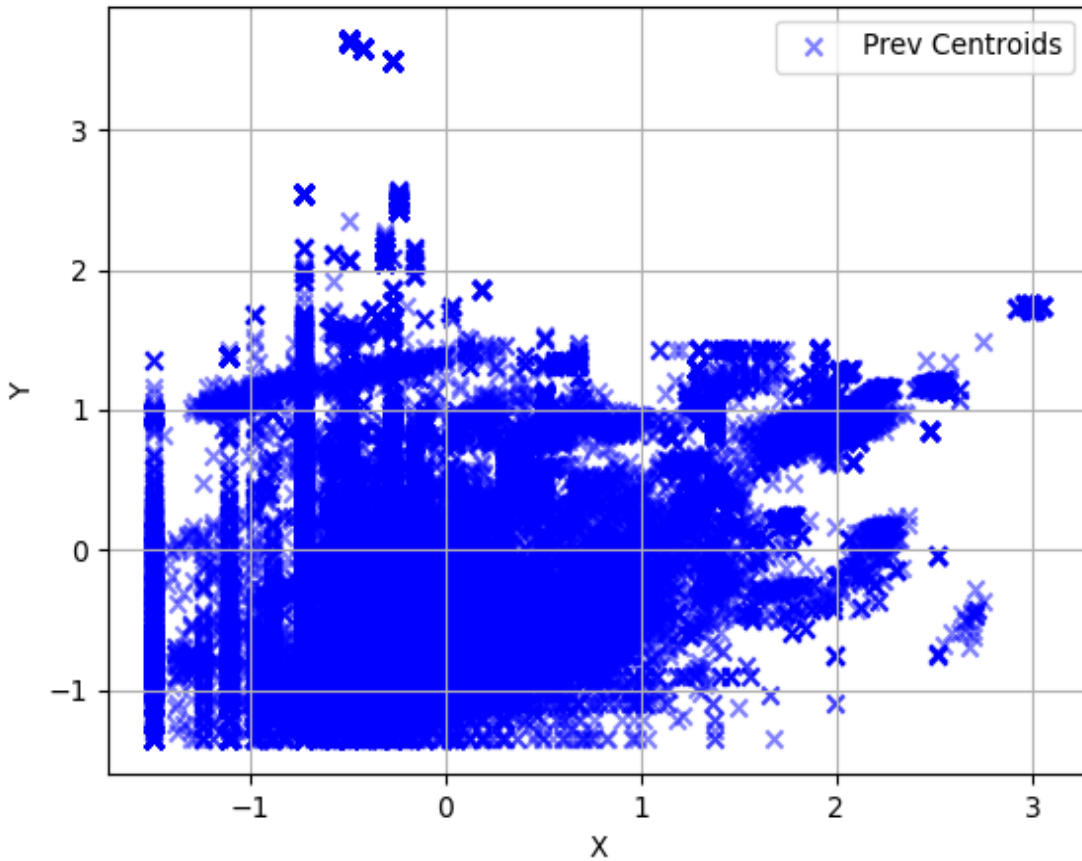
	a	b	c	d	e
b	0				
c	0.0909	0.0909			
d	0.2307	0.2307	0.2857		
e	0.7435	0.7435	0.75	0.7619	
f	1	1	1	0.8333	0.8947

# Z-score log scaled centroid

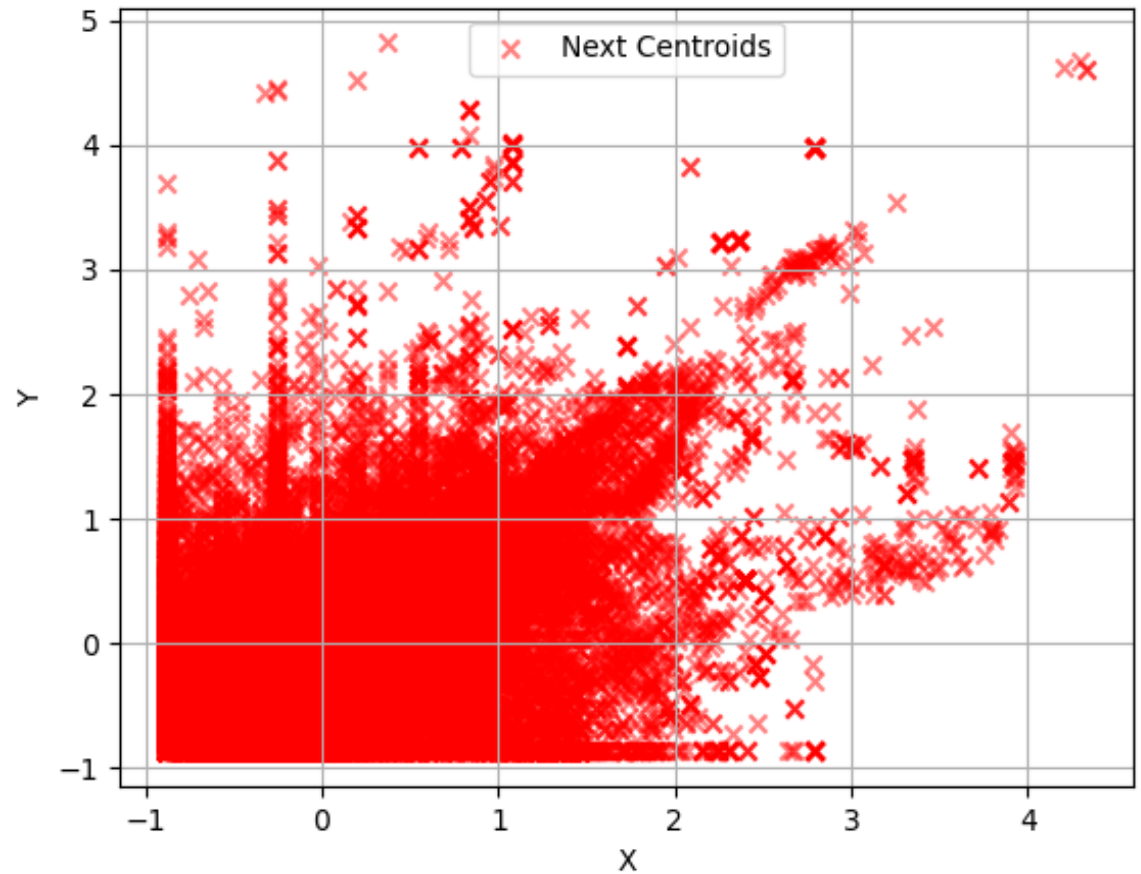


# Z-score log scaled centroid

Z-score log centers



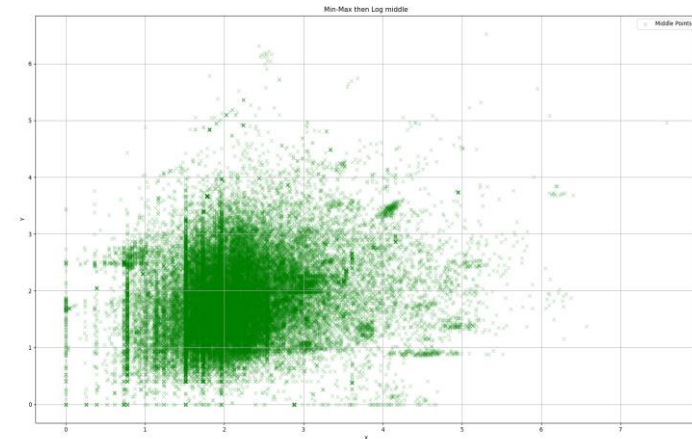
Z-score log centers



# Neighbours, not observed IP



Step one neighbour centroids



Step two

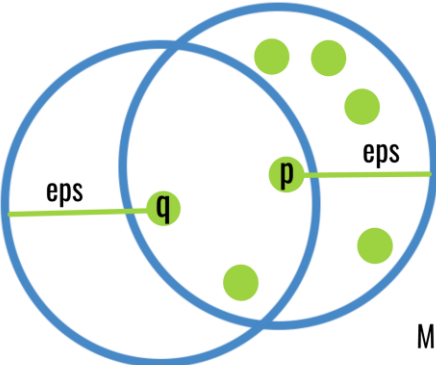


Step three

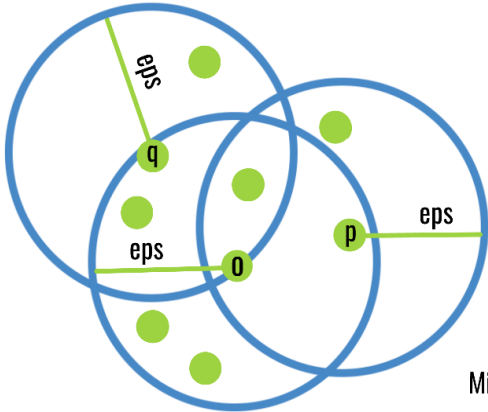


# DBSCAN

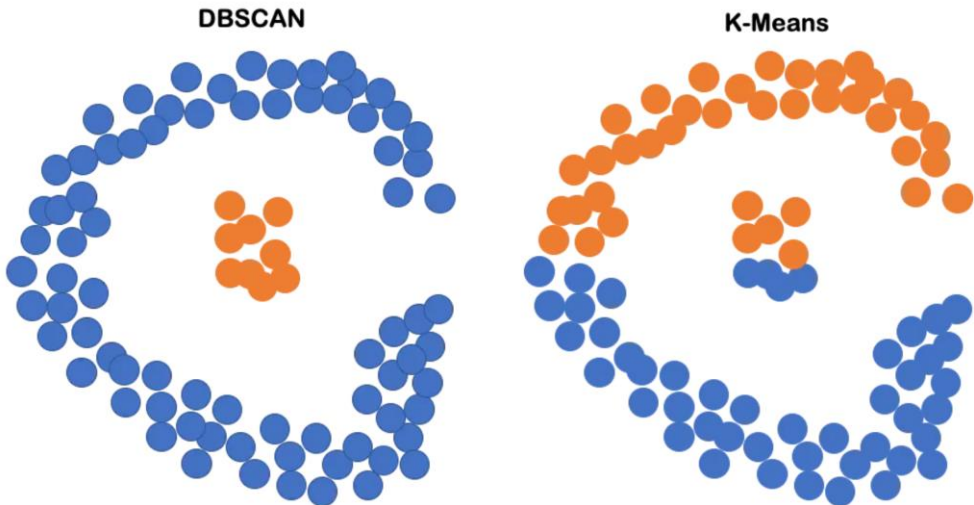
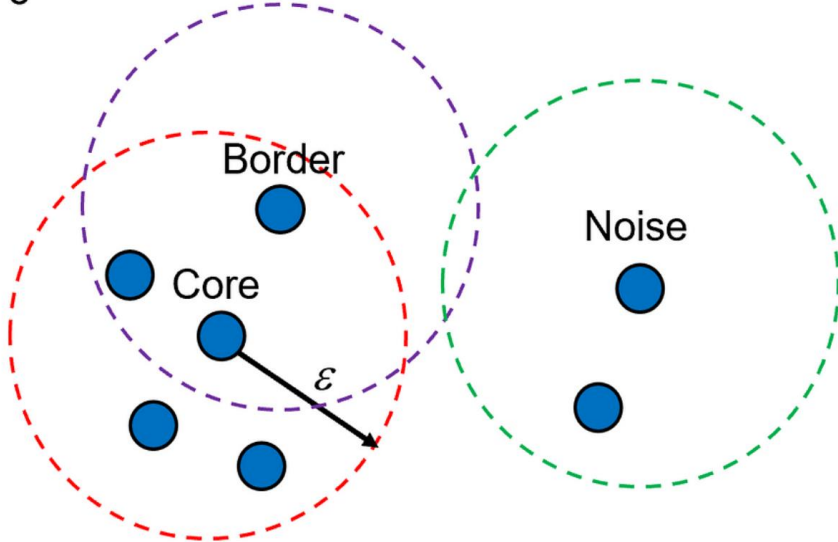
MinPts = 5



MinPts = 6

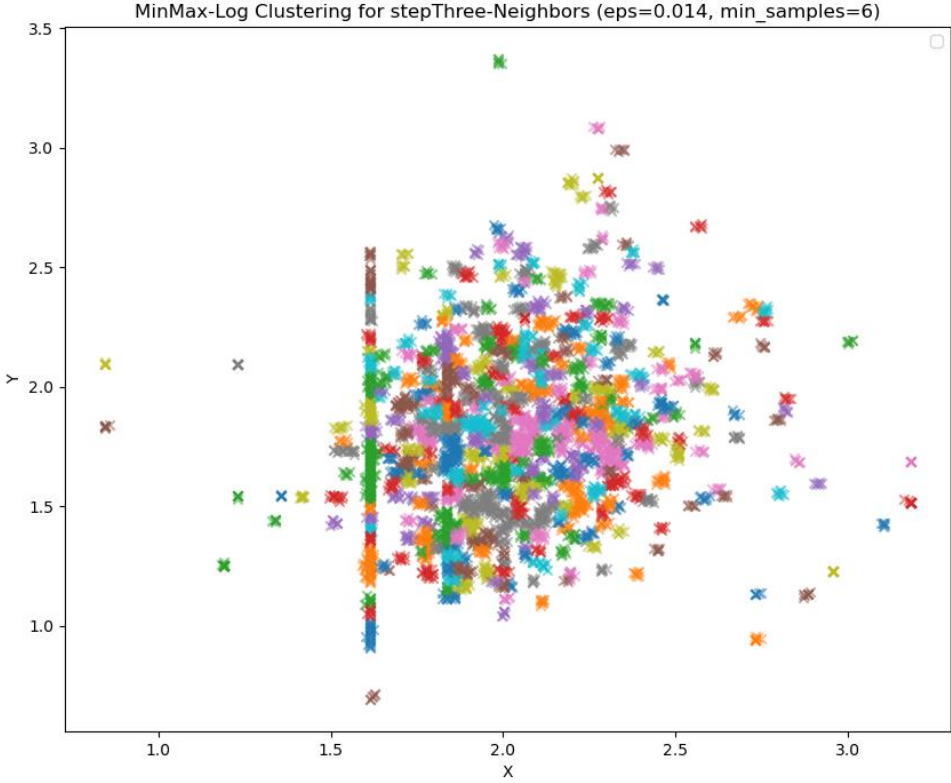
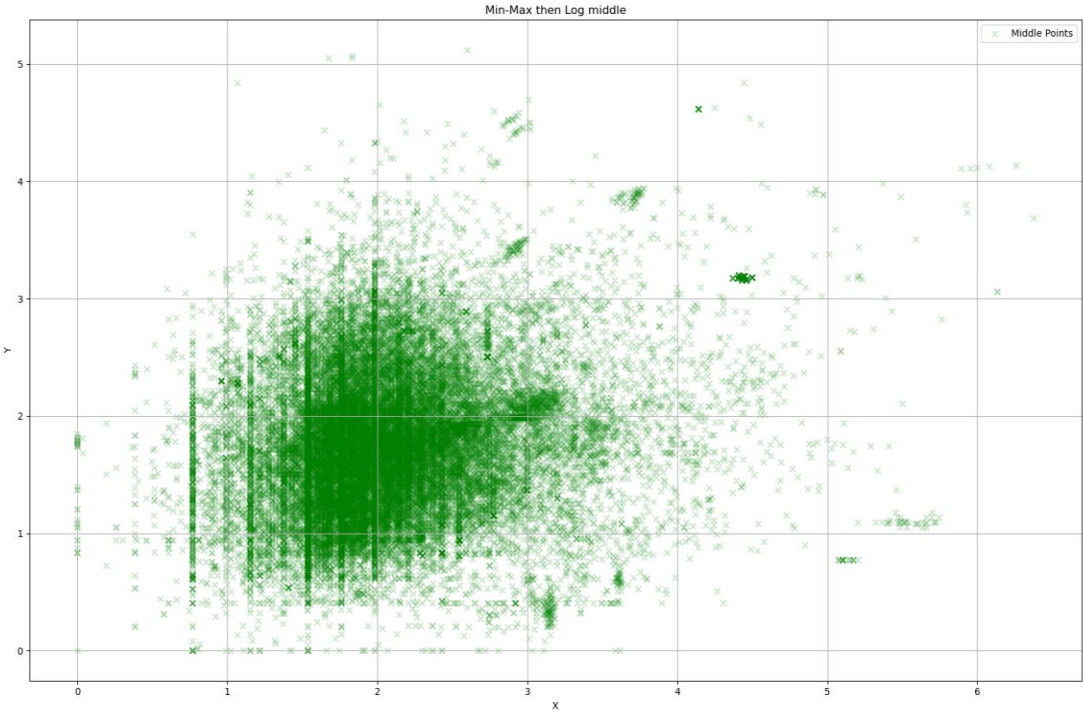


MinPts = 6

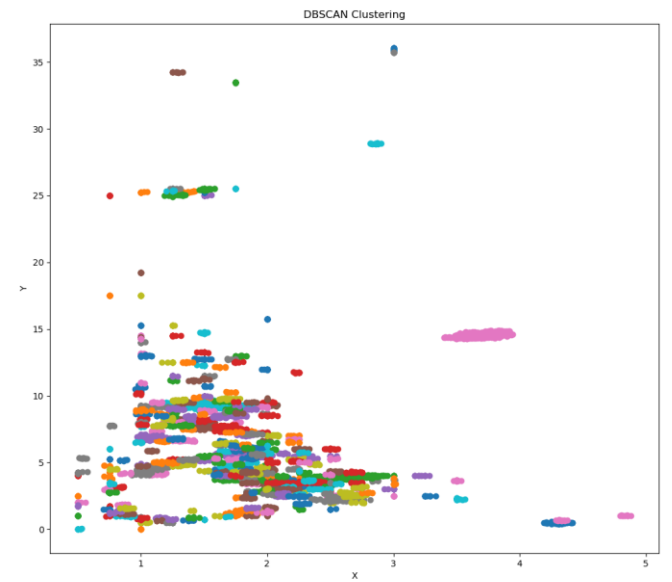
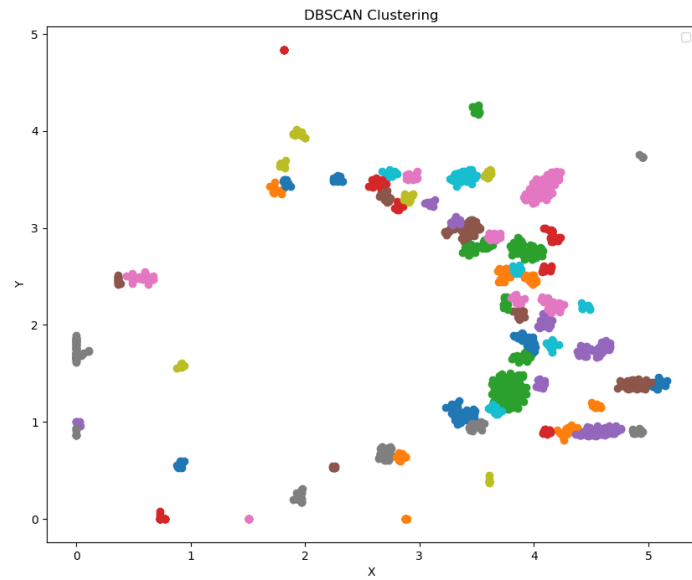
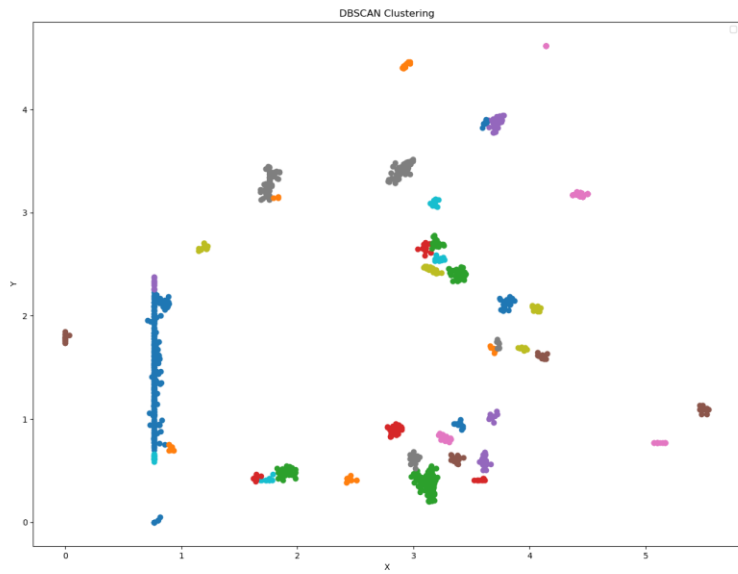


DBSCAN vs K-means. Dots are samples, the X-axis is feature 1 and the Y-axis is feature 2. (Image by author)

# Apply DBSCAN



# Apply DBSCAN with different radii



# Apply DBSCAN iteratively and with RTT

$\epsilon$  from large to small

---

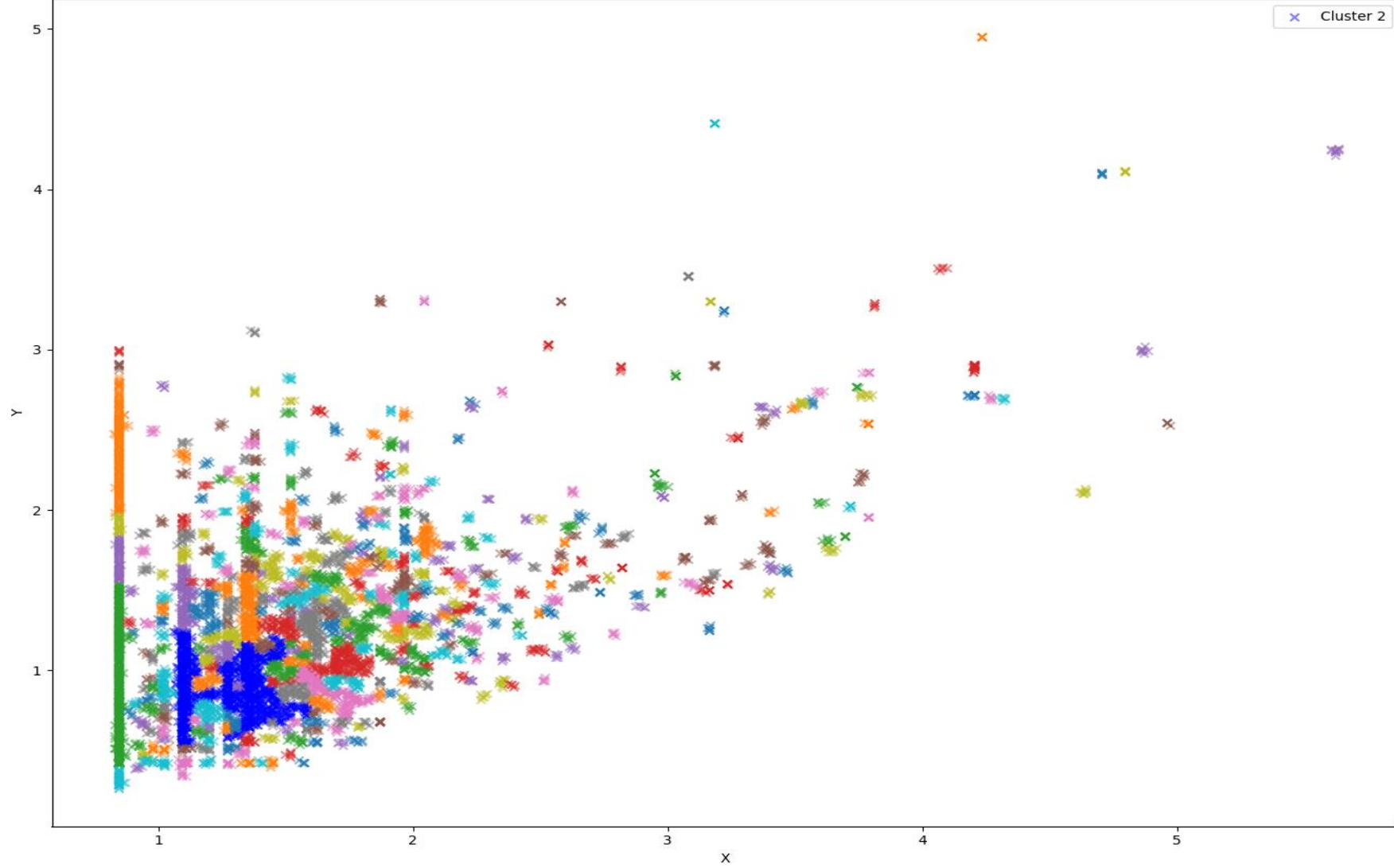
## Algorithm 3 Iterative RTT Enhanced DBSCAN

---

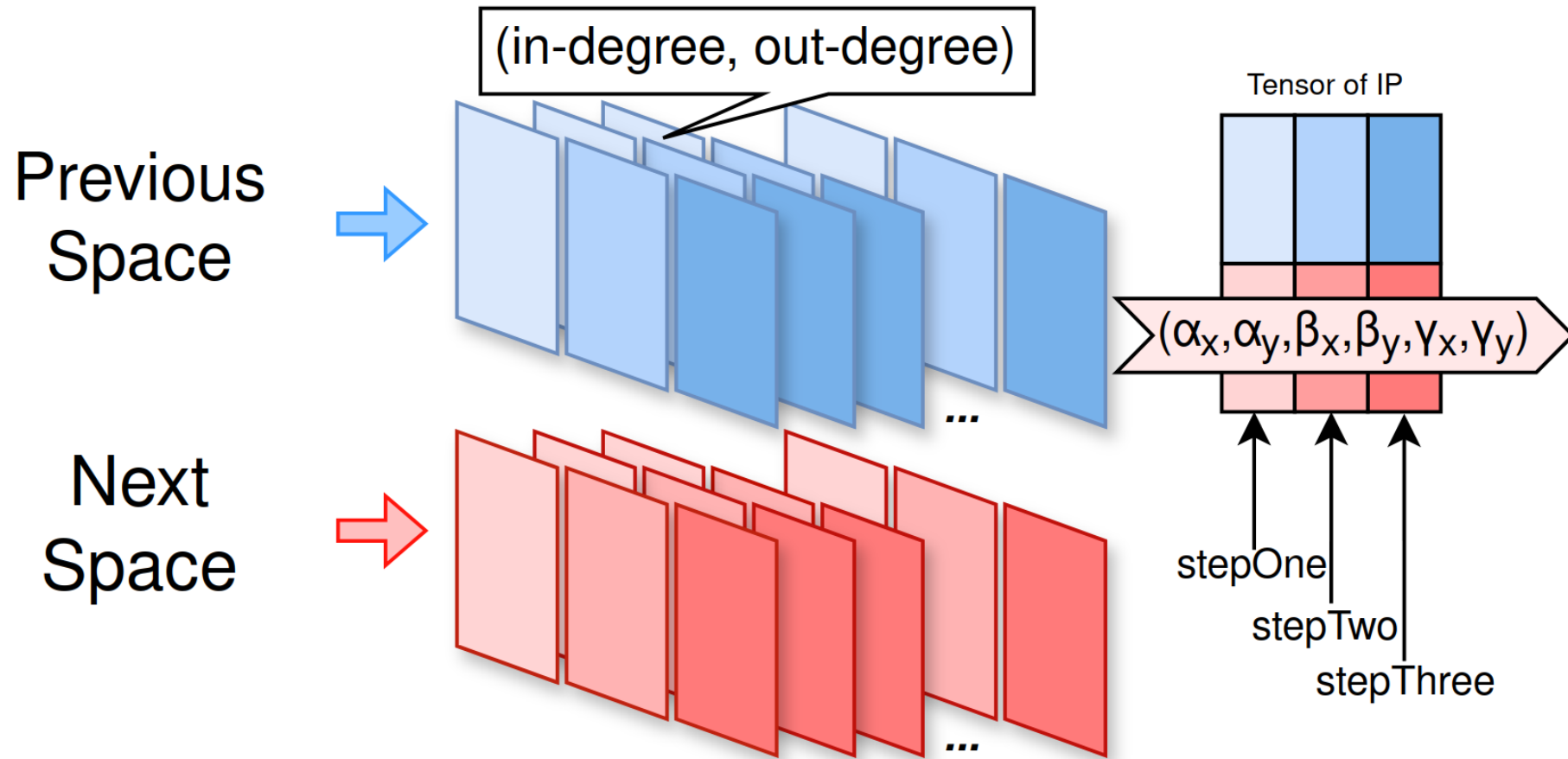
```
1: function ITERATIVERTTENHANCEDDBSCAN( $D, P, \epsilon, \epsilon_{rtt}$ )
2:   Initialize value of threshold  $N$ 
3:   Initialize value of step size  $step$ 
4:    $Pool \leftarrow D$ 
5:    $\Delta \leftarrow 1$ 
6:   while  $\Delta > 0$  do
7:      $PoolSizeBefore \leftarrow$  size of  $Pool$ 
8:      $Clusters, Noise \leftarrow$  RTT-DBSCAN( $Pool, MinPts, \epsilon, \epsilon_{rtt}$ )
9:     for cluster in  $Clusters$  do
10:      if size of  $cluster \leq N$  then
11:        add cluster to  $Result$ 
12:      else
13:        add cluster to  $LargeClusters$ 
14:      end if
15:    end for
16:     $Pool \leftarrow Noise \cup LargeClusters$ 
17:     $\Delta = PoolSizeBefore -$  size of  $Pool$ 
18:     $\epsilon = \epsilon - step$ 
19:  end while
20:  return  $Result$ 
21: end function
```

---

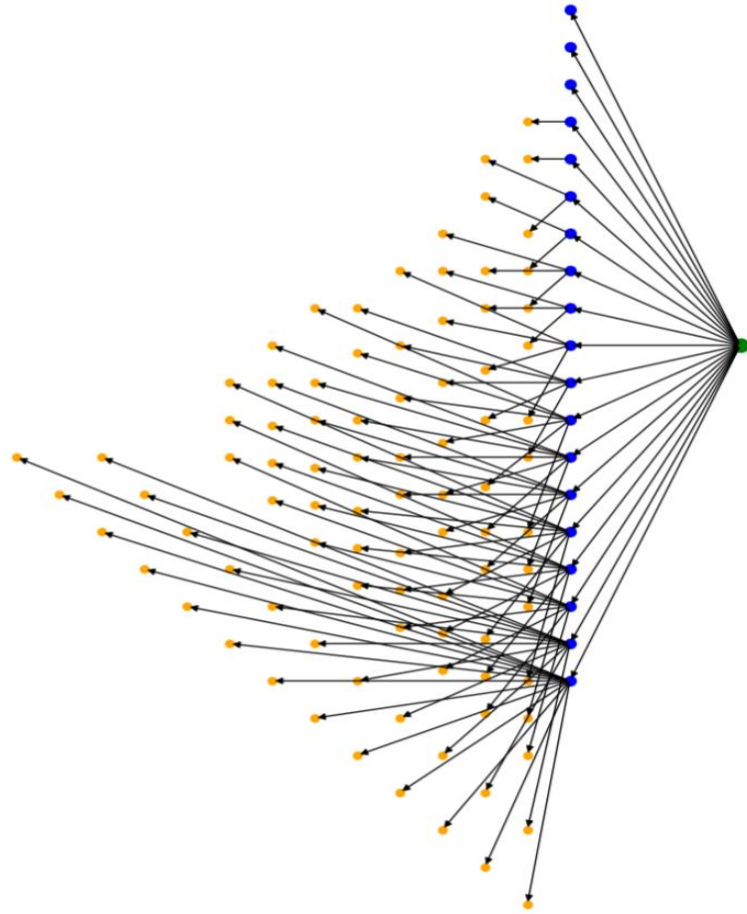
MinMax-Log Clustering for stepOne-Neighbors (eps=0.019, min\_samples=5)



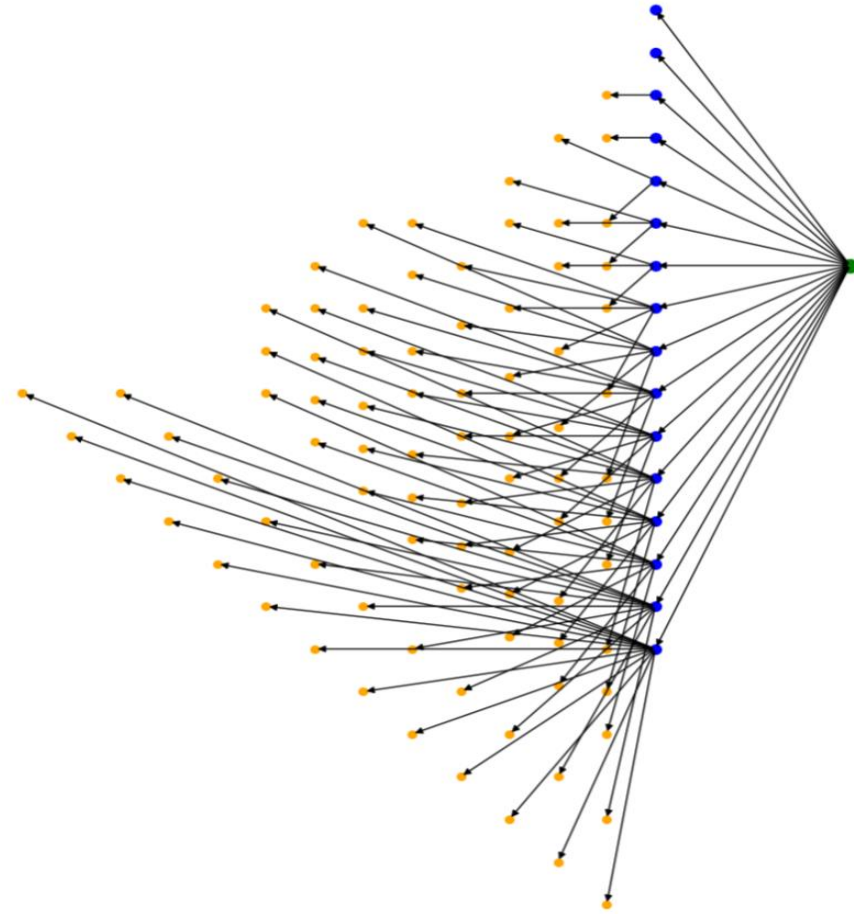
# My Solution - Local Topology



# Results



66.109.3.233



209.18.43.73

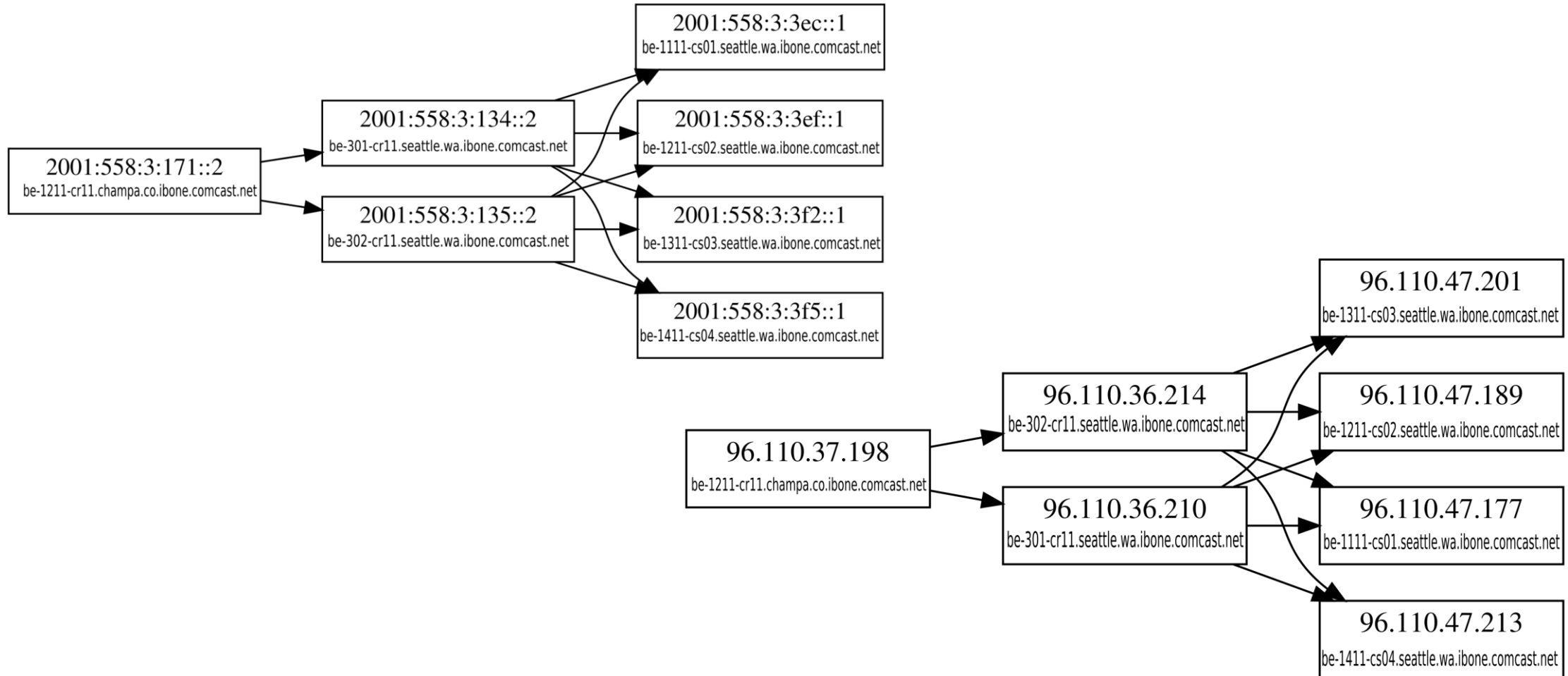
# DNS

Table 8: Cluster "62" IP and DNS

	IP Address	DNS Name
<b>1</b>	2001:558:3:170::2	be-1111-cr11.champa.co.ibone.comcast.net
<b>2</b>	2001:558:3:172::2	be-1311-cr11.champa.co.ibone.comcast.net
<b>3</b>	96.110.37.202	be-1311-cr11.champa.co.ibone.comcast.net
<b>4</b>	96.110.37.206	be-1411-cr11.champa.co.ibone.comcast.net
<b>5</b>	58.138.106.18	osk004ipgw01.IIJ.Net
<b>6</b>	96.110.37.198	be-1211-cr11.champa.co.ibone.comcast.net
<b>7</b>	2001:558:3:173::2	be-1411-cr11.champa.co.ibone.comcast.net
<b>8</b>	96.110.37.194	be-1111-cr11.champa.co.ibone.comcast.net
<b>9</b>	2001:558:3:171::2	be-1211-cr11.champa.co.ibone.comcast.net



# IPv4 , IPv6



# Future Works

## 1. Pattern explicit definition

- „Bridge, dumbbell, star“ .
- Router typ definition on Patterns.
- Routing Rules Detection/ Prediction
- Relationships between ASes.

## 2. Supervised Learning

## 3. Dynamic Sampling

- Update, Delete...

## 4. Fault Detection

- Pattern - „Should be a EBGP, act like a IGP, why?“
- Dynamy – „It was a EGP yesterday, looks like a IGP today, why?“

## 5. #ALLHands