

UNIVERSITAT POLITÈCNICA DE
CATALUNYA (UPC) – BARCELONATECH

FACULTAT D'INFORMÀTICA DE BARCELONA (FIB)

MASTER IN INNOVATION AND RESEARCH IN
INFORMATICS

DATA SCIENCE

Title

Author:
Sergi CAROL BOSCH

Supervisor:
Dr. Albert CABELLOS -
Computer Architecture
Co Supervisor:
Dra. Marta ARIAS -
Computer Science

June 12, 2019



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH

Abstract

Today, network operators still lack functional network models able to make accurate predictions of end-to-end Key Performance Indicators (e.g., delay or jitter) at limited cost. RouteNet is a Graph Neural Network that aims to provide this lacking functionality which models network schemes and is able to predict the key performance indicators. Thanks to its GNN architecture that operates over graph-structured data, RouteNet reveals an unprecedented ability to learn and model the complex relationships among topology, routing and input traffic in networks. In previous related works it was proved that RouteNet was able to generalize multiple computer networks by training the model with various topologies of different sizes. This thesis aims to prove how RouteNet is able to model a computer network, but also, and more importantly, wants to provide the ground work for routing sampling creation, which allows to create enough routing combinations to provide a representative sample from which RouteNet would be able to be trained by only simulating a single topology.

Proving this last objective would permit future research to be done much faster, since simulation and training of multiple topologies is a time costly procedure. It would also allow us to model any other network by only using the routing samples from one simulated network.

Contents

1	Introduction	7
1.1	Related Work	8
1.2	Objectives	8
1.3	Document Structure	9
2	State of the Art	10
2.1	Notation	10
2.2	Graph Neural Networks	10
2.3	RouteNet characteristics	11
2.4	RouteNet	12
3	Methodology	15
3.1	Data Gather	15
3.2	Model Evaluation	15
4	Dataset and Technical Details	16
4.1	Introduction	16
4.2	Dataset	16
4.2.1	Data Gather - Simulations	16
4.2.2	Dataset	17
4.2.3	Routing Creation	18
4.2.4	Traning and Evaluation	21
4.2.5	Data Features	22
4.2.6	50 Nodes Alternative	22
4.2.7	50 Nodes	26
4.2.8	24 Nodes	27
4.2.9	14 Nodes	29
4.3	Architecture	30
4.3.1	Recurrent Neural Networks	31
4.3.2	Gated Recurrent Unit	31
4.3.3	Readout function	32
4.3.4	Loss minimization function	34
4.3.5	Performance Metrics	35
5	Evaluation	36
5.1	Introduction	36
5.2	Model Performance	36
5.3	Metrics	37
5.3.1	50 Nodes Alternative	37
5.3.2	50 Nodes	40
5.3.3	24 Nodes	41
5.3.4	14 Nodes	42

6	Conclusions	45
6.1	Future Work	48
7	References	50
8	Appendix	52
8.1	True Delay vs Predicted 50 Nodes alternative	52
8.2	True Delay vs Predicted 50 Nodes	53
8.3	True Delay vs Predicted 24 Nodes	53
8.4	True Delay vs Predicted 14 Nodes	54
8.5	Table of results for 14 Nodes with λ 15	55

List of Figures

1	Molecules have a deep relationship between the atoms that compose them, and are a clear example of Graph with explicit relationship between nodes. [1]	10
2	A message m is from v_i to v'_i .	11
3	v'_i updates its state using the values received from adjacent nodes.	11
4	The Graph updates its state.	11
5	RouteNet Architecture.	14
6	Evolution of delay over time for 50 nodes for each traffic intensity λ .	16
7	Evolution of delay over time for 100 nodes for each traffic intensity λ .	16
8	Network graph with 14 nodes.	17
9	Network graph with 24 nodes.	17
10	Network graph with 50 nodes.	17
11	Example of routing with loops.	20
12	Histogram of bandwidth distribution	23
13	Histogram of delay distribution	24
14	Histogram of packets distribution	25
15	Histogram of drops distribution	25
16	Delay vs Bandwidth	25
17	Delay vs Drops	25
18	Dataset Summary for 50 Nodes Alternative	26
19	Distribution of the bandwidth for dataset with 50 nodes.	26
20	Distribution of the delay for dataset with 50 nodes	26
21	Distribution of the packets for dataset with 50 nodes.	27
22	Delay vs Drop for 50 Nodes	27
23	Distribution of the delay for dataset with 24 nodes.	28
24	Distribution of the bandwidth for dataset with 24 nodes	28
25	Delay vs Bandwidth for dataset with 24 nodes.	29
26	Distribution of the bandwidth for dataset with 14 nodes.	30
27	Distribution of the delay for dataset with 24 nodes	30
28	Delay vs Bandwidth for dataset with 14 nodes.	30
29	Delay vs Drop for dataset with 14 nodes.	30
30	Gradient vanishing problem	31
31	Gated Recurrent Unit [17]	32
32	Model Structure	33
33	Adam performance compared to other models as described in [21]	34
34	GRU Kernel candidate distribution.	36
35	GRU gate value distribution.	36
36	GRU bias for the output.	37
37	Loss (MSE) for training and evaluation data. The blue line indicates the evolution of the training data, the grey line indicates the evolution of the evaluation data.	37
38	R_2 for training and evaluation data.	38

39	MRE for training and evaluation data.	38
40	MAE for training and evaluation data.	38
41	ρ for training and evaluation data.	39
42	True delay vs Predicted for 50 nodes Alternative.	40
43	True delay vs Predicted for 50 nodes.	41
44	True delay vs Predicted for 24 nodes.	42
45	True delay vs Predicted for 14 nodes.	43
46	Cumulative density function for mean relative error.	44
47	Results from the old model trained with two different topologies and evaluated with the 24 Nodes topology.	46
48	Results for the 24 node topology with the model trained with the alternative routing dataset.	47
49	λ 9 routing 8	52
50	λ 15 routing 8	52
51	λ 9 routing 9	52
52	λ 9 routing 58	52
53	λ 9 routing 9	53
54	λ 12 routing 46	53
55	λ 15 routing 51	53
56	λ 12 routing 20	53
57	λ 8 routing 17	53
58	λ 10 routing 38	53
59	λ 10 routing 1 k 8	54
60	λ 8 routing AL 1 k 5	54
61	λ 8 routing 2 k 5	54
62	λ 8 routing AL 2 k 3	55
63	λ 8 routing 2 k 5	55

List of Tables

1	Number of routings per topology	18
2	Simulations per lambda and nodes	21
3	Total amount of samples per each routing.	21
4	Training and evaluation split.	22
5	Correlation between features for the 50 Nodes Alternative dataset	26
6	Table of results for evaluation dataset	39
7	Metric results for 50 Nodes.	41
8	Table of results for 24 nodes.	41
9	Table of results for 14 nodes without the data with λ 15.	43
10	Results for 14 nodes with λ 15 included in the dataset	55

1 Introduction

Graph Neural Networks have been a hot topic in the recent months, with the first publication of the Graph Neural Network topic released on 4 Jun 2018 by the Deep Mind team at Google [1]. Since this publication the use of graph neural networks to solve problems that had a deep relation with structured data, mainly graphs, has escalated greatly. Before the publication of the mentioned before paper, the modeling of systems which rallied in structured forms of data was problematic, eg: routing network systems, optimal traffic routing for cars, or chemical reactions prediction.

With the introduction of systems such as SDN [2], and more importantly KDN [3], the paradigm of routing optimization has been made easier thanks to the simplification provided by these systems to observe and control the network, which allows traffic engineers to make better and faster decisions regarding the state of the network.

Yet, it is believed with the introduction of the concept of *Knowledge Define Networking*, amortization, or at least some degree of it, can be achieved. Therefore research on using machine learning concepts to networking routing has been an interesting topic for many researches and companies. With this interest in mind, development towards the creation of a model that can be used to optimize or automatize traffic engineering has made significant progress. These kinds of optimization's are really interesting for big companies in the market since it would mean that the burden of network and traffic engineering is no longer relied to solely on persons but the algorithms can be the ones to decide the different combinations of options or at least can help the engineers to achieve better solutions.

The aim of this thesis will be outlined in the objectives section, but as a quick explanation the thesis objectives are to introduce a Graph Neural Network model benchmark to the problem of network routing optimization. Routing optimization has been one of the important topics for computer network engineers. A good routing policy can increase significantly the performance of any network system, and has been deeply researched by different institutions [4] [5]. Since it is not an easy topic and there is no perfect solution for this kinds of problems due to the changing nature of internet traffic networks, specialized personnel is needed in order to create optimized routing in the networks and configure the different devices that live in the internet network appropriately.

The objective is also to prove how the creation of alternative routing combinations helps generalize any kind of topology. meaning that we do not need to train a model with multiple topologies in order to achieve a good model performance.

1.1 Related Work

In a previous paper presented in the SIGCOMM conference myself and the research team I work with introduced a paper with the aim to challenge the generalization capabilities of Graph Neural Networks. In the mentioned paper our research team demonstrated that a Graph Neural Network is able to predict the delay of any kind of network provided that the training data comes from big enough representative sample with the network size. Yet this research was done on the basis that in order to be able to generalize a computer network, multiple topologies from different sizes were needed.

1.2 Objectives

The aim of this thesis is to prove the usefulness of Graph Neural Networks to accurately predict the delay of a computer network, more over instead of using a technique similar to the one done in the previous work, the paper presented in the *SIGCOMM* conference, which involved training the model with multiple topologies with different sizes, this thesis aims to prove that this is not necessary, but instead only a representative sample of routing is needed to provide training data for a model to be able to train a Graph Neural Network that is capable to correctly generalize any kind of computer network regardless of size. To sum up, the objectives are the following:

- Prove that a Graph Neural Network is able to correctly model a computer network.
- Prove that in order to obtain a good model it is not needed to use different topologies sizes as train data.
- Prove that in order to obtain an accurate Graph Neural Network, only a representative sampling of different routing from one topology is able to create a model which is capable of correctly represent any computer network.

If the thesis objectives proves satisfactory, it would mean that we are able to produce a model which is able to generalize any kind of network without the need to simulate a multitude of different sized networks, which is time consuming just to simulate, and even more time consuming in order to train a new model. This would mean that with one topology thoroughly explored we could be able to model any network no matter the size.

1.3 Document Structure

This document will be structured in the following way: it will begin by explaining the state of the art model, which will show the theoretical point of view for RouteNet, then it will follow the methodology used during the thesis, after the technical explanation/ implementation of the routing algorithms, dataset introduction, and RouteNet architecture will be described. With the technical and theoretical model explained, the fifth section will look at the results of the model evaluation and will explain them in detail.

Finally I will end the document with the conclusions of the thesis as well as some future work that this thesis opens up.

2 State of the Art

The proposed model is an *state of the art* Graph Neural Network, called **RouteNet**, which is able to understand the complex characteristics of a network model by interpreting it as a graph. The model is able to understand the inner structure of a network, in which the total delay of an IP packet depends on the state of the whole network. As such, the computation of the delay has to take into account all the other packets in the network by calculating the state of the links and the paths of the graph.

2.1 Notation

A computer network is represented by a set of links $N = l_i$, and a routing scheme is composed by a number of different paths $R = p_k$. Each path p_k is formed by a set of links $p_k = (l_{k_1}, l_{k_2}, \dots, l_{k_{(|p_k|)}})$. The properties of the paths and links are denoted \mathbf{x}_{p_i} and \mathbf{x}_{l_i} . The state of a path is called h_{p_i} , and the state of a path is stated as h_{l_i} . d is understood as the delay of a link or a path. During the course of the explanations I will use both the terms graph and computer network indistinctly, since they are representatives of the same thing in the terms of the architecture of the Graph Neural Network.

2.2 Graph Neural Networks

Graph Neural Networks (**GNN**) are a recently introduced concept in the Neural Network paradigm. GNN work by comprehending the architecture of a graph, and the deep relationship that exists between nodes, as such GNN are able to understand the explicit representations of nodes and edges and also find rules for computing such relations. Graph Neural Networks have been previously successfully used to perform computations for Quantum Chemistry in order to predict quantum properties for organic molecules, using a GNN is able to reduce dramatically the calculation time of this kind of really computer intensive calculations [6].

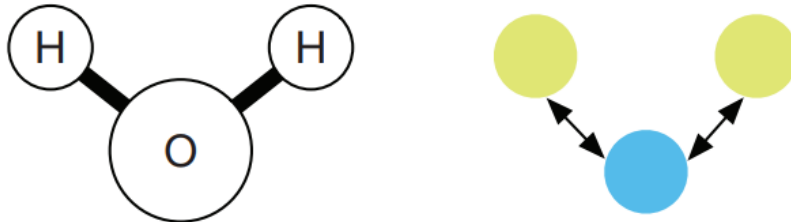


Figure 1: Molecules have a deep relationship between the atoms that compose them, and are a clear example of Graph with explicit relationship between nodes. [1]

Due to the explained deep relation between the different nodes that conform a graph, GNN take into consideration these relationships in the form of message passing m between the different connected nodes in a graph. The sequence of events that a GNN uses to work with the mentioned relationship between nodes is the following:

1. A node v_i in a graph will send a message m to another adjacent node v'_i
2. The node v'_i will then update its value using its current value, and the value from $m_i, i \in (1, 2, \dots, \text{adj}(v'_i))$.
3. Once the update of the nodes in the graph has been performed, update the update for the global attribute of the graph.

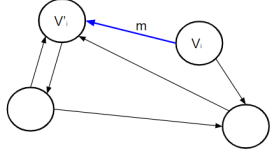


Figure 2: A message m is from v_i to v'_i .

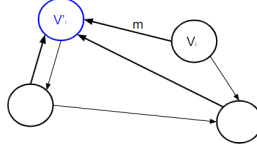


Figure 3: v'_i updates its state using the values received from adjacent nodes.

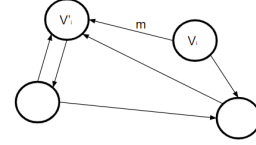


Figure 4: The Graph updates its state.

2.3 RouteNet characteristics

Our state of the art model differs from one standard GNN, as explained before a normal GNN updated the state of a node using the state of the adjacent nodes. Unlike this approach, in our case what we need is a model that has two main characteristics, **paths** and **links**. The links are the edges between two nodes in a graph, the paths are the set of links that make up a routing between a source node and a destination node. As such our model is implemented with the idea that in order to understand the network model the links, which would be the edges in normal graph, become the nodes in said graph, since the links are the ones that carry the delay information.

So, in order to know the state of a link h_{l_i} we need to know the state of the paths $h_{p_{1..n}}$ that use h_{l_i} , and in order to know the state of a path h_p it is needed to know the state of the links h_l that form that path. to sum up:

- h_p depends on the $h_{l_1, \dots, k_{(l_{p_k})}}$ of the path.
- h_l depends on the h_p that contain that link.

These two principals can also be expressed in a more mathematical form:

$$h_{p_k} = g(h_{l_{k1}}, h_{l_{k2}}, \dots, h_{l_{k(|p_k|)}}) \quad (1)$$

$$h_{l_i} = f(h_{p_1}, \dots, h_{p_j}), l_i \in p_k, k = 1, 2, \dots, j \quad (2)$$

Where f and g are a nonlinear system of equations, where the states are hidden variables, each function depends on the routing used, and the dimensionality of each function is is incredibly large.

Moreover the order on which the links are evaluated is essential in order to have a good performing model. If a link is lossy (produces drop packets) this will impact on the state of the coming links. As such a path is composed from a set of ordered links $p_k = (l_0, l_1, \dots, l_n)$, and the delay of a path is indeed the sum of delays on every link of the path, $\sum id(l_k i)$, where $d(l_k i)$ is the sum of the delays on every link. The order from which the paths are calculated is not relevant.

2.4 RouteNet

The presented architecture is able to provide a routing invariant model, which at the same time is able to understand how a routing in a computer network behaves. RouteNet is able to overcome the problems raised by equations (1) and (2), by creating a GNN based on message passing between the links and paths, and updating the states of paths and links (h_p, h_l) using said message passing architecture. The algorithm below explains the architecture of RouteNet.

Data: x_p, x_l, R
Result: h_p, h_l, y_p

```

1 for  $p \in R$  do
2    $h_p = [x_p, 0, 0, \dots, 0];$ 
3 end
4 for  $l \in N$  do
5    $h_l = [x_l, 0, 0, \dots, 0];$ 
6 end
7 for  $t$  from  $[0, 1, \dots, T]$  do
8   for  $p \in R$  do
9     for  $l \in p$  do
10       $h_p^t = RNN_t(h_p^t, h_l^t);$ 
11       $m_p, l^t + 1 = h_p^t;$ 
12    end
13     $h_p^t + 1 = h_p^t;$ 
14  end
15  for  $l \in N$  do
16     $m_l^{t+1} = \sum_{p: l \in p} m_p, l^t + 1;$ 
17     $h_l, l^{t+1} = U_t(h_l^t, m_l^{t+1});$ 
18  end
19 end
20  $y = F_p(h_p)$ 

```

The algorithm begin by initializing the states of both the paths and the links to 0, then proceeds to iterate T times, in order to achieve convergence, over the state vectors h_p , and h_l . This allows to reach a fixed point of a function from the initial states, which solves the problem of defining implicit functions. In order to overcome the second issue with routing invariance, the proposed model allows to perform message passing which unites topology and state vectors representation. This is done through the message passing solution, which are stated in lines 11 and 16, that allows links and paths to exchange information.

The message sent from the paths is the information of the current state h_p extracted by the RNN , which used the current state of the path and the state of the links that the path uses (defined in line 9). The use of the RNN is because there exists a sequential dependence between the links on a path, which allows to propagate the information about the losses in a path. While the message from the links is the sum of h_p that go through that link, we can use a simple summation in this case due to the order of the paths not being important, finally the state of the links for the next step is extracted from using another RNN from the current h_l and the received sum of the path messages, as stated in line 17.

Finally, the calculation of the predicted variable y_p is a function of the link and path features x_l, x_p , the predicted variable is extracted with a neural network.

In the figure below we can see a more graphical representation of the inner workings of the neural network. In the figure it begins with the addition of the

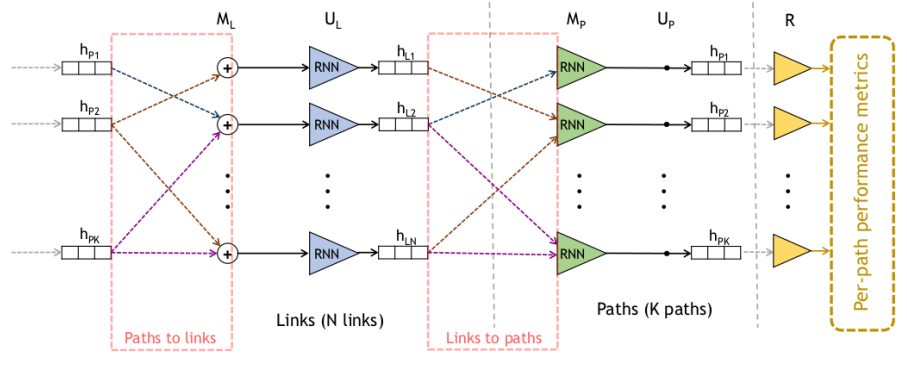


Figure 5: RouteNet Architecture.

states h_p being added for each link that contains them, then the state of the link h_l is calculated using an *RNN*, with the value of h_l from the output of the *RNN*, it we proceed to calculate the h_p using another *RNN*. The output of the h_p is directly the value of the message from that path.

3 Methodology

Ficar mes
cozes aqui?

3.1 Data Gather

The data generated is completely synthetic, the topologies are computer generated, and the simulations performed are packet-by-packets simulations.

3.2 Model Evaluation

In order to train the Model an stratified sampling has been performed in order to split the data between training and evaluation. An stratified sampling is a method for sampling a population which allows to partition said population into sub-populations.

In this case the stratified sampling is done over the value of λ , which will be explained later. The split between train and evaluation is 80% training, and 20% evaluation.

4 Dataset and Technical Details

4.1 Introduction

In the dataset section I will go through how the data is gathered, simulations are performed, which different techniques have been used to generate the data, and how this data generation impacts the model.

The hypothesis we are trying to prove is that by creating an algorithm that us able to generate a multitude of routings which are a representative samples of routings. Then using the model of RouteNet as prof that we are able to generalize any kind of network in order to calculate the delay. Finally I will explain in more detail the technical implementation of RouteNet.

4.2 Dataset

4.2.1 Data Gather - Simulations

In order to obtain the dataset to train with Routenet it is first needed to simulate a real network environment, to do so the simulation tool Omnet++ **CITA** is used to simulate networks and produce the desires dataset. Each simulation is done in a given routing by four different traffic intensities λ , each simulation is repeated 500 times in order to achieve a proper mean value for the delays in each path. In our case all nodes are connected to all the others nodes in the network. More over each simulation used to take 16050 seconds. As you can imagine this meant that the simulations took a long time to complete.

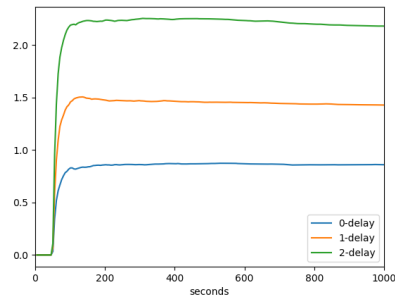


Figure 6: Evolution of delay over time for 50 nodes for each traffic intensity λ .

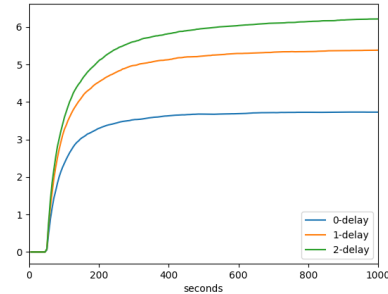


Figure 7: Evolution of delay over time for 100 nodes for each traffic intensity λ .

It became quite clear for us (myself and the research team) that we were simulating way more time than what was necessary in order to achieve an stable value of a given delay, as such it was decided to cut the simulation time in order

to achieve a more reasonable time for the simulations.

Once the simulation time was reduced to a reasonable amount of time to perform correct simulations, we were able to get the necessary data to create the model.

4.2.2 Dataset

In order to create a sizable data sample simulations were performed on three different network graphs, with different sizes, each with 14, 24, and 50 nodes. In the figures below we have the different topologies.

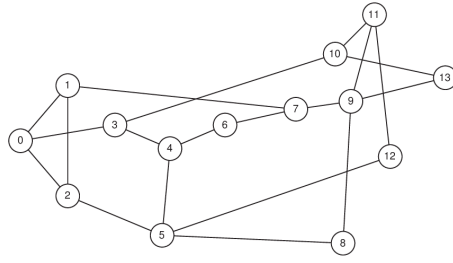


Figure 8: Network graph with 14 nodes.

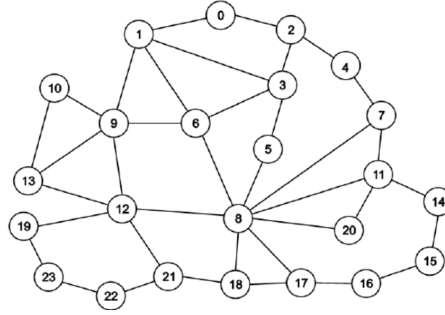


Figure 9: Network graph with 24 nodes.

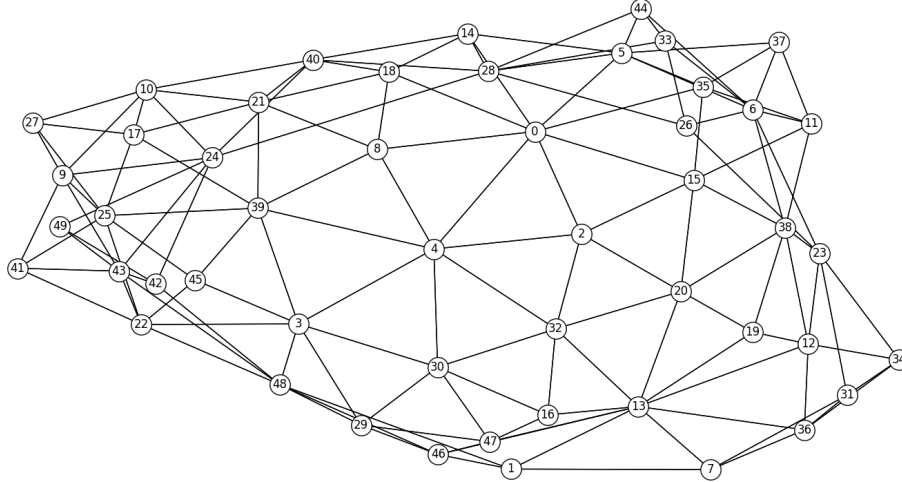


Figure 10: Network graph with 50 nodes.

These three topologies are synthetic (meaning that they do not come from

real networks but instead are computer made) [7] [8]. For each topology shortest path routings algorithms were performed for every pair of nodes, meaning that one given node in a network graph of N has a total of $N - 1$ routings, since no self loops are allowed. The reason why no self loops are allowed is because in a network environment they make no sense. This means that a network topology has a total of $N * (N - 1)$ routings in total. In the following table we can see how many routings have each of the network topologies that we are using in this research. This means that we in theory have 3 different datasets, one for each topology type.

	14 Nodes	24 Nodes	50 Nodes
Routings	212	552	2450

Table 1: Number of routings per topology

4.2.3 Routing Creation

In order to perform the shortest path algorithm a *Deep-first search*, which was presented by *Charles Pierre Trémaux* as an strategy for solving mazes, the version of the algorithm used is the one explained in the book Introduction to algorithms [9], which establishes a cutoff value in order to keep the algorithm to a reasonable time performance. The idea is to use the algorithm to search all the possible paths shorter than the cutoff value, from one source to a destination. The cutoff value has always a base value of the shortest path distance between the two nodes The algorithm works in the following manner:

Ficar-ho
aqui o al
state of the
art?

Data: The topology as G , s as a vertex of origin, and d as the destination, cut as the cutoff value

Result: Routings for source - destination pairs

```

1 Function  $DSF(G, s, d, cut)$  is
2   let  $S$  be a stack of all the adjacent nodes to  $S$ ;
3   let  $V$  be the visited nodes;
4   while  $S$  is not empty do
5      $v = S.pop()$ ;
6     neighbour = Adjacent nodes of  $v$ 
7     if  $c$  is empty then
8        $S.pop()$ ;
9        $V.pop()$ ;
10    end
11    else if length  $V \geq cut$  then
12      if neighbour is the  $d$  then
13        return  $V + d$ ;
14      end
15      else if neighbour not in  $V$  then
16         $V.push(neighbour)$ ;
17         $S.push(Adjacent\ nodes\ of\ neighbour)$ ;
18      end
19    end
20    else cutoff value reached
21       $S.pop()$ ;
22       $V.pop()$ ;
23    end
24  end
25 end

```

This algorithm is able to generate a full valid routing for the whole network topology in under 1.5 seconds, which means that we are able to generate quite a big amount of routings in a short span of time. When I refer to a valid routing it means a routing that does not have loops in the routing procedure. It is important to understand that in networking routing policy is usually destination only, meaning that the source node does not have any impact on the route to use, in other words, when a IP packet arrives on a network node, it only takes into consideration the destination of the packet and forwards it to the node port specified in the IP table. This can lead to loops in the network, for example: Here we have an small topology consisting of 4 nodes, lets say we have 2 routing, one having the source node as node 0, and destination node as node 4, and the other routing with source 2 and destination 4. The path for each routing are:

$R1 : Node0 -> Node1 -> Node3 -> Node4$

$R2 : Node2 -> Node3 -> Node1 -> Node4$

In this case we have a loop, since the nodes only care for the destination of the packet, in $R1$ once the packets gets to node 1, it will be sent to node 3,

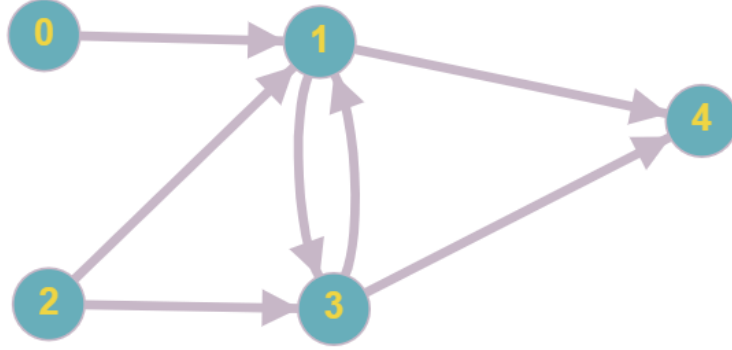


Figure 11: Example of routing with loops.

but due to *R2* establishing that packets that arrive to Node 3 with destination Node 4 have to go to Node 1, the packets will be sent once again to Node 1, which in turn will send it once again to Node 3, thus creating a loop.

Moreover the algorithm take a parameter which we will call α , that indicates the algorithm how it should attempt to place the paths on the graph. At the moment the α parameter takes three possible values, *easy*, *normal*, and *hard*.

	Routing placement
easy	Sparse placement
normal	Random placement
hard	Overlapping placement

The algorithm works by assigning weights on each link, depending on the α chosen the algorithm will choose paths that contain the links with less, more, or random weight. This allows us to generate a representative sample of routing, with different values of delay and bandwidth.

With this algorithm we are able to generate a multitude of algorithms that explore the topology in a bigger number of ways than an standard shortest path algorithm would, thus creating the basis for the main objective of this thesis, the theory that in order to be able to predict the delay of any model the number of nodes on the topology is irrelevant, but instead the data should be a sample representative enough of routings.

Due to time constrains, we were not able to implement the proper simulation policy, which works with a source-destination model for implementing the

routing algorithm. This causes that the algorithm is only able to create combinations of shortest path, which is not ideal, but is enough to create an small representative sample.

Once the routing have been created, each is simulated for different values of traffic intensities, called λ , these λ represent the amount of traffic that is produced by each source destination pair. In table 2 .

λ	8	9	10	11	12	13	14	15	Total
14 Nodes	156		174		91			174	595 aprox 3.8GB
24 Nodes	174		170		172			174	693 aprox13.3GB
50 Nodes		200			200			200	600 aprox 42.8GB
50 Nodes Alternative		139			138			138	414 aprox 36.3 GB

Table 2: Simulations per lambda and nodes

The difference between *50 Nodes* and *50 Nodes Alternative* are the way the routing are created. The routing for 14, 24, and 50 Nodes are standard SP routing, while the ones created for *50 Nodes Alternative* are the ones we aim to prove our theory in the thesis, which is creating a representative sample enough of routings in order to be able to predict any other routing regardless of the topology size.

Each combination of routing per λ is then simulated 500 times for the 14, 24, and 50 Nodes, and 300 times for the 50 Nodes alternative. In the end the total amount of samples for each topology is the the total number of routings, multiplied by the total number of λ , multiplied by the total amount of simulations performed for each combination of routing and λ .

	14 Nodes	24 Nodes	50 Nodes	50 Nodes Alternative
Samples	297.500	346.500	300.000	124.200

Table 3: Total amount of samples per each routing.

4.2.4 Traning and Evaluation

As explained in the methodology the split between training and evaluation is done through and stratified sampling using the λ as the parameter to create the subpopulations from. The split for train and evaluation is 80% and 20 %.

The Model is **only** trained using the 50 Nodes Alternative data. The other samples are used as evaluation only. The table blow shows the split in sample numbers for traning and evaluation.

Samples	Training	Evaluation
	99.360	24.840

Table 4: Training and evaluation split.

4.2.5 Data Features

The dataset is composed of quite a large number of features. A single sample is a whole topology with all the delays for each source destination pair routing, meaning that, for example, a Routing with $N = 50$, has a total amount of 49×50 *delays*, *bandwidth*, *jitter*, *links*, *Number of packets*, *drops*, *sequences*, and *paths*.

- **Delays:** Target feature. Its a continuous variable containing the delay of each path in the network.
- **Bandwidth:** Also known as a traffic matrix, a continuous variable which is the set of bandwidth that each path has.
- **Jitter:** Variation of the delay. Can also be the target variable, in the case it is not the target, it is not used.
- **Drops:** Continuous variable. Number of packets lost for each source destination combination.
- **Number of packets:** Continuous variable. Number of packets sent between source and destination.
- **links:** Number of hops that the path takes from source to destination. (Calculated at runtime)

4.2.6 50 Nodes Alternative

The figures below indicate some of the statistics of the dataset for 50 Nodes Alternative, which is the dataset used for the training of the model. The figures do not use all the data in the dataset, but only use 1800 samples. Figures 12, 13, 14, 15 show the distribution of values for the bandwidth, delay, packets, and drop packets.

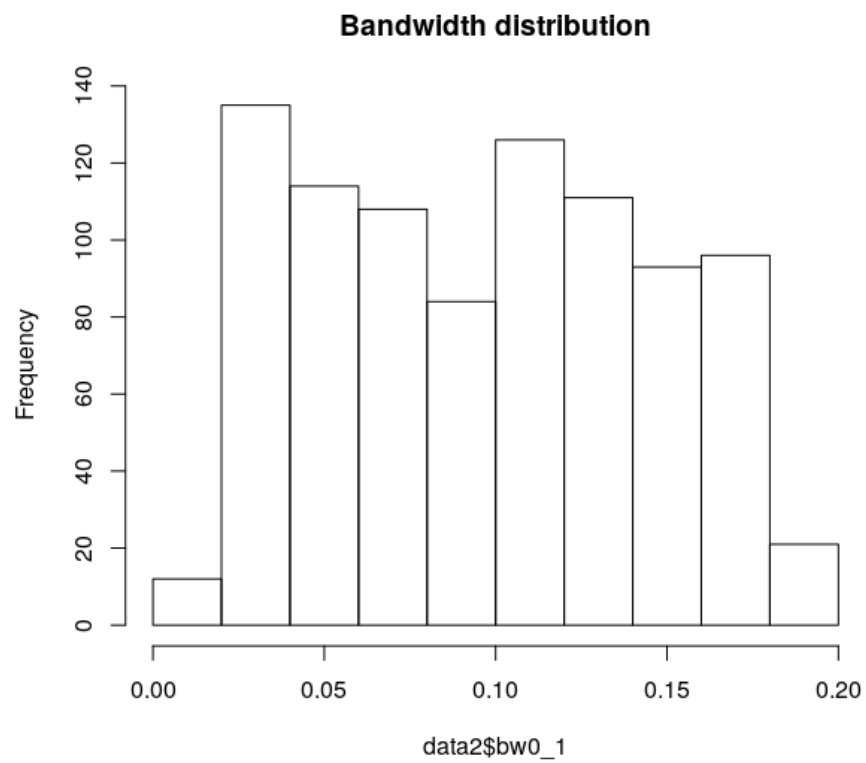


Figure 12: Histogram of bandwidth distribution

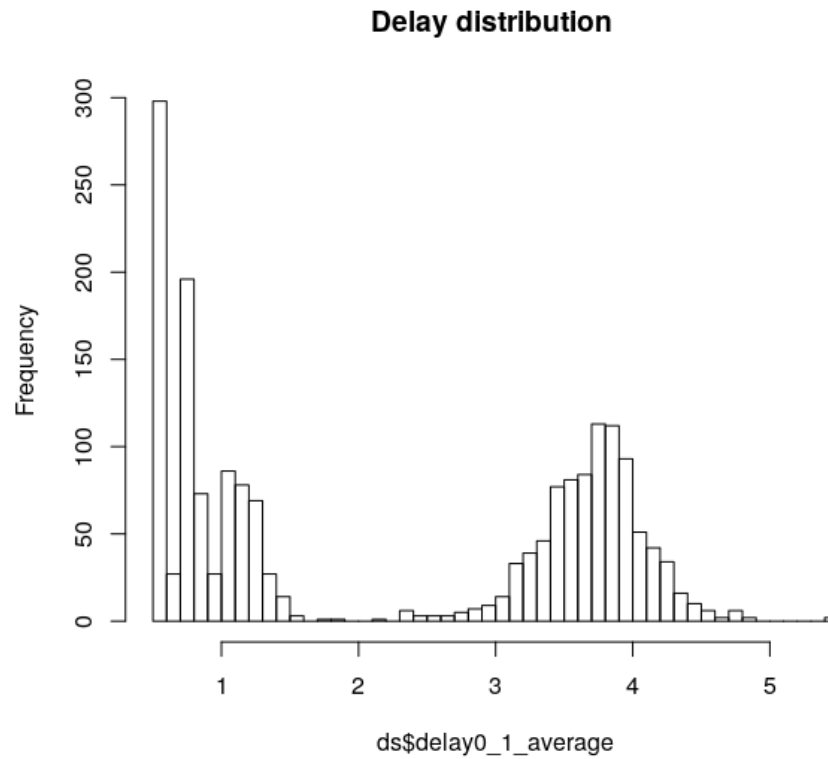


Figure 13: Histogram of delay distribution

This is an important figure, delay is indeed the target variable, so the distribution of values that our dataset has is quite important in order to know which range of values RouteNet is able to predict, values outside the range shown in the histogram will be hard to predict for the model.

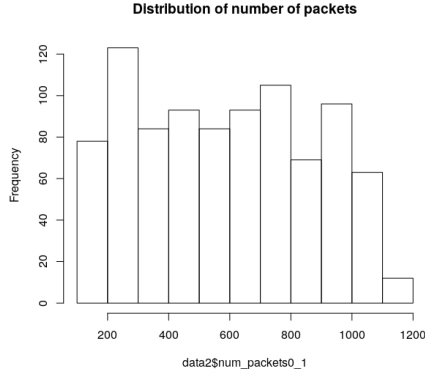


Figure 14: Histogram of packets distribution

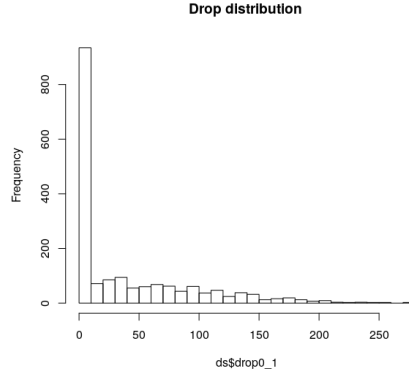


Figure 15: Histogram of drops distribution

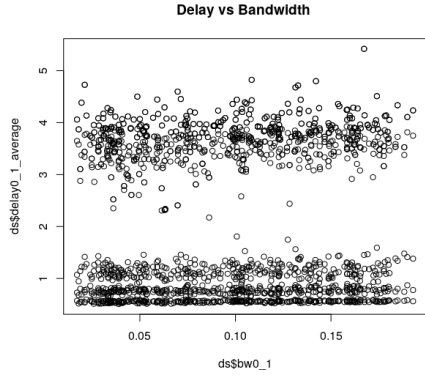


Figure 16: Delay vs Bandwidth

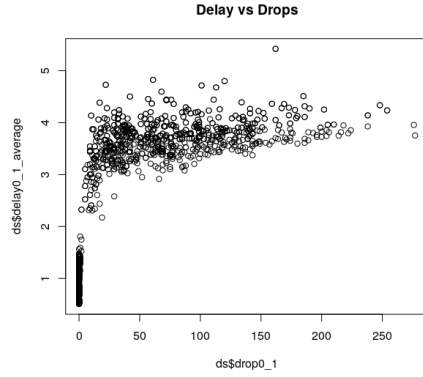


Figure 17: Delay vs Drops

Figure 15 indicates that we are losing a few packets during the network transmissions. Figure 16 shows the relation between the delay and the bandwidth, we can see how the delay is not affected much as the bandwidth increases. If we look at figure 17, we can see how the delay is indeed affected by the number of drops in the network. The data seems to show that the delay vs drops packages follows a logarithmic function.

The table below 5 shows the correlation between each feature. We can see how the bandwidth, number of packets, and dropped packets all have some degree of correlation to the delay.

Finally the figure below shows some statistics for the dataset.

	BW	# Packets	Drops	Delay
BW	1.0000000	0.9982516	0.6240850	0.1521073
# Packets	0.9982516	1.0000000	0.6251060	0.1488149
Drops	0.6240850	0.6251060	1.0000000	0.3540376
Delay	0.1521073	0.1488149	0.3540376	1.0000000

Table 5: Correlation between features for the 50 Nodes Alternative dataset

bw0_1	num_packets0_1	drop0_1	delay0_1_average
Min. :0.01713	Min. : 100.0	Min. : 0.00	Min. :0.5067
1st Qu.:0.05292	1st Qu.: 321.0	1st Qu.: 0.00	1st Qu.:0.7300
Median :0.09964	Median : 596.5	Median : 0.00	Median :1.1557
Mean :0.09744	Mean : 585.8	Mean : 31.98	Mean :1.9595
3rd Qu.:0.13815	3rd Qu.: 840.2	3rd Qu.: 51.00	3rd Qu.:3.5673
Max. :0.19293	Max. :1121.0	Max. :277.00	Max. :5.4211

Figure 18: Dataset Summary for 50 Nodes Alternative

4.2.7 50 Nodes

Although the dataset from 50 Nodes and 50 Nodes alternative might seem similar at first glance, but the reality is that the only thing that both datasets share is the same topology. The main difference between the two is that the data for the 50 Nodes dataset has been scaled by a factor that takes into consideration the size of the network. This procedure reduces drastically the traffic in the network, put it also produces no packets dropped.

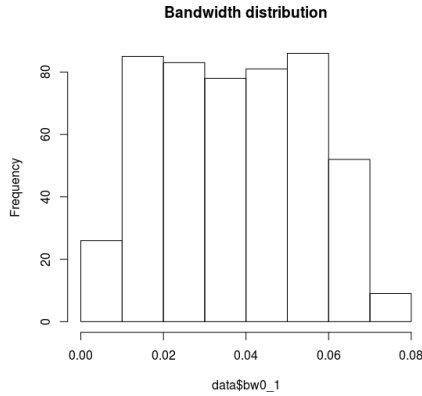


Figure 19: Distribution of the bandwidth for dataset with 50 nodes.

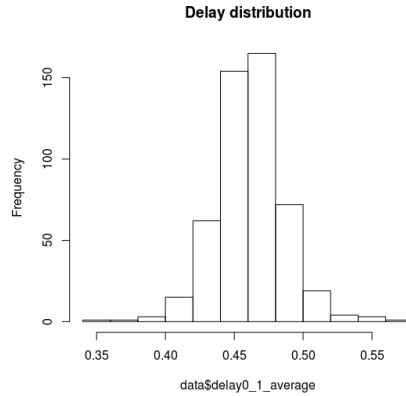


Figure 20: Distribution of the delay for dataset with 50 nodes

If we compare the bandwidth histogram for 50 Nodes Alternative dataset 12

and the one for the current dataset 19 we can clearly see how the bandwidth for the current dataset is approximately a magnitude 10 times smaller than the one for the 50 Nodes Alternative dataset

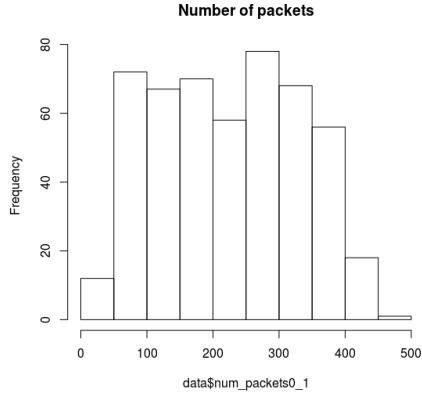


Figure 21: Distribution of the packets for dataset with 50 nodes.

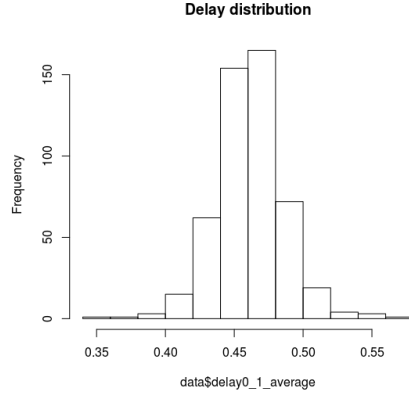


Figure 22: Delay vs Drop for 50 Nodes

If we compare again the number of packets in the network 14 21 it is also clear how there are much more packets in the dataset 50 Nodes Alternative than the 50 Nodes dataset. This clearly shows that there is less traffic in the network.

4.2.8 24 Nodes

The below figures are the statistics for the topology with 24 nodes, as with the data from the 50 nodes alternative dataset, we are not using all the data samples, but are only part of it. Also, the dataset of the topology for 24 nodes does not contain the same information as the other topologies due to a change in the data format that occurred at the end of the thesis period, the main difference is that there is not attribute for the total amount of packets from source - destination.

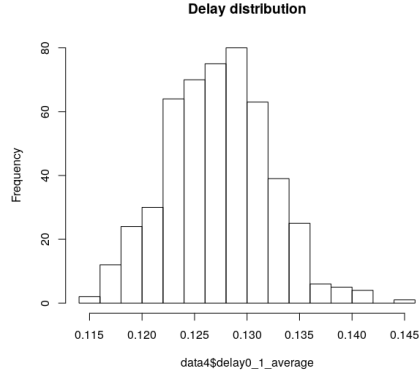


Figure 23: Distribution of the delay for dataset with 24 nodes.

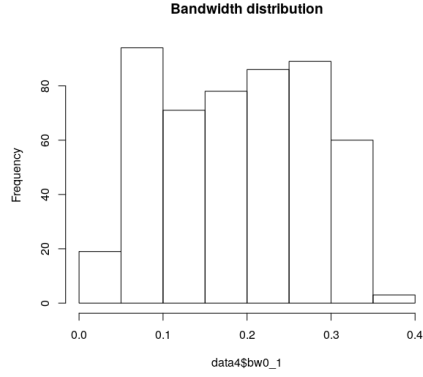


Figure 24: Distribution of the bandwidth for dataset with 24 nodes

The main difference that we find in this subpart of the dataset compared to the figures from the train data is the values that the delay achieves, if we look at figure 13 and compare it to figure 23, it is clear that the topology of 24 nodes achieves lower rates of delay compares to the ones from the 50 Nodes. This is believed to be a cause of more traffic going through each link in the topology of 50 nodes, since we have way more routings, also, the fact that more hops are needed due to the diameter of the graph being higher is also a cause of higher delay in the network.

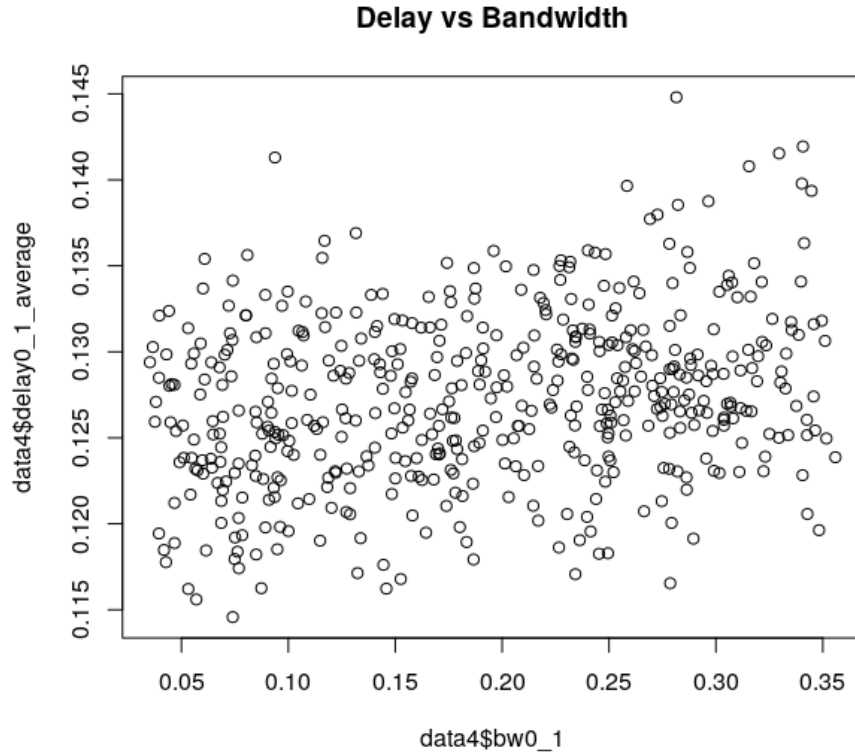


Figure 25: Delay vs Bandwidth for dataset with 24 nodes.

With the delay vs bandwidth there seems to be no clear relation between the two, but it seems like the delay increases a bit as the bandwidth goes up. The drop vs bandwidth plot is not shown due to there not being any packets drops in the simulation.

4.2.9 14 Nodes

Finally the topology with 14 Nodes, due to this being the topology with the least amount of nodes, it is also the topology with the lower amount of delay in the system. The delay vs bandwidth 26 also seems to indicate an small increase of the delay when the bandwidth increases. Figure 27 also indicates how the delay spikes when the dropped packets increase.

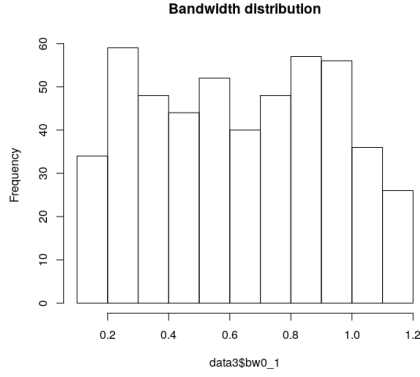


Figure 26: Distribution of the bandwidth for dataset with 14 nodes.

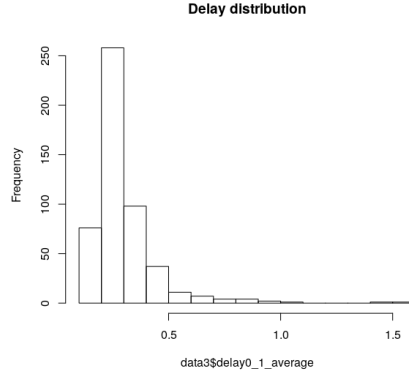


Figure 27: Distribution of the delay for dataset with 24 nodes

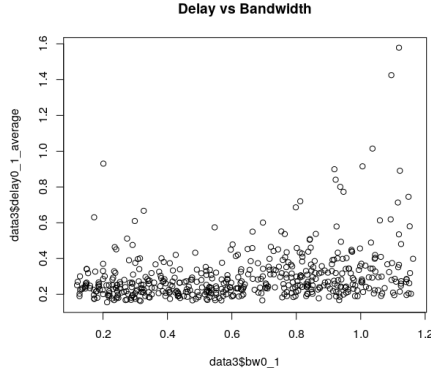


Figure 28: Delay vs Bandwidth for dataset with 14 nodes.

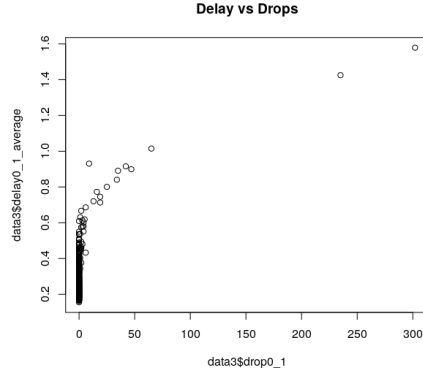


Figure 29: Delay vs Drop for dataset with 14 nodes.

4.3 Architecture

Routenet is implemented using the Tensorflow library. Tensorflow is, as described by Tensorflow creators, "An end-to-end open source platform for machine learning. It has a comprehensive, flexible ecosystem of tools, libraries and community resources that lets researchers push the state-of-the-art in ML and developers easily build and deploy ML powered applications" [10].

Tensorflow eases the creation of machine learning models, by providing efficient and custom model creation, implements performance metrics, and provides with an easy to use tool to check the results of the models created in the form of

tensorboard [11].

4.3.1 Recurrent Neural Networks

The model is composed of two Recurrent Neural Networks (RNN). Recurrent Neural Networks have been extensively used for multiple task which require contextual information in order to understand the current input of the Neural Network, some of the examples of RNN usage are: text processing, handwriting recognition or speech recognition [12] [13]. Graph Neural Networks use RNN due to the fact that they are able to retain information from previous states in order to update the current state of the hidden neuron, this mechanism allows to carry a bias for locality in the sequence via their Markovian structure.

4.3.2 Gated Recurrent Unit

The RNN are used one for the path and the other for the links, the RNN gating mechanism chosen for this implementation are Gated Recurrent Units (**GRU**) [14], the use of GRU over other mechanisms such as Long Short Term Memory (LSTM), is mainly due to achieving similar performance while using less parameters and thus being simpler than an LSTM gate [15] [16]. The need for using this kinds of gates is due the whats called *vanishing gradients* problem, this occurs when a gradient shrinks as the backpropegation mechanism proceeds, if a gradient value becomes small after some training, the model stops learning, and due to how the backpropagation mechanism works, small gradients usually appear in the beginning layers of the Neural Network, thus, stopping all learning of the model. Figure 30 shows an example of vanishing gradients.

new weight = weight - learning rate*gradient

$$\boxed{2.0999} = \boxed{2.1} - \boxed{0.001}$$

Not much of a difference update value

Figure 30: Gradient vanishing problem

In the figure below (31) we can see the inner structure of a GRU gate.

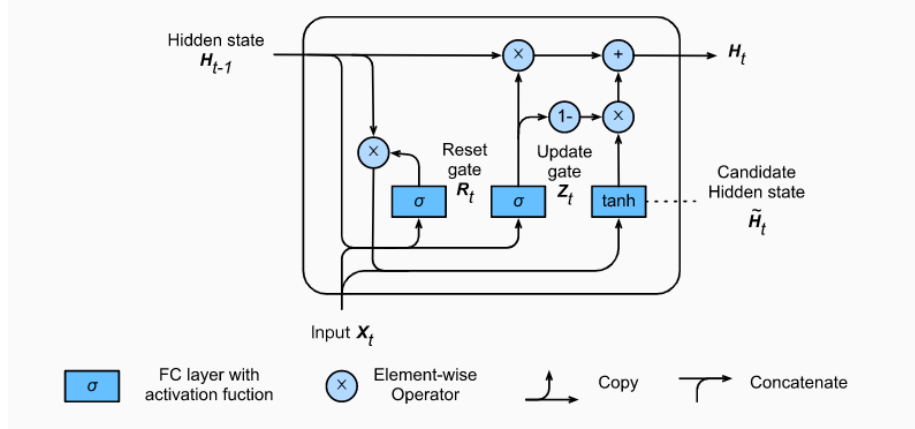


Figure 31: Gated Recurrent Unit [17]

Both the reset gate and the update reset gates are both *sigmoid* functions, which squish the value that go through it between 0 and 1, this allows the neural network to keep or forget data, a value which is converted to 0 is forgotten, while the others value become more or less relevant as time goes on. This mechanism allows the gate to keep or forget important and unimportant data. The *tanh* activation function is used to regulate the flow of data in the network, the output of the *tanh* function are the values squished between -1 and 1.

A quick summary of what each gates does, the update gate helps the network how much of the past information needs to be passed along to the next steps. Similarly the reset gate is used how much of the past information will be forgotten.

4.3.3 Readout function

The readout function is an standard neural network with two hidden layers which uses a *Scaled Exponential Linear Units* or *SELU* for short, the main characteristic of SELUs are the internal normalization [18] and also for solving the vanishing gradients problem. The readout neural network also uses an **L1** regularization, also know as Lasso regularization, which is useful not only as a regularization procedure in order to reduce the weights of the features that are not important, but also as a feature selection procedure in case there is a feature that is not interesting for the model. Choosing an L1 regularization is due to mostly not performing a feature selection beforehand, and as such it is performed implicitly on the model [19].

The readout functions are mixed with two dropout layers, in the training phase,

they prevent overfitting, while during the inference, they are used for Bayesian posterior approximation [20]. This is an important feature of the model as it allows us to assess the confidence of network prediction, what this means is that when a Neural Network is trained and optimized for a specific output, the solution of the optimization can be for which a network is too optimistic, meaning that a normal neural network without the dropout layers does not take into consideration the uncertainty that another Bayesian model could achieve, yet due to the infeasibility due to the computational power required to train such model.

Quoting [20] “Using a dropout layer allows to extract the necessary information from the model that has been thrown away so far. This mitigates the problem of representing uncertainty in deep learning without sacrificing either computational complexity or test accuracy.”

Repeating inference multiple times with dropout layers being active gives a distribution of results. The spread of this distribution is a measure of network confidence about the prediction and the optimal point.

The figure below shows the model structure as implemented in tensorflow.

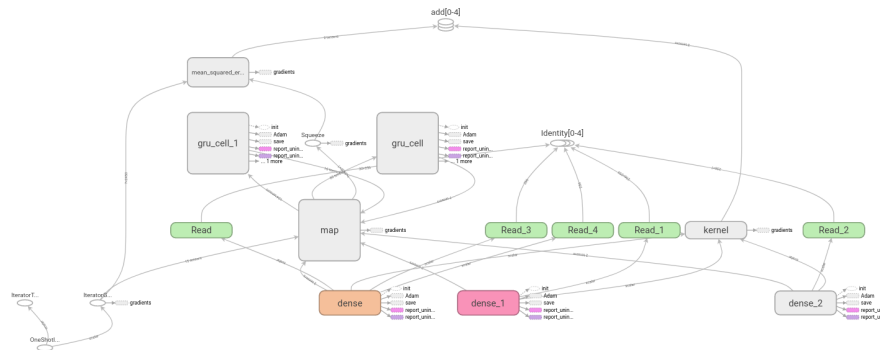


Figure 32: Model Structure

An other important item to choose is the algorithm used in order to minimize the loss function. The loss function is the metric that is going to be used as the performance metric for the model to know if the back-propagation algorithm is working or not, meaning that if the loss function increases, the model is not working properly, but if instead decreases, it means that the weights computed by the back-propagation algorithm are going in the right direction.

Explicar una mica mes això?

4.3.4 Loss minimization function

The algorithm used in order to minimize the loss is the *Adam Optimizer* [21]. Adam is defined by its creator as a combination of both RMSprop and Stochastic Gradient Descent with momentum, both of which are also used as optimizer for neural networks. The attractiveness for using Adam Optimizer are the following [22]:

- Straightforward to implement.
- Computationally efficient.
- Little memory requirements.
- Invariant to diagonal rescale of the gradients.
- Well suited for problems that are large in terms of data and/or parameters.
- Appropriate for non-stationary objectives.
- Appropriate for problems with very noisy/or sparse gradients.
- Hyper-parameters have intuitive interpretation and typically require little tuning.

All of these characteristics are the reasons why Adam Optimizer was chosen over other kinds of optimizers. Adam Optimizers maintains a learning rate that improves performance on problems with sparse gradients and also are adapted based on the average of recent magnitudes of the gradients for the weight.

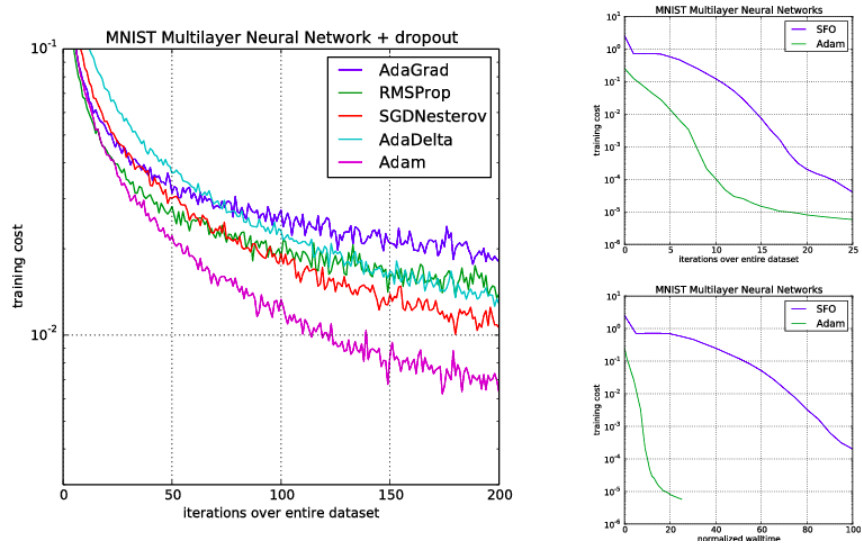


Figure 33: Adam performance compared to other models as described in [21]

4.3.5 Performance Metrics

With the model created it is needed to define some performance metrics, not only to use during the valuation process but to also define which should be the metric that the algorithm should try to optimize, also called the loss, and to choose which method should be used in order to optimize it.

The metric implemented are:

- Mean Absolute Error: which measures the difference between two continuous variables.
- ρ : which is a Pearson correlation, which indicates the extend on which two linear variables are correlated.
- Mean Relative Error: which is a measure of the uncertainty of measurement compared to the size of the measurement.
- r_2 : which is the standard metric to evaluate the goodness of a model.

In our model we have chosen the Mean Absolute Error as the *loss* metric, meaning that this is the metric that RouteNet Adam optimizer will try to minimize, the other metrics will be used as evaluation.

Besides the metrics mentioned there will also be used as metrics for the goodness of the fit the True delay vs predicted delay scatter plot with the error for each point and distribution of points as an histogram. There will also the cumulative distribution function of the mean relative error.

5 Evaluation

5.1 Introduction

In this section I will proceed to explain the results of our experiments and check if the hypothesis stated in the objectives of this thesis are achieved. We will do this by showing the results of the metrics explained in the previous section. I will also evaluate the performance of the model.

5.2 Model Performance

In this section I will evaluate the model performance by showing a variety of graphs that assets how well the model is created.

The figures below show the evolution of the distribution of values that come as an input, figure 34, and the output, figure 35. From the graphs it easily to see

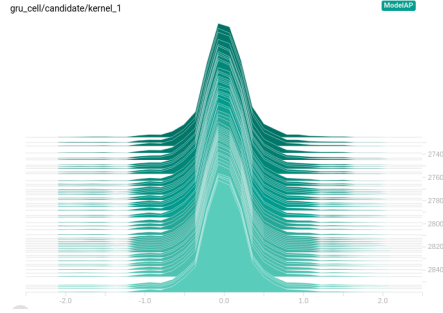


Figure 34: GRU Kernel candidate distribution.

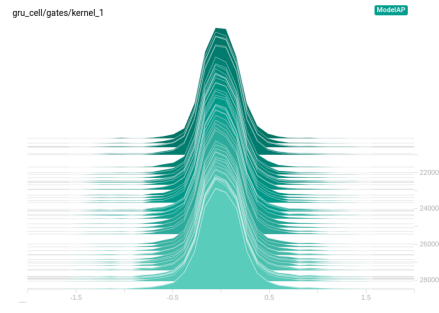


Figure 35: GRU gate value distribution.

how the distribution from the input and the output is similar over time. In the figure below we can see the bias of the GRU.

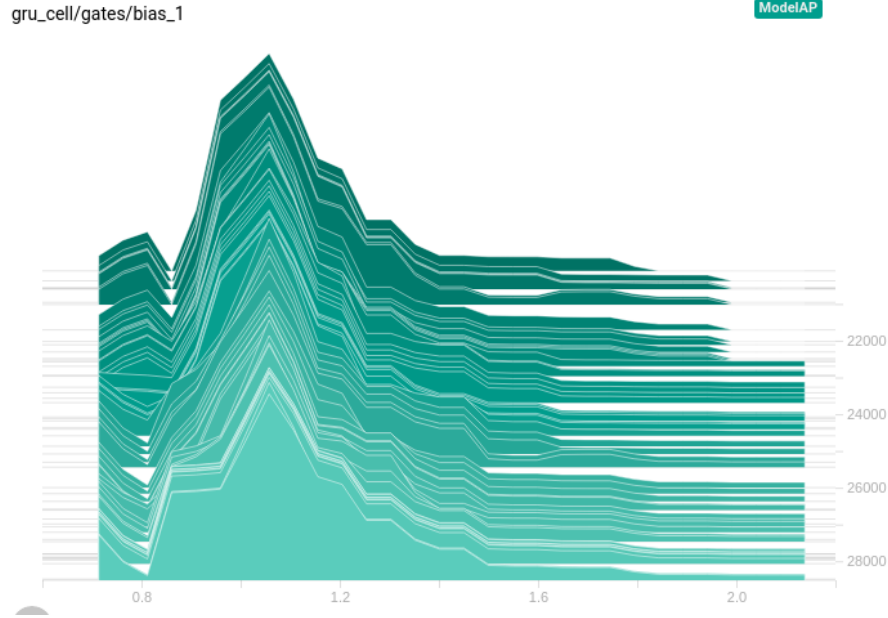


Figure 36: GRU bias for the output.

5.3 Metrics

5.3.1 50 Nodes Alternative

In this section I will discuss the metrics results for our model. I will begin by showing the evolution of the different metrics for the training and evaluation dataset over time, commenting the results. Finally I will show the results for the other topologies.

Figure 37 shows the evolution of the loss (Mean squared error) metric, which is the metric that the model tries to minimize.

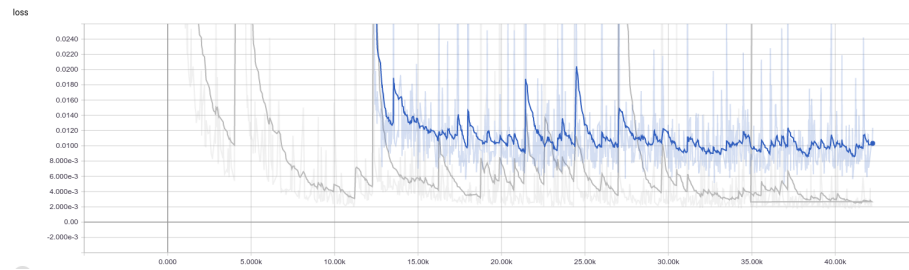


Figure 37: Loss (MSE) for training and evaluation data. The blue line indicates the evolution of the training data, the grey line indicates the evolution of the evaluation data.

The training line starts later than the evolution due to tensorflow only storing the past 30K steps for training. The loss for the training has a final value of *0.01063*.

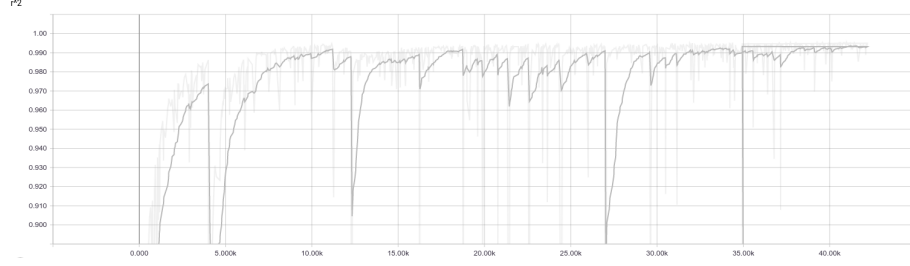


Figure 38: R_2 for training and evaluation data.

The drops seen in the evolution are due to the training not being done all at once, but instead training it for a number of steps, and then restarting the training. The following plots show the evolution of the performance metrics for the evaluation dataset. All the results below are for the part of the evaluation data for the dataset called 50 Nodes Alternative.

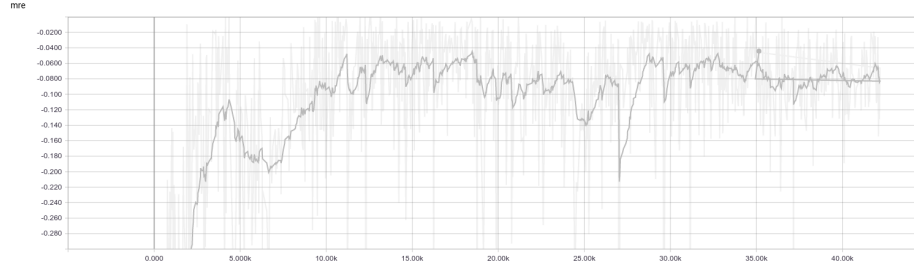


Figure 39: MRE for training and evaluation data.

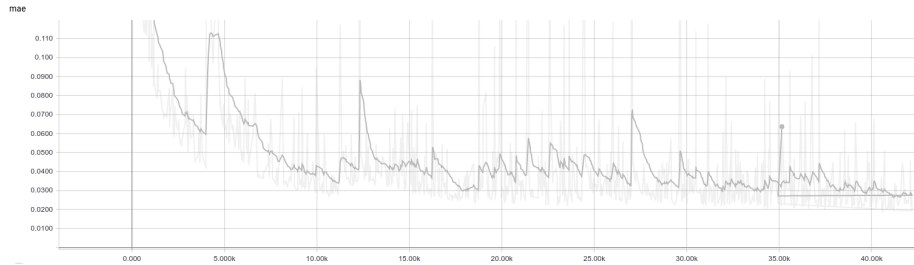


Figure 40: MAE for training and evaluation data.

It is clear that we are achieving a really good model for our first dataset, 50 Nodes Alternative, all metric indicate a low error, the ρ metric shows that

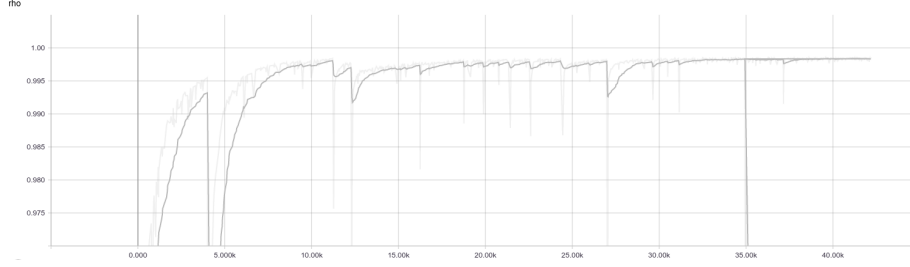


Figure 41: ρ for training and evaluation data.

	R^2	ρ	MAE	MRE	MSE
Evaluation	0.9944	0.9985	0.02685	-0.0786	0.00289

Table 6: Table of results for evaluation dataset

there is indeed a great linear correlation between the predicted delay values and the true delay. The model also has a really high R^2 value, showing that our model explains a great amount of variation of the response variable (delay) is explained by our model. Finally RouteNet achieves a low amount of error of the predictions (mre, mae, mse).

Finally we can use a simple plot of True delay vs Predicted delay for a more visual representation of the accuracy of the results and the model. The histograms on the side of the plots indicate the distribution of the samples.

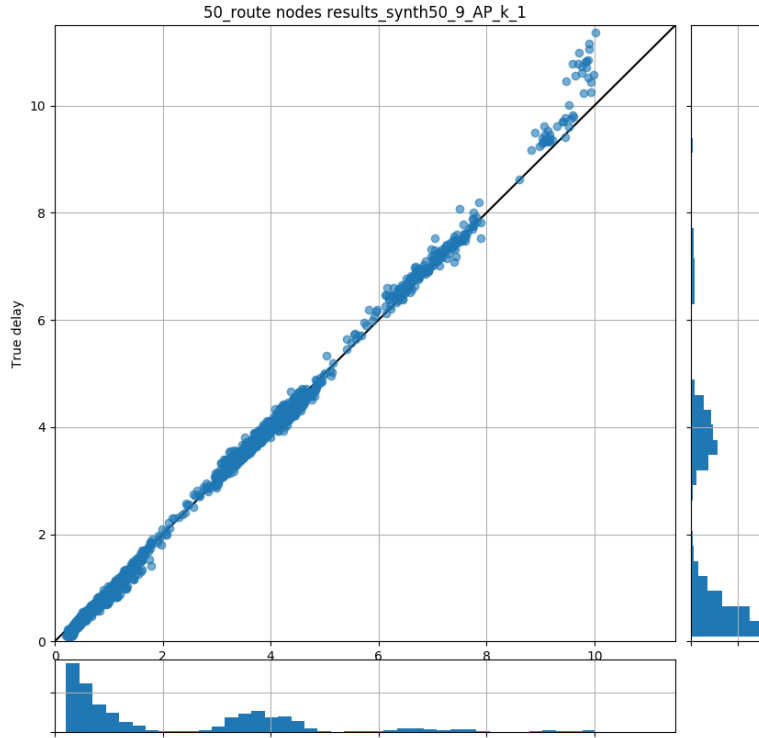


Figure 42: True delay vs Predicted for 50 nodes Alternative.

The model predicts nearly perfectly the values of the delays, the only exception seems to be for really high values of delay, which are not realistic in a computer network.

More True delay vs Predicted scatter plots can be seen in the appendix 8, All these results above are for the dataset for 50 Nodes Alternative, below I will show the results for the rest of the datasets.

5.3.2 50 Nodes

The dataset for 50 Nodes is similar to the one used for training, but as explained before, there is less traffic in the network, this leads to lower values of delay but also way lower values of bandwidth.

50 Nodes	R^2	ρ	MSE	MRE	MAE
	0.81	0.96	0.00414	0.0377	0.047

Table 7: Metric results for 50 Nodes.

The metrics provide good results, we achieve a lower value of R^2 than on the original dataset, but that to be expected.

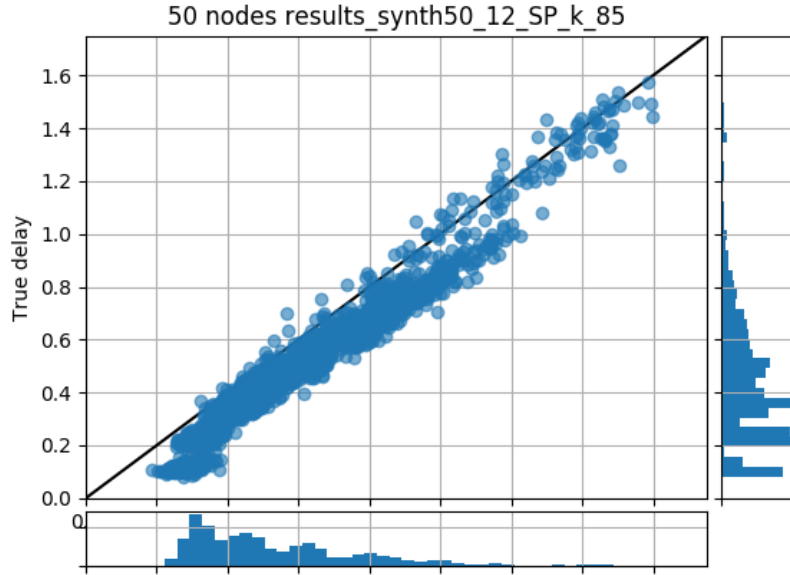


Figure 43: True delay vs Predicted for 50 nodes.

The model seems to predict delays to be somewhat higher than their actual value. Yet, overall the results seem quite good, there are no values which seem far off of the regression line. It appears that RouteNet tends to predict values with higher delay value than the real delay. We can say that the model is able to predict correctly the values for this particular topology.

5.3.3 24 Nodes

Evaluation	R_2	ρ	MAE	MRE	MSE
	0.946	0.9861		0.115	0.65

Table 8: Table of results for 24 nodes.

Explicar resultats 24 nodes

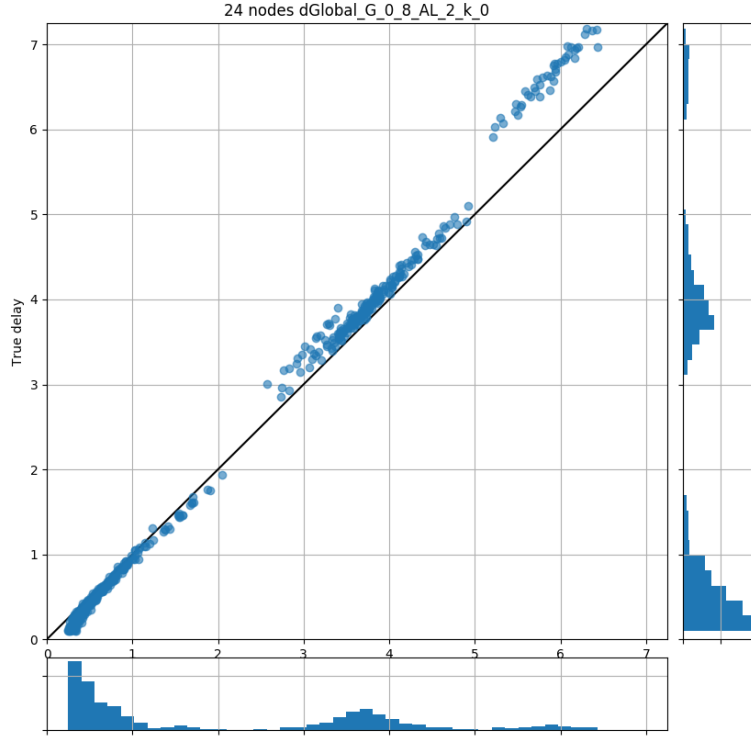


Figure 44: True delay vs Predicted for 24 nodes.

Yet again, we see how the model correctly predicts the values of delay, except for the really high values, which once again, are not realistic.

5.3.4 14 Nodes

The 14 Nodes is an special case and deserves more explanation than the other topologies. The main issue with the 14 Nodes topology is the simulations done with a λ equal to 15. These simulations due to being with a higher value of traffic require more bandwidth to work with no dropped packets, yet the delay is more or less the same as the other simulations with different λ . This causes the model to believe that the delay that should be predicted is way higher than the real one, if we look again at the correlation between variables 5, we see how the bandwidth takes an important role, yet what it is not shown in this table is that in the case where the dataset has no dropped packets, the bandwidth

becomes the most important feature in order to predict the delay.
 In the case of a λ equal to 15, we have no dropped packets at the cost of increasing the bandwidth, thus the model thinks that the predicted delay should be high.

So in order to not have the results for this topology be relevant I will omit the data from *lambda* 15. The table of results with 15 nodes can be found in the appendix on table 10.

	R_2	ρ	MAE	MRE	MSE
Evaluation	0.9504	0.9873		0.09	0.111

Table 9: Table of results for 14 nodes without the data with λ 15.

Explain the metrics

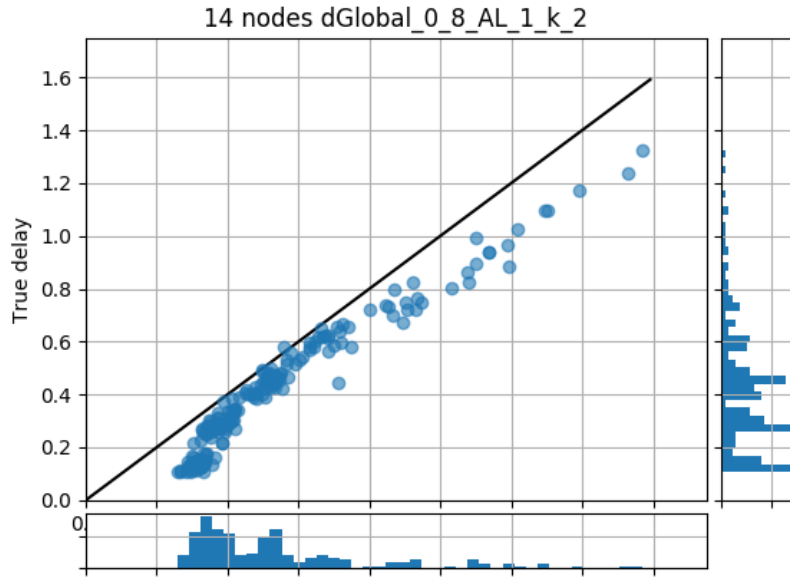


Figure 45: True delay vs Predicted for 14 nodes.

In the case of the 14 nodes, once again RouteNet tends to predict higher values of delay than the real values. But are still reasonable errors. Even more so, if we look at the metrics it is quite clear that the values are good enough to consider the evaluations to be correct. The 14 nodes topology is clearly the one that performs the worst, still I believe that due to the small experiment size, with more training routing samples the issues of over optimistic prediction

would be solved. Still the only issue in the case of this figure 45 are the samples lower than 0.2 seconds of delay, which are a bit off of the regression line, the other values, especially the ones between 0.2 and 0.8 are all close to the regression line. If we look at the previous figures, RouteNet in general seems to have issues trying to predict delays that have a lower values than 0.2 seconds, in both 24 and 50 nodes the predictions don not follow the regression line in the cases where the true delay is lower than 0.2.

As a final figure to show the performance, below we can see a cumulative density function for the mean relative error, which shows that the three datasets used for evaluation are indeed quite similar in terms of the mean relative error.

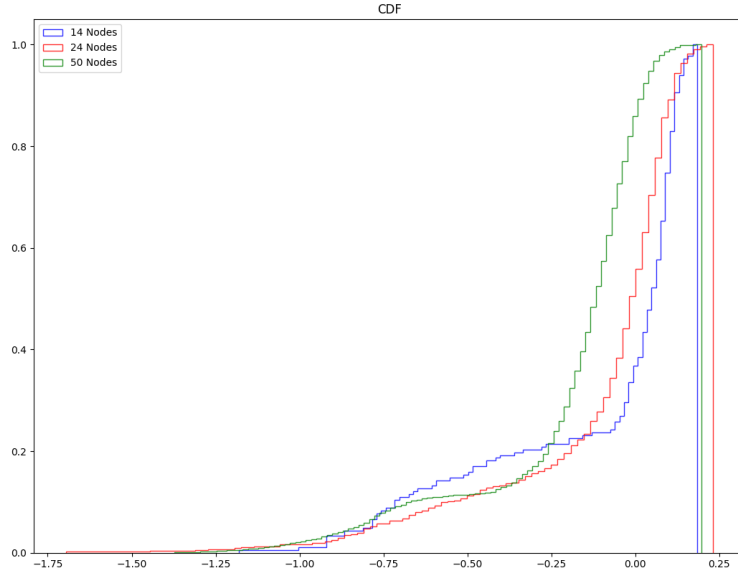


Figure 46: Cumulative density function for mean relative error.

This Figure shows how the mean relative error is accumulated, we can see how the majority of the errors are close to 0, meaning that we are achieving a good, low, error for all models. The evaluation for 14 nodes seems to increase faster than the other topologies, meaning that is most likely the worst performing topology of all of them, yet less than 20% of the values have a significant error, which is still an error percentage that it is acceptable. Still the results are good enough that we can say that our model generalizes correctly.

6 Conclusions

If we look at the initial objectives of the thesis laid down in the objectives section 1.2:

1. Prove that a Graph Neural Network is able to correctly model a computer network
2. Prove that in order to obtain a good model it is not needed to use different topologies sizes as train data.
3. Prove that in order to obtain an accurate Graph Neural Network, only a representative sampling of different routing from one topology is able to create a model which is capable of correctly represent any computer network.

We can clearly see that objective (1) has been accomplished by looking at the evaluation section. The results clearly show that RouteNet is able to model a computer network, achieving good performance results on the metrics and True Delay vs Predicted.

For the second and third objective we can make a comparison between the results of our previous paper presented at *SIGCOMM* and the results from this thesis. The best way to compare them would be to use the results for 24 nodes dataset using the two models.

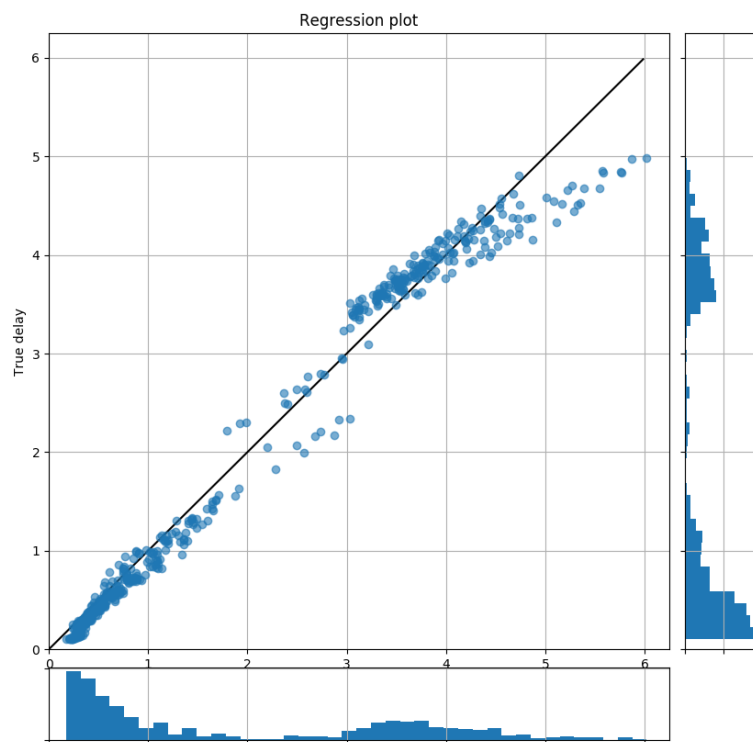


Figure 47: Results from the old model trained with two different topologies and evaluated with the 24 Nodes topology.

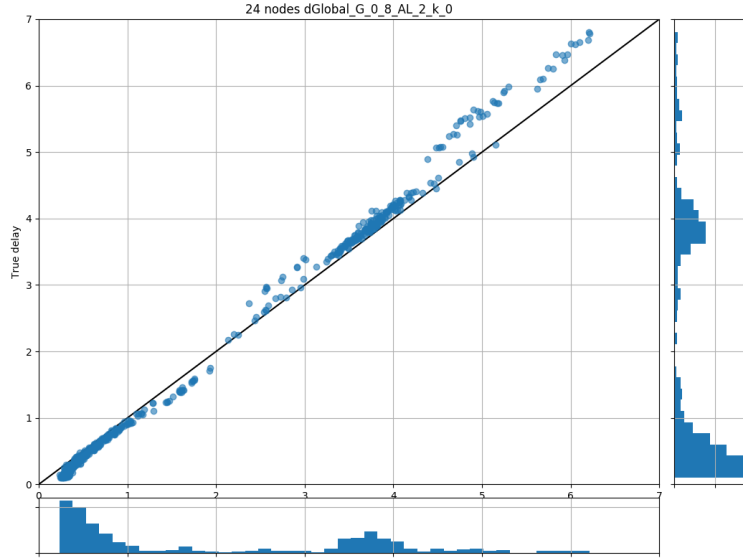


Figure 48: Results for the 24 node topology with the model trained with the alternative routing dataset.

We can see how figure 48 has the data points closer to regression line, this is an indication that the model is better, or at least equal, to the SIGCOMM model, which once again, is trained with 14 nodes, and 50 nodes topology. Even if we achieve equal results to the SIGCOMM model we can still call the thesis a success, since it means that we are not only able to predict the 24 nodes dataset, but we are able to do it with less training data than the SIGCOMM model. If we take a look at the other results from the evaluation, for 14, and 50 nodes, it is also clear how we are able to achieve good performance results in all topologies.

These results prove that we are indeed able to model any network independently of the size by training the RouteNet model with only a routing sample from one topology, yet there are some issues that need to be looked at.

- Not predicting properly really high values of delay.
- Unable to produce good predictions in the case of 14 nodes with λ 15.
- Incorrect predictions for low values of delay (lower than 0.2 seconds).

The issue with values with high delay is the least worrying of all due to these kinds of values not being realistic, I would consider them not relevant for

evaluation proposes. What should be more considered and fixed is the issue with 14 nodes and a high traffic intensity value (λ 15), but the issue should be fixed by generating more data to train the model with.

The other issue that needs to be fixed is the one with RouteNet not predicting correctly the delays lower than 0.2 seconds. This is an issue that should be solved by providing more training data, by taking a look at the summary table for the training data set on figure 25 we can see how the minimum delay is around 0.5 seconds, meaning that the model has not been trained to predict such low values of delay. Adding samples with lower values of average delay would prove very beneficial to RouteNet, and would most likely solve the issue of not being able to predict these kinds of small delays.

Reducing the need to train RouteNet with multiple topologies solves the issue of needing to simulate a quite a high number of topologies and training them. This thesis proves that by only using one topology, but extensively exploring the routing possibilities in a single topology produces enough data to train RouteNet which is then able to model other routings from different topologies.

In conclusion we are able model any kind of network by only training RouteNet with one topology.

As said in the objectives section in the beginning of this document this means that we do not need to simulate multiple topologies with different sizes, as was stated in the SIGCOMM paper, thus reducing the amount of time to train a new RouteNet model.

6.1 Future Work

In order to produce better results, more data would be needed, mainly to work with the edge cases. Also it is likely that changes in the simulation would be needed in order to produce more realistic data, for example, all the traffic in the system is more or less the same. The number of packets is also usually the same also.

Simulations for topologies with a higher amount of nodes would also be needed, the higher the node degree the more paths are possible to generate, thus producing better, and more numerous, training samples.

In order to explore the topologies even more to create more routing combinations, which are needed in order to obtain even better results by creating a model which is able to generalize routings even more, the problem of creating valid routings need to be solved. As of right now all the routings are combinations of shortest path routings, yet it is preferable to create routing schemes that are a combination of shortest path, and non shortest path. To do so the

simulation program needs to be changed to support source-destination routing, which allows to send packets in which the source of the packets is taken into consideration in order to deliver the packet. This would solve the issues of having loops in the routing creation.

7 References

References

- [1] Peter W. Battaglia et al. “Relational inductive biases, deep learning, and graph networks”. In: (2018), pp. 1–40. ISSN: 0031-1820. DOI: 10.1017/S0031182005008516. arXiv: 1806.01261. URL: <http://arxiv.org/abs/1806.01261>.
- [2] Giorgio Stampa et al. *A Deep-Reinforcement Learning Approach for Software-Defined Networking Routing Optimization*. 2017. arXiv: 1709.07080. URL: <http://arxiv.org/abs/1709.07080>.
- [3] Albert Mestres et al. “Knowledge-Defined Networking”. In: *CoRR* abs/1606.0 (2016). arXiv: 1606.06222. URL: <http://arxiv.org/abs/1606.06222>.
- [4] N Wang et al. “An overview of routing optimization for internet traffic engineering”. In: *IEEE Communications Surveys Tutorials* 10.1 (2008), pp. 36–56. ISSN: 1553-877X. DOI: 10.1109/COMST.2008.4483669.
- [5] B Fortz and M Thorup. “Internet traffic engineering by optimizing OSPF weights”. In: *Proceedings IEEE INFOCOM 2000. Conference on Computer Communications. Nineteenth Annual Joint Conference of the IEEE Computer and Communications Societies (Cat. No.00CH37064)*. Vol. 2. 2000, 519–528 vol.2. DOI: 10.1109/INFCOM.2000.832225.
- [6] Justin Gilmer et al. “Neural Message Passing for Quantum Chemistry”. In: (2017). ISSN: 00295981. DOI: 10.1002/nme.2457. arXiv: 1704.01212. URL: <http://arxiv.org/abs/1704.01212>.
- [7] Fernando Barreto. “Fast Emergency Paths Schema to Overcome Transient Link Failures in OSPF Routing”. In: *International journal of Computer Networks Communications* 4.2 (2012), 17–34. ISSN: 0975-2293. DOI: 10.5121/ijcnc.2012.4202. URL: <http://dx.doi.org/10.5121/ijcnc.2012.4202>.
- [8] Xiaojun Hei et al. “Wavelength converter placement in least-load-routing-based optical networks using genetic algorithms”. In: *J. Opt. Netw.* 3.5 (2004), pp. 363–378. DOI: 10.1364/JON.3.000363. URL: <http://jon.osa.org/abstract.cfm?URI=jon-3-5-363>.
- [9] Thomas H. Cormen et al. *Introduction to Algorithms*. 2nd. Cambridge, MA, USA: MIT Press, 2001. ISBN: 0-262-03293-7, 9780262032933.
- [10] 2019. URL: <https://www.tensorflow.org/>.
- [11] 2019. URL: https://www.tensorflow.org/guide/summaries_and_tensorboard.
- [12] A. Graves et al. “A Novel Connectionist System for Unconstrained Handwriting Recognition”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 31.5 (2009), pp. 855–868. ISSN: 0162-8828. DOI: 10.1109/TPAMI.2008.137.

- [13] Xiangang Li and Xihong Wu. “Constructing Long Short-Term Memory based Deep Recurrent Neural Networks for Large Vocabulary Speech Recognition”. In: *CoRR* abs/1410.4281 (2014). arXiv: 1410.4281. URL: <http://arxiv.org/abs/1410.4281>.
- [14] Kyunghyun Cho et al. “Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation”. In: *CoRR* abs/1406.1078 (2014). arXiv: 1406.1078. URL: <http://arxiv.org/abs/1406.1078>.
- [15] Junyoung Chung et al. “Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling”. In: *CoRR* abs/1412.3555 (2014). arXiv: 1412.3555. URL: <http://arxiv.org/abs/1412.3555>.
- [16] 2019. URL: <http://www.wildml.com/2015/10/recurrent-neural-network-tutorial-part-4-implementing-a-grulstm-rnn-with-python-and-theano/>.
- [17] 2019. URL: https://www.d2l.ai/chapter_recurrent-neural-networks/gru.html.
- [18] Günter Klambauer et al. “Self-Normalizing Neural Networks”. In: *CoRR* abs/1706.02515 (2017). arXiv: 1706.02515. URL: <http://arxiv.org/abs/1706.02515>.
- [19] F. Santosa and W. Symes. “Linear Inversion of Band-Limited Reflection Seismograms”. In: *SIAM Journal on Scientific and Statistical Computing* 7.4 (1986), pp. 1307–1330. DOI: 10.1137/0907087. eprint: <https://doi.org/10.1137/0907087>. URL: <https://doi.org/10.1137/0907087>.
- [20] Yarin Gal and Zoubin Ghahramani. “Dropout as a bayesian approximation: Representing model uncertainty in deep learning”. In: *international conference on machine learning*. 2016, pp. 1050–1059.
- [21] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. 2014. arXiv: 1412.6980 [cs.LG].
- [22] Jason Brownlee. *Gentle Introduction to the Adam Optimization Algorithm for Deep Learning*. 2019. URL: <https://machinelearningmastery.com/adam-optimization-algorithm-for-deep-learning/>.

8 Appendix

8.1 True Delay vs Predicted 50 Nodes alternative

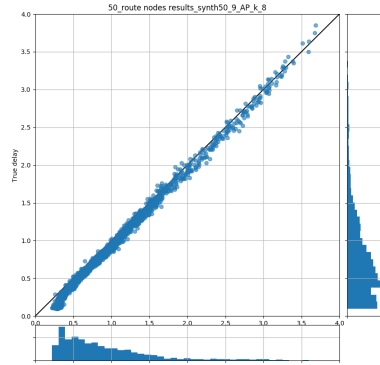


Figure 49: λ 9 routing 8

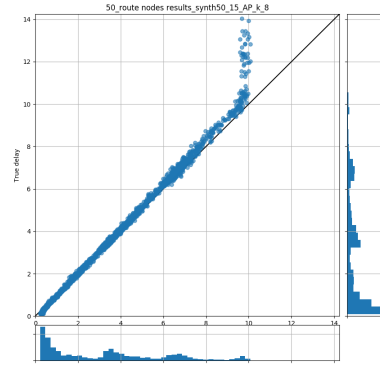


Figure 50: λ 15 routing 8

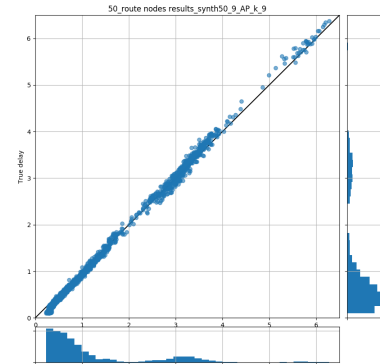


Figure 51: λ 9 routing 9

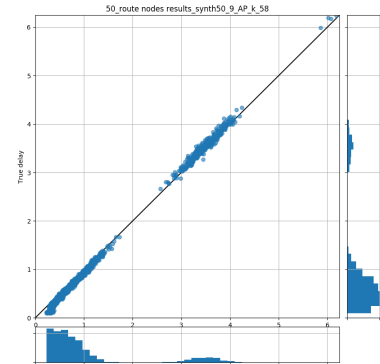


Figure 52: λ 9 routing 58

8.2 True Delay vs Predicted 50 Nodes

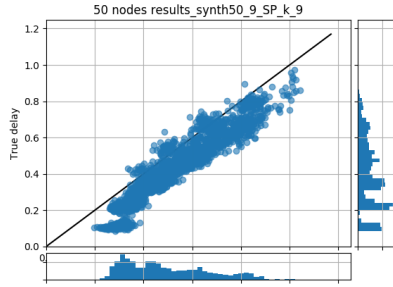


Figure 53: λ 9 routing 9

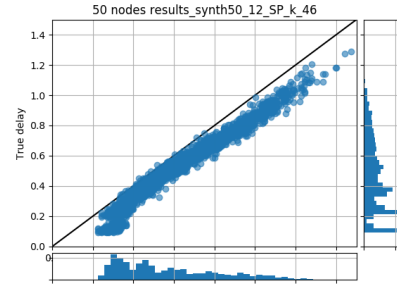


Figure 54: λ 12 routing 46

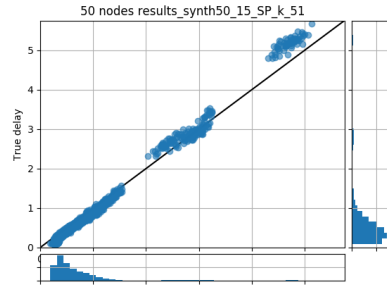


Figure 55: λ 15 routing 51

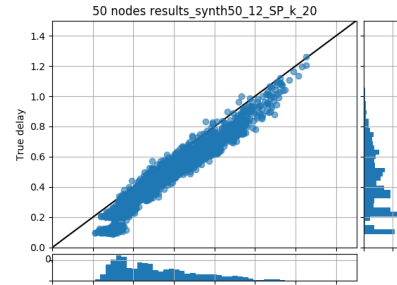


Figure 56: λ 12 routing 20

8.3 True Delay vs Predicted 24 Nodes

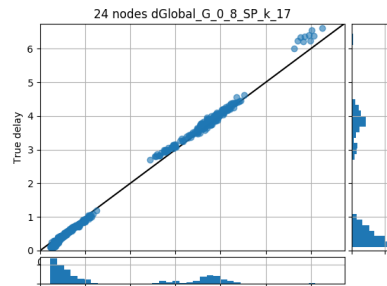


Figure 57: λ 8 routing 17

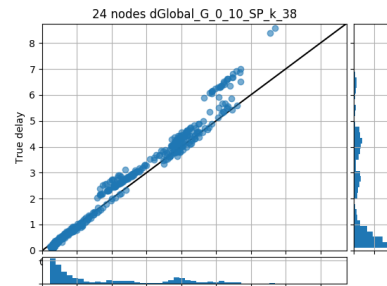


Figure 58: λ 10 routing 38

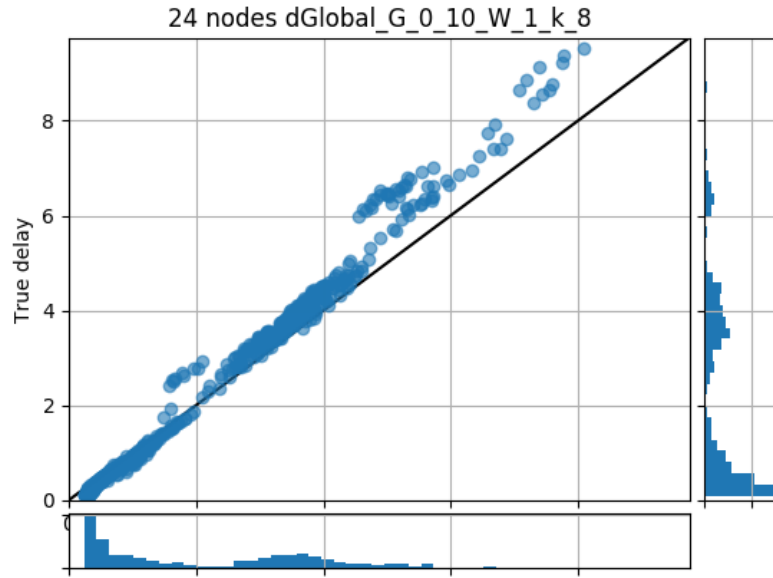


Figure 59: λ 10 routing 1 k 8

8.4 True Delay vs Predicted 14 Nodes

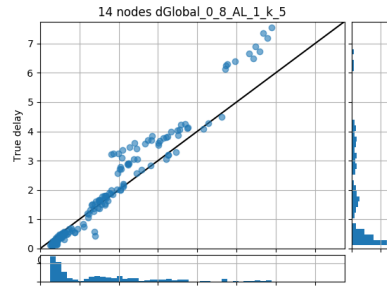


Figure 60: λ 8 routing AL 1 k 5

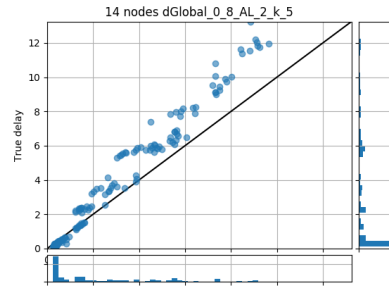


Figure 61: λ 8 routing 2 k 5

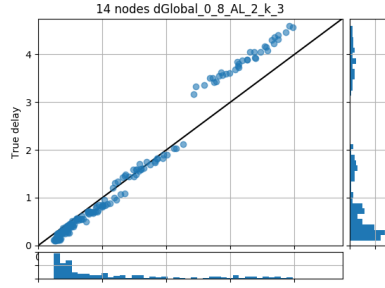


Figure 62: λ 8 routing AL 2 k 3

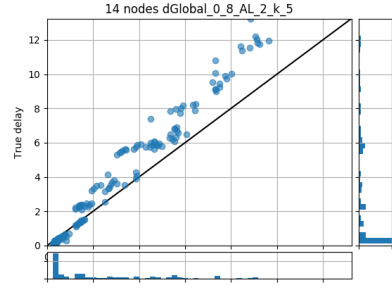


Figure 63: λ 8 routing 2 k 5

8.5 Table of results for 14 Nodes with λ 15

	R^2	ρ	MSE	MRE	MAE
14 Nodes Full	0.23	0.49	1.80	0.161	0.71

Table 10: Results for 14 nodes with λ 15 included in the dataset